

IMAGE TAMPERING DETECTION AND LOCALIZATION USING DEEP LEARNING TECHNIQUES

TEAM ID : FYP2026P1T15

BY:

JOTHI SRI S 2022103049

KEERTHIKA B 2022103552

NALINIKSHA P 2022103553

PUNITHA K 2022103561

PROJECT GUIDE :

Dr. ANGELIN GLADSTON,

Associate Professor, Department of Computer Science and Engineering.

INTRODUCTION

THE PROBLEM OF IMAGE MANIPULATION

- Image manipulation is widely used to create fake images for misinformation, fraud, and fake news.
- Common manipulation types include:
 - Splicing – Copying content from one image into another.
 - Copy-move – Duplicating and moving regions within the same image.
 - Removal – Erasing or inpaint-removing objects from an image.
- These fake images are often visually blended into authentic backgrounds, making them hard to detect manually.

INTRODUCTION

DETECTION CHALLENGES & OUR FOCUS

- Detecting splicing, copy-move, and removal edits is difficult because:

Manipulated areas are fine-tuned to match the background.

Low-level traces (noise, lighting) are often erased.

However, boundaries of tampered regions often contain subtle forensic cues that can be key to detection.

- Our work introduces a boundary-guided deep model that:

Leverages edge-aware attention to focus on manipulation boundaries.

Uses contrastive learning to separate tampered vs. authentic features.

Can also be adapted for image-level fake vs. real classification.

OBJECTIVES

- **Detect & localize** tampered regions at pixel-level.
- **Leverage boundary cues** to improve manipulation detection.
- **Distinguish** between object boundaries and tampered boundaries.
- **Enhance robustness** with image-level classification (authentic vs. manipulated).
- Achieve **state-of-the-art performance** on benchmark datasets.

LITERATURE SURVEY

Paper Title, Authors, Year, Journal name, Vol, Issue & PP	Methodology/ Technique/ Algorithm used	Issues / Gaps/ Limitations	Idea for Adoption
CFL-Net: Image Forgery Localization Using Contrastive Learning Fahim Faisal Niloy, Kishor Kumar Bhaumik, Simon S. Woo, 2023, IEEE Proceedings of WACV 2023, pp. 1518-1527	<ol style="list-style-type: none">Combines contrastive loss with cross-entropy loss for better feature separation.Learns distinct embeddings for tampered vs untampered regions.	<ol style="list-style-type: none">Lacks advanced feature fusion mechanisms (e.g., attention or transformers).Focused mainly on pixel-level separation, less on boundary refinement.	<ol style="list-style-type: none">Apply contrastive learning to enhance tampered–untampered feature discrimination.Fuse RGB and noise-based features for improved localization.
ObjectFormer for Image Manipulation Detection and Localization Junke Wang, Zuxuan Wu, Jingjing Chen, Xintong Han, Abhinav Shrivastava, Ser-Nam Lim, Yu-Gang Jiang, 2023, IEEE Transactions on Information Forensic.	<ol style="list-style-type: none">Combine RGB and high-frequency features for patch embeddings.Use learnable object prototypes for object-level consistency.Stacked encoders/decoders refine patch and object embeddings.	<ol style="list-style-type: none">Depends on frequency feature quality.Robustness against adversarial edits unclear.Generalization to other datasets not tested.	<ol style="list-style-type: none">Automated image tampering detection.Integrate into digital forensics tools.Extend to video tampering detection.

LITERATURE SURVEY

Paper Title, Authors, Year, Journal name, Vol, Issue & PP	Methodology/ Technique/ Algorithm used	Issues / Gaps/ Limitations	Idea for Adoption
Image Manipulation Detection by Multi-View Multi-Scale Supervision Xinru Chen, Chengbo Dong, Jiaqi Ji, Juan Cao, Xirong Li, 2021, IEEE Proceedings of ICCV 2021, pp. 14165-14173	<ol style="list-style-type: none">Dual-branch CNN architecture.Uses multi-view feature learning from noise and edge cues.Employs multi-scale supervision for robust training.	<ol style="list-style-type: none">Faces difficulty in differentiating natural edges from tampered ones.Image-level loss reduces pixel-level localization performance.	<ol style="list-style-type: none">Use edge-based supervision to strengthen boundary detection.Apply multi-scale learning for better feature representation.Combine pixel-level and image-level tasks for balanced accuracy.
ManTra-Net: Manipulation Tracing Network for Detection and Localization of Image Forgeries Yue Wu, Wael AbdAlmageed, Premkumar Natarajan, 2019 ,IEEE Proceedings of CVPR 2019, pp. 9535-9544	<ol style="list-style-type: none">End-to-end FCN for image forgery detection and localization.Self-supervised learning on 385 manipulation types.Uses Z-score anomaly detection and LSTM for local tampering analysis.	<ol style="list-style-type: none">Lacks boundary-awareness for precise edge localizationComputationally heavy due to extensive manipulation-type training.Localization accuracy reduces on fine-grained or blended manipulations.	<ol style="list-style-type: none">Apply self-supervised learning to capture manipulation patterns.Use local anomaly detection for subtle tampering spots.Extend model to handle unknown or mixed manipulations effectively.

LITERATURE SURVEY

Paper Title, Authors, Year, Journal name, Vol, Issue & PP	Methodology/ Technique/ Algorithm used	Issues / Gaps/ Limitations	Idea for Adoption
Bi-Stream Multiscale Hamhead Networks with Contrastive Learning for Image Forgery Localization G. Li et al., 2023, Springer	<div><div>1.</div><div>Bi-stream multiscale network combining RGB and frequency domain features</div></div> <div><div>2.</div><div>Contrastive learning to enhance feature discrimination</div></div>	<div><div>1.</div><div>Performance may vary on unseen datasets</div></div> <div><div>2.</div><div>Computational cost high due to multiscale processing</div></div>	<div><div>1.</div><div>Automated image forgery detection in digital forensics</div></div> <div><div>2.</div><div>Media authenticity verification tools</div></div>
EC-Net: General Image Tampering Localization Network Based on Edge Distribution Guidance and Contrastive Learning Q. Hao et al., 2024, Elsevier	<div><div>1.</div><div>Edge distribution guidance to capture boundary-level tampering</div></div> <div><div>2.</div><div>Contrastive learning to improve feature representation</div></div>	<div><div>1.</div><div>May struggle with subtle or low-contrast manipulations</div></div> <div><div>2.</div><div>Limited evaluation on diverse datasets</div></div>	<div><div>1.</div><div>Accurate image tampering detection in digital forensics</div></div> <div><div>2.</div><div>Can be used in media authenticity verification platforms</div></div>

SUMMARY OF ISSUES

1. Boundary Confusion

- Natural object edges often mistaken as tampered boundaries.
- Leads to high false positives in textured regions.

2. No Image-Level Decision

- Base model outputs only pixel masks.
- Cannot classify whether the image itself is authentic or manipulated.

3. Weak Boundary Discrimination

- Contrastive learning treats all boundaries similarly.
- Tampered vs natural edges are not explicitly separated.

4. Limited Robustness & Generalization

- Struggles with subtle / low-contrast manipulations.
- Performance drops significantly across different datasets.

PROPOSED SYSTEM

BOUNDARY-GUIDED TAMPERING DETECTION WITH DUAL-EDGE LEARNING

A novel deep learning model that localizes spliced, copy-moved, and removed regions by leveraging boundary priors and contrastive feature separation.

Input & Preprocessing:

Multi-source benchmark data → resizing → normalization → pseudo-edge generation

Feature Encoding:

ResNet-50 extracts multi-scale forensic representations

Dual-Branch Boundary Decoder:

Tamper-aware branch (BAM) + natural-edge branch + boundary discriminator

Multi-Task Output:

Pixel-level forgery mask + image-level authenticity classifier

PROPOSED SYSTEM

KEY INNOVATIONS & CONTRIBUTIONS

Problem Identified:

Prior works such as CFL-Net, ObjectFormer, ManTra-Net, and Multi-View Multi-Scale Supervision (MVSS-Net) often confuse tampered boundaries with natural object edges, leading to localization errors and false positives.

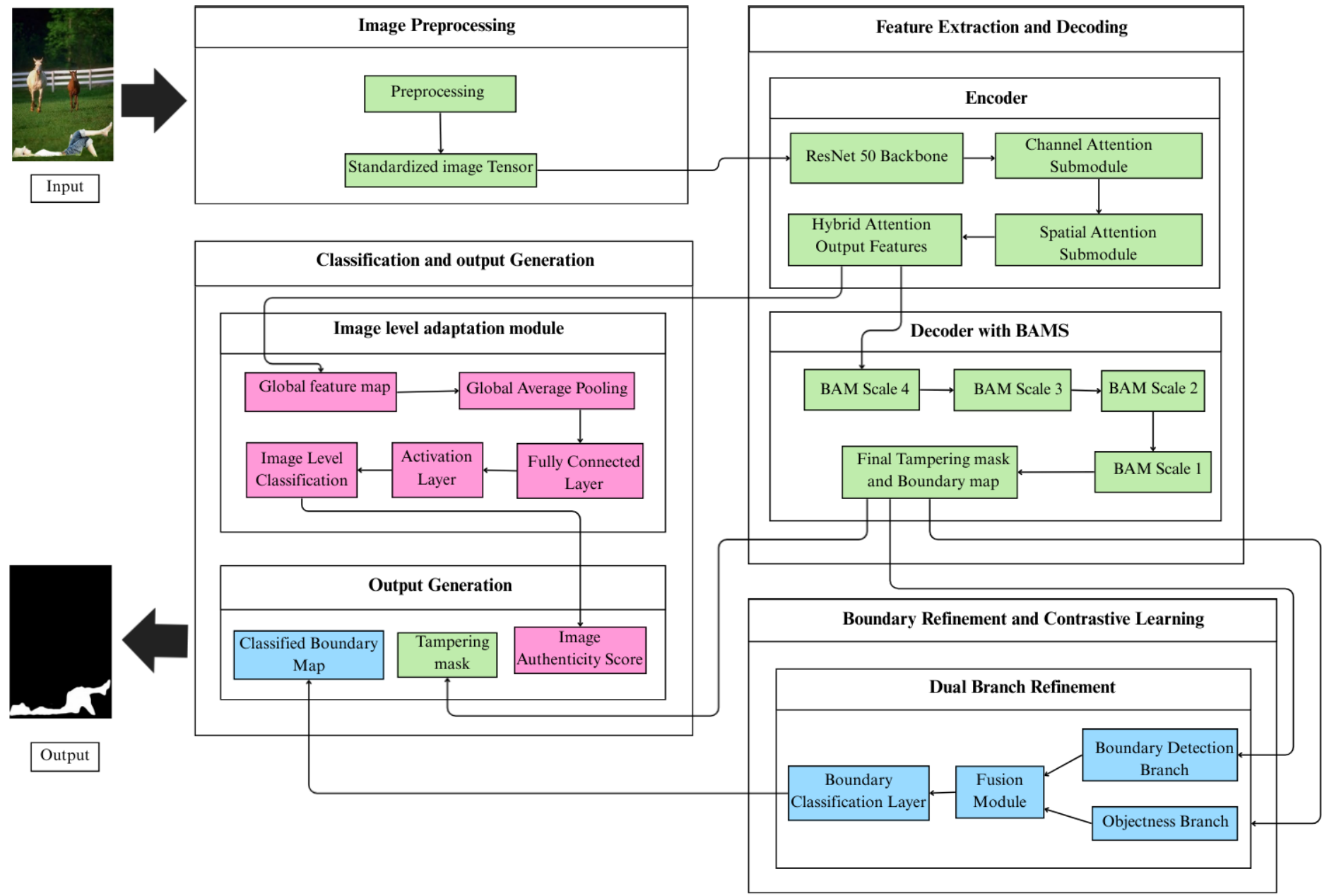
Our Solution:

- Dual-Branch Edge Learning explicitly separates tampered and natural boundaries.
- Boundary-Aware Attention focuses on forensic cues in manipulated regions.
- Boundary-Guided Contrastive Learning enhances feature discriminability.
- Unified Framework supports both pixel- and image-level detection.

Expected Outcome:

- More accurate and robust localization.
- Better generalization across forgery types and datasets.
- Reduced false positives from natural edges and complex backgrounds.

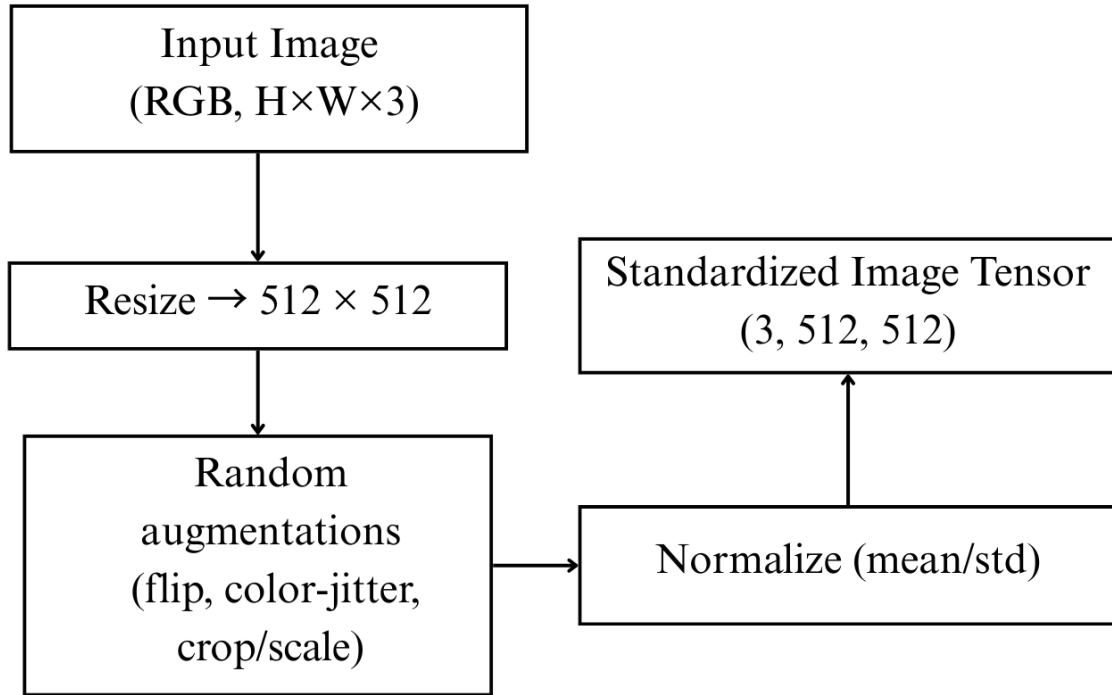
SYSTEM ARCHITECTURE FOR IMAGE TAMPERING DETECTION AND PIXEL-LEVEL LOCALIZATION



DETAILS OF MODULES

1. Preprocessing Module
2. Encoder Module (Feature Extraction using ResNet-50)
3. Decoder with Boundary-Aware Attention Module (BAM)
4. Dual-Branch Boundary Refinement Module
5. Boundary-Guided Contrastive Learning Module
6. Image-Level Adaptation Module

1. PREPROCESSING MODULE



Purpose: Standardize input images and masks for model consistency, augment data for generalization, and generate boundary ground truth for boundary-guided learning.

Functions: Resize images/masks; apply training augmentations (jitter, flip); normalize to ImageNet stats; convert mask to binary tensor; compute boundary via morphological operations.

Input: PIL Image (RGB tampered image), PIL Mask (grayscale ground truth).

Output: Normalized image tensor (3x512x512), binary mask tensor (1x512x512), boundary tensor (1x512x512).

Algorithm:

Algorithm 1: Image Preprocessing

Input: RGB Image I

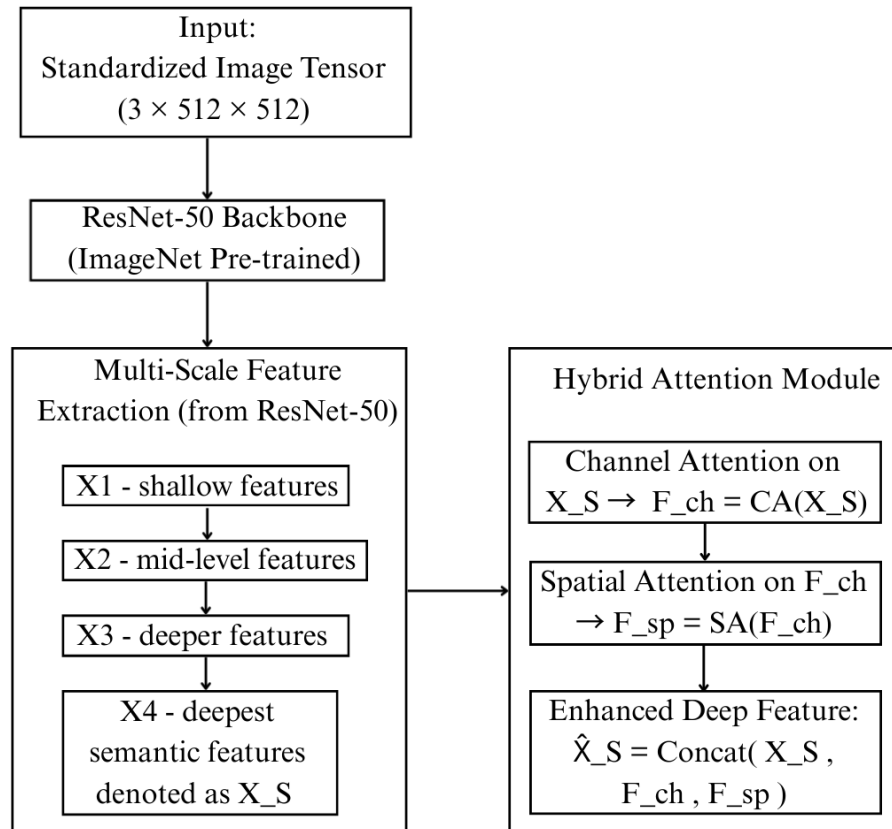
Output: Standardized Image Tensor T

1. Start
2. Read the input RGB image I
3. Resize the image to 512×512
4. Apply random augmentations (flip, color jitter, crop/scale)
5. Normalize the image using mean and standard deviation
6. Convert the image into tensor format $(3, 512, 512)$
7. Output the standardized image tensor T
8. Stop

Process:

1. Input RGB image is collected from the dataset
2. Image is resized to a fixed resolution of 512×512
3. Random augmentations are applied during training
4. Image is normalized using mean and standard deviation
5. Image is converted into a standardized tensor
6. Preprocessed tensor is forwarded to the encode

2. ENCODER MODULE (FEATURE EXTRACTION USING RESNET-50)



Purpose:

Extract multi-scale visual features and enhance them using channel and spatial attention to better capture manipulation-related patterns.

Functions:

Generate hierarchical features via ResNet-50; apply channel attention to highlight important feature channels; apply spatial attention to emphasize key regions; produce enriched deep features for decoding.

Input:

Normalized image tensor $(3 \times 512 \times 512)$.

Output:

Multi-scale encoder feature maps $(X1, X2, X3, X_S)$ and hybrid-attention enhanced deep feature (\hat{X}_S) .

Algorithm:

Algorithm 2: Encoder with Hybrid Attention

Input: Standardized image tensor I

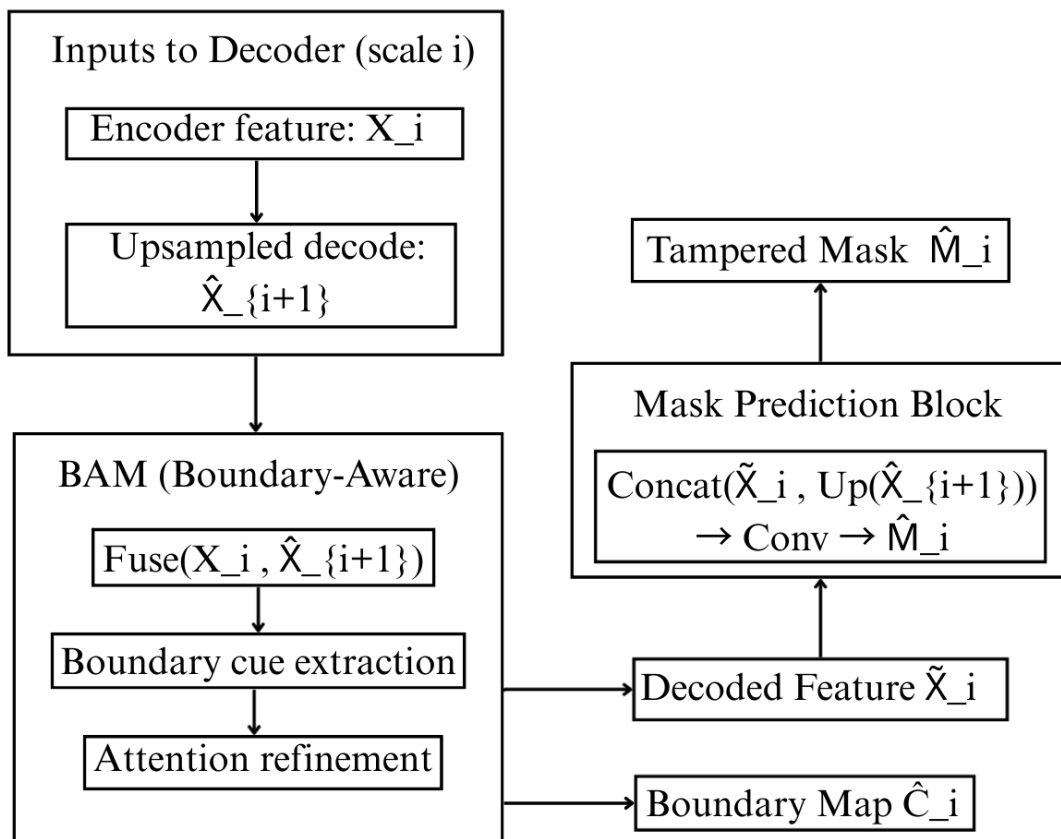
Output: Enhanced deep feature map \hat{X}_S

1. Start
2. Read input image tensor I
3. Pass I through ResNet-50 backbone
4. Extract feature maps X_1, X_2, X_3, X_4
5. Select deepest feature $X_S = X_4$
6. Apply channel attention on X_S
7. Apply spatial attention on channel-refined feature
8. Concatenate attention-refined features to obtain \hat{X}_S
9. Output \hat{X}_S
10. Stop

Process:

1. Preprocessed image tensor is passed to ResNet-50 backbone
2. Multi-scale features are extracted at different depths
3. Deepest semantic feature is selected
4. Channel attention enhances important feature channels
5. Spatial attention highlights informative regions
6. Attention-refined features are fused for decoding

3. DECODER WITH BOUNDARY-AWARE ATTENTION MODULE (BAM)



Purpose:

Enhance the deepest encoder feature map by highlighting important channels and spatial regions, improving manipulation-related feature representation.

Functions:

Compute channel attention via global pooling; compute spatial attention via pooled spatial descriptors; refine the feature map through sequential channel and spatial attention operations.

Input:

Deep encoder feature map

$$X_S \in \mathbb{R}^{B \times C \times H \times W}$$

Output:

Hybrid-attention enhanced feature map

$$\hat{X}_S \in \mathbb{R}^{B \times (3C) \times H \times W}$$

Algorithm:

Algorithm 3: Boundary-Aware Decoder with Mask Prediction

Input: Encoder feature X_i and upsampled decoded feature \hat{X}_{i+1}

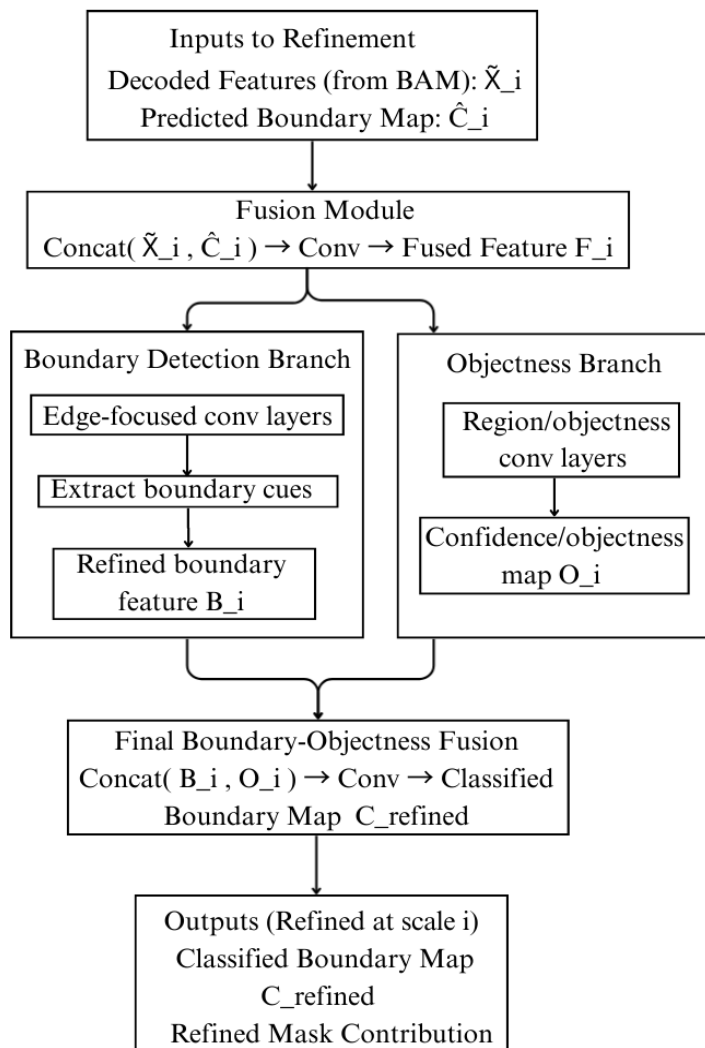
Output: Tampered mask \hat{M}_i and boundary map \hat{C}_i

1. Start
2. Receive encoder feature X_i at decoder scale i
3. Receive upsampled decoded feature \hat{X}_{i+1}
4. Fuse X_i and \hat{X}_{i+1} using Boundary-Aware Module (BAM)
5. Extract boundary cues from fused features
6. Apply attention refinement to obtain decoded feature \hat{X}_i
7. Generate boundary map \hat{C}_i from BAM output
8. Concatenate \hat{X}_i with upsampled \hat{X}_{i+1}
9. Apply convolution to predict tampered mask \hat{M}_i
10. Output \hat{X}_i , \hat{M}_i , and \hat{C}_i
11. Stop

Process:

1. Encoder features and upsampled decoder features are received
2. Features are fused using the Boundary-Aware Module (BAM)
3. Boundary cues are extracted and refined using attention
4. Decoded feature is generated at the current scale
5. Tampered mask is predicted using convolution layers
6. Boundary map is produced for further refinement

4. DUAL-BRANCH BOUNDARY REFINEMENT MODULE



Purpose:

Refine coarse boundary predictions by integrating edge-focused boundary cues and region-level objectness information, producing a cleaner and more reliable boundary representation.

Functions:

Combine aggregated multi-scale boundary maps with decoder-derived objectness features; extract edge-sensitive responses through a boundary stream; extract regional confidence through an objectness stream; and fuse both refined representations to generate an accurate boundary map.

Input:

Refined decoder feature map and multi-scale boundary predictions

$$C_i, X_i^{dec}$$

Output:

Final refined boundary map

$$C_{\text{refined}} \in \mathbb{R}^{B \times 1 \times H \times W}$$

Algorithm:

Algorithm 4: Boundary–Objectness Refinement Module

Input: Decoded feature \hat{X}_i and predicted boundary map \hat{C}_i

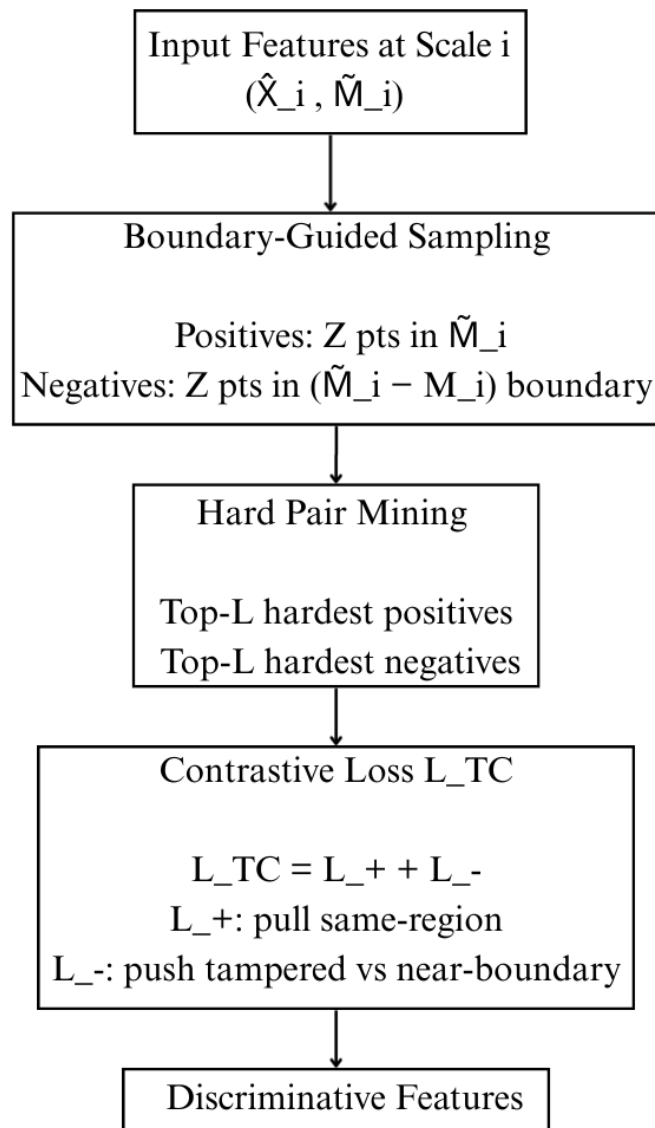
Output: Refined boundary map \tilde{C}_i and refined mask contribution

1. Start
 2. Receive decoded feature \hat{X}_i from BAM
 3. Receive predicted boundary map \hat{C}_i
 4. Concatenate \hat{X}_i and \hat{C}_i and apply convolution to obtain fused feature F_i
 5. Feed F_i into boundary detection branch
 6. Apply edge-focused convolution layers to extract boundary cues
 7. Generate refined boundary feature B_i
 8. Feed F_i into objectness branch
 9. Apply region-focused convolution layers to generate objectness map O_i
 10. Concatenate B_i and O_i and apply convolution for boundary–objectness fusion
 11. Obtain classified boundary map \tilde{C}_i
 12. Output refined boundary map \tilde{C}_i and refined mask contribution
 13. Stop
-

Process:

1. Decoded features and predicted boundary map are fused
2. Boundary detection branch extracts fine edge details
3. Objectness branch estimates region confidence
4. Boundary and objectness features are combined
5. Refined boundary map is generated
6. Refined output contributes to accurate mask prediction

5. BOUNDARY-GUIDED CONTRASTIVE LEARNING MODULE



Purpose:

Improve feature discriminability around tampered boundaries by enforcing separation between tampered and near-boundary non-tampered pixels, leading to clearer and more reliable mask predictions.

Functions:

Construct a boundary-guided negative sampling region, extract feature samples from tampered and boundary-adjacent regions, identify hard positive and hard negative feature pairs, and apply a contrastive objective that strengthens feature separation across decoder scales.

Input:

Decoder feature map: $F \in \mathbb{R}^{B \times C \times H \times W}$

Ground-truth mask: $G \in \mathbb{R}^{B \times 1 \times H \times W}$

Output:

Contrastive loss value: $L_{TC} \in \mathbb{R}$

used to guide boundary-sensitive feature learning during training.

Algorithm:

Algorithm 5: Boundary-Guided Contrastive Learning

Input: Decoded feature \hat{X}_i and predicted mask \hat{M}_i

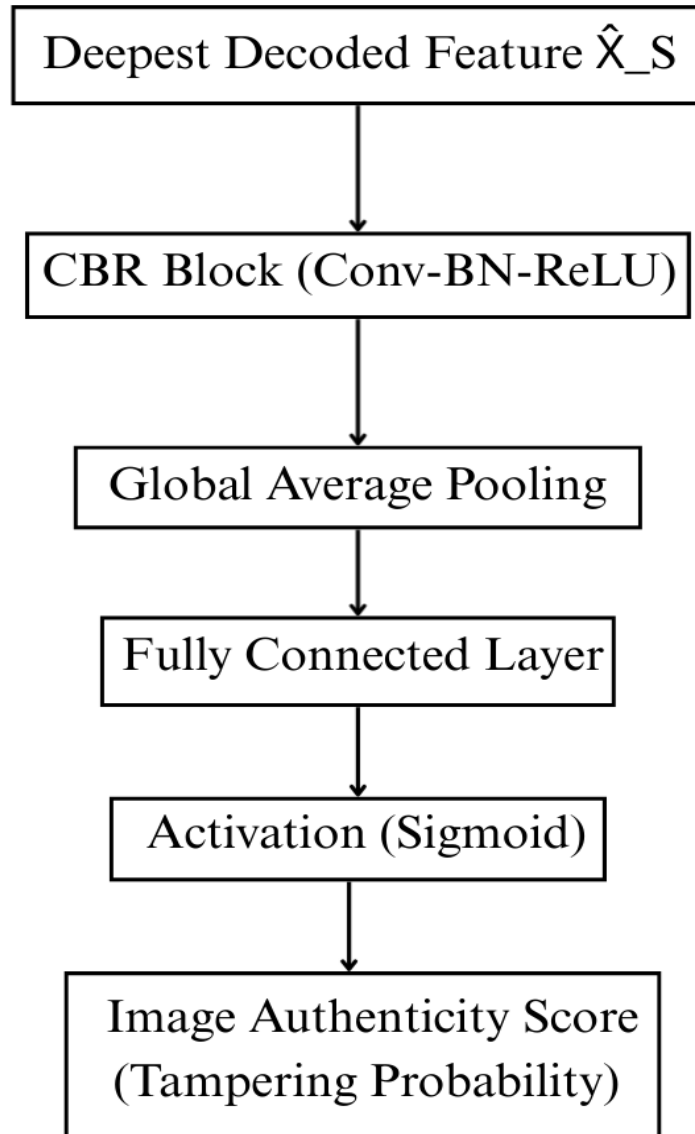
Output: Discriminative feature representations

1. Start
2. Receive decoded feature map \hat{X}_i at scale i
3. Receive predicted tampered mask \hat{M}_i
4. Perform boundary-guided sampling on \hat{M}_i
5. Sample Z positive points from tampered region \hat{M}_i
6. Sample Z negative points from near-boundary region $(\tilde{M}_i - \hat{M}_i)$
7. Apply hard pair mining on sampled points
8. Select top- L hardest positive samples
9. Select top- L hardest negative samples
10. Compute contrastive loss $L_{TC} = L_+ + L_-$
11. Minimize L_+ to pull same-region features closer
12. Minimize L_- to push tampered and near-boundary features apart
13. Obtain discriminative feature representations
14. Stop

Process:

1. Decoded features and predicted mask are taken as input
2. Positive samples are selected from tampered regions
3. Negative samples are selected from near-boundary regions
4. Hard positive and negative samples are mined
5. Contrastive loss pulls similar regions together
6. Contrastive loss pushes boundary-confusing regions apart
7. Discriminative features are learned

6. IMAGE-LEVEL ADAPTATION MODULE



Purpose:

Provide a global tampering prediction to complement pixel-level detection, allowing the network to learn image-level manipulation cues and strengthen overall feature representation..

Functions:

Aggregate deep decoder features into a global descriptor using average pooling, transform them through a lightweight classification head, and produce a single authenticity score indicating whether the entire image is manipulated.

Input:

Deep decoder feature map

$$X_{deep} \in \mathbb{R}^{B \times C \times H \times W}$$

Output:

Image-level tampering probability

$$p_{img} \in [0, 1]$$

.

Algorithm:

Algorithm 6: Image-Level Authenticity Classification

Input: Deepest decoded feature \hat{X}_S

Output: Image authenticity score (tampering probability)

1. Start
2. Receive deepest decoded feature \hat{X}_S
3. Apply CBR block (Convolution–BatchNorm–ReLU) on \hat{X}_S
4. Perform global average pooling on the refined feature
5. Pass pooled feature through a fully connected layer
6. Apply sigmoid activation function
7. Compute image authenticity score representing tampering probability
8. Output image authenticity score
9. Stop

Process:

1. Deepest decoded feature is selected
2. Feature is refined using Conv-BN-ReLU block
3. Global average pooling aggregates spatial information
4. Fully connected layer maps features to a score
5. Sigmoid activation produces tampering probability
6. Final image authenticity score is obtained


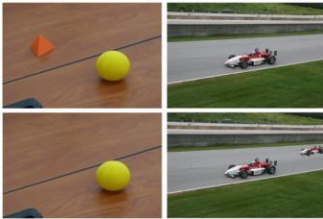

IMPLEMENTATION DETAILS

(20 Percentage)



DATASET

Purpose:

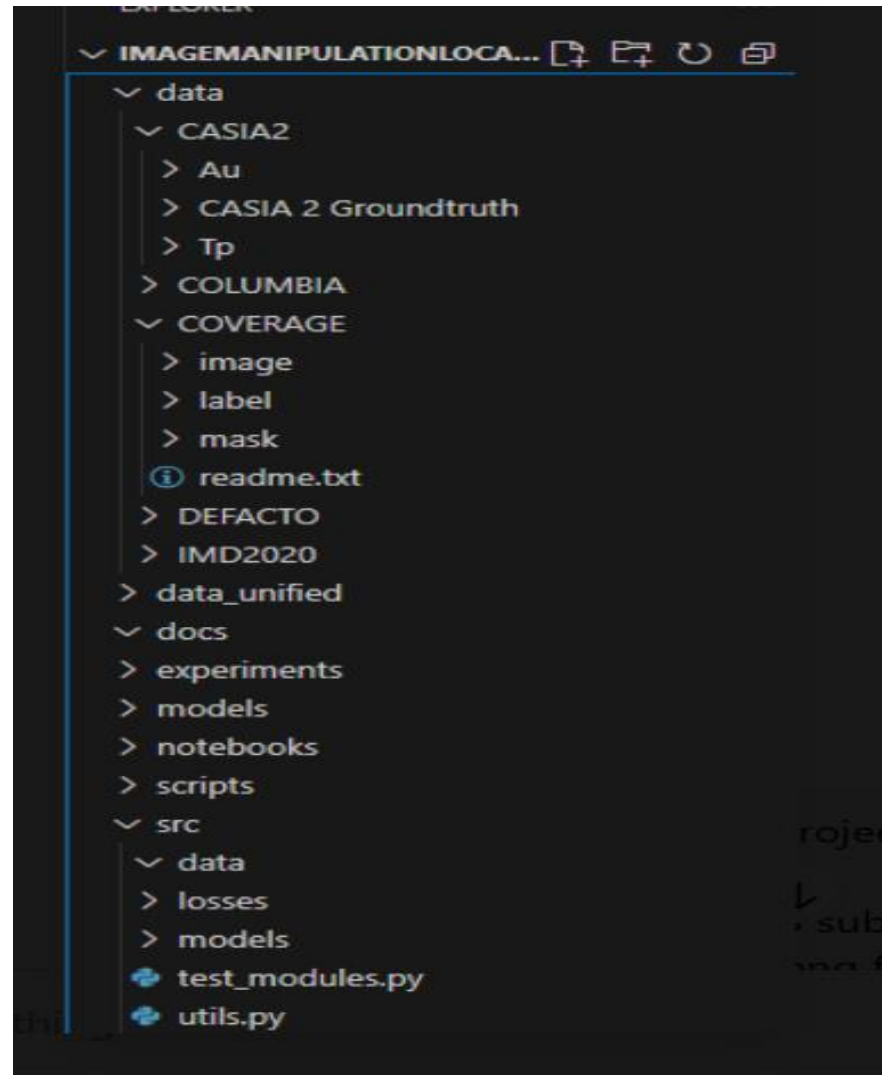
To train and evaluate the model on diverse types of image manipulations such as **splicing**, **copy-move**, and **inpainting** using benchmark datasets.

Dataset Name	# of Images	Manipulation Type	Sample Image	Purpose
CASIA v2.0	~12,000	Splicing, Copy-Move		Training & Evaluation
Columbia	363	Splicing		Cross-Dataset Testing
Coverage	100	Copy-Move		Boundary Evaluation

DATASET

Dataset Name	# of Images	Manipulation Type	Sample Image	Purpose
IMD2020	2,010	Real-World Manipulations		Cross-Dataset Testing
DeFacto	~149,000	Synthetic (Copy-Move, Splicing, Inpainting)		Extended Evaluation

DATASET & FOLDER STRUCTURE



1. CASIA2 Dataset Loading & Transformation

Loads CASIA2 images, applies necessary transformations, and prepares complete ground-truth inputs required for tampering localization.

```
3 import os
4 from PIL import Image
5 from torch.utils.data import Dataset
6 from .transforms import Preprocessor
7
8 class CASIA2Dataset(Dataset):
9     def __init__(self, root, input_size=512, train=True):
10
11         self.tp_dir = os.path.join(root, "Tp")
12         self.gt_dir = os.path.join(root, "CASIA 2 Groundtruth")
13         self.gt_map = {}
14         for fname in os.listdir(self.gt_dir):
15             if fname.lower().endswith("_gt.png"):
16                 base = fname[:-7] # remove "_gt.png"
17                 mask_path = os.path.join(self.gt_dir, fname)
18                 self.gt_map[base] = mask_path
19         self.files = []
20
21         for img_name in os.listdir(self.tp_dir):
22             if img_name.lower().endswith((".jpg", ".jpeg", ".png", ".tif")):
23                 base = os.path.splitext(img_name)[0] # remove extension
24                 if base in self.gt_map:
25                     img_path = os.path.join(self.tp_dir, img_name)
26                     mask_path = self.gt_map[base]
27                     self.files.append((img_path, mask_path))
28
29         print(f"[CASIA2] Loaded {len(self.files)} paired tampered samples.\n")
30         self.preproc = Preprocessor(input_size=input_size, train=train)
31
32     def __len__(self):
33         return len(self.files)
34     def __getitem__(self, idx):
35         img_path, mask_path = self.files[idx]
36
37         img = Image.open(img_path).convert("RGB")
38         mask = Image.open(mask_path).convert("L")
39
40         img_t, mask_t, boundary_t = self.preproc(img, mask)
41
42         return {
43             "image": img_t,
44             "mask": mask_t,
45             "boundary": boundary_t,
46             "img_path": img_path,
47             "mask_path": mask_path
48         }
```

```
11 def pil_to_tensor(img_pil):
12     """Convert PIL RGB image to Torch tensor (C,H,W), float32 in [0,1]."""
13     arr = np.asarray(img_pil).astype(np.float32) / 255.0 # H,W,C
14     if arr.ndim == 2:
15         arr = np.expand_dims(arr, axis=-1)
16     arr = arr[..., :3]
17     arr = arr.transpose(2,0,1)
18     return torch.from_numpy(arr).float()
19
20 def normalize_tensor(tensor):
21     """Normalize ImageNet-style."""
22     mean = torch.tensor(IMAGE_MEAN).view(-1,1,1)
23     std = torch.tensor(IMAGE_STD).view(-1,1,1)
24     return (tensor - mean) / std
25
26 class Preprocessor:
27     def __init__(self, input_size=512, train=True):
28         self.input_size = input_size
29         self.train = train
30
31     def _resize(self, img, mask):
32         img = img.resize((self.input_size, self.input_size), Image.BILINEAR)
33         mask = mask.resize((self.input_size, self.input_size), Image.NEAREST)
34         return img, mask
35
36     def _color_jitter(self, img):
37         if torch.rand(1) < 0.7:
38             img = ImageEnhance.Brightness(img).enhance(1 + (torch.rand(1).item() - 0.5) * 0.4)
39         if torch.rand(1) < 0.7:
40             img = ImageEnhance.Contrast(img).enhance(1 + (torch.rand(1).item() - 0.5) * 0.4)
41         return img
42
43     def _random_flip(self, img, mask):
44         if torch.rand(1) < 0.5:
45             img = img.transpose(Image.FLIP_LEFT_RIGHT)
46             mask = mask.transpose(Image.FLIP_LEFT_RIGHT)
47         return img, mask
48
49     def _mask_to_boundary(self, mask_tensor, k=5):
50         """Dilation - Erosion to produce boundary GT."""
51         x = mask_tensor.unsqueeze(0) # (1,1,H,W)
52         pad = k // 2
53
54         dil = F.max_pool2d(x, kernel_size=k, stride=1, padding=pad)
55         inv = 1 - x
56         ero_inv = F.max_pool2d(inv, kernel_size=k, stride=1, padding=pad)
57         ero = 1 - ero_inv
```

[tampering] Loaded 1502 paired tampered samples

--- SAMPLE INFORMATION ---

Image path: ../data/CASIA2\Tp\Tp_D_CNN_M_N_ani00052_ani00054_11130.jpg

Mask path: ../data/CASIA2\CASIA 2 Groundtruth\Tp_D_CNN_M_N_ani00052_ani00054_11130_gt.png

Image shape: torch.Size([3, 512, 512])

Mask shape: torch.Size([1, 512, 512])

Boundary shape: torch.Size([1, 512, 512])

(tampering) C:\Users\UG\ImageManipulationLocalization\src>.

Preprocessed Image



Mask GT



Boundary GT



2. ResNet-50 Feature Extraction Backbone

Extracts multi-scale deep features from input images, providing the essential encoded representations required for accurate tampering localization.

```
17 class ResNet50Backbone(nn.Module):
18     def __init__(self, pretrained=True):
19         super().__init__()
20         # Load torchvision resnet50
21         resnet = models.resnet50(pretrained=pretrained)
22
23         # Keep the initial layers (conv1, bn1, relu, maxpool)
24         self.stem = nn.Sequential(
25             resnet.conv1, # stride=2
26             resnet.bn1,
27             resnet.relu,
28             resnet.maxpool # reduces to 1/4
29         )
30
31         # ResNet blocks
32         self.layer1 = resnet.layer1 # C2
33         self.layer2 = resnet.layer2 # C3
34         self.layer3 = resnet.layer3 # C4
35         self.layer4 = resnet.layer4 # C5
36
37         self._out_channels = {
38             "c2": 256, # layer1 output channels for resnet50
39             "c3": 512, # layer2
40             "c4": 1024, # layer3
41             "c5": 2048 # layer4
42         }
43
44     def forward(self, x):
45
46         # stem -> reduces to 1/4 resolution
47         x = self.stem(x)
48         c2 = self.layer1(x) # 1/4
49         c3 = self.layer2(c2) # 1/8
50         c4 = self.layer3(c3) # 1/16
51         c5 = self.layer4(c4) # 1/32
52
53         return {"c2": c2, "c3": c3, "c4": c4, "c5": c5}
54
55     def out_channels(self):
56         return self._out_channels
57
```

```
[INFO] Loaded sample for backbone test
Image shape: torch.Size([3, 512, 512])
Using device: cuda
C:\Users\UG\anaconda3\envs\tampering\lib\site-packages\torchvision\models_utils.py:208: UserWarning: The parameter 'pretrained' is deprecated since 0.13 and may be removed in the future, please use 'weights'
instead.
  warnings.warn(
C:\Users\UG\anaconda3\envs\tampering\lib\site-packages\torchvision\models_utils.py:223: UserWarning: Arguments other than a weight enum or 'None' for 'weights' are deprecated since 0.13 and may be removed i
n the future. The current behavior is equivalent to passing `weights=ResNet50_Weights.IMAGENET1K_V1`. You can also use `weights=ResNet50_Weights.DEFAULT` to get the most up-to-date weights.
  warnings.warn(msg)

----- BACKBONE FEATURE MAP SHAPES -----
c2: (1, 256, 128, 128)
c3: (1, 512, 64, 64)
c4: (1, 1024, 32, 32)
c5: (1, 2048, 16, 16)

Expected sizes for 512x512 input:
c2 -> 128x128
c3 -> 64x64
c4 -> 32x32
c5 -> 16x16

===== BACKBONE MODULE TEST COMPLETE =====

(tampering) C:\Users\UG\ImageManipulationLocalization\src\test>
```

3. Hybrid Attention Module (HAM)

The Hybrid Attention Module refines the deepest encoder feature (C5) by applying channel and spatial attention to emphasize manipulation-relevant regions. It produces an enhanced feature map (c5_hat) that provides stronger, more informative input to the decoder for accurate tampering localization.

```
19
20 class ChannelAttention(nn.Module):
21     def __init__(self, channels, reduction=16):
22         super().__init__()
23         self.gap = nn.AdaptiveAvgPool2d(1)
24         mid = max(1, channels // reduction)
25         self.fc = nn.Sequential(
26             nn.Conv2d(channels, mid, kernel_size=1, bias=True),
27             nn.ReLU(inplace=True),
28             nn.Conv2d(mid, channels, kernel_size=1, bias=True),
29             nn.Sigmoid()
30         )
31
32 class NonLocalSpatialAttention(nn.Module):
33     def __init__(self, in_channels, inter_channels=None):
34         super().__init__()
35         if inter_channels is None:
36             inter_channels = max(1, in_channels // 2)
37
38         self.g = nn.Conv2d(in_channels, inter_channels, kernel_size=1, bias=False)
39         self.theta = nn.Conv2d(in_channels, inter_channels, kernel_size=1, bias=False)
40         self.phi = nn.Conv2d(in_channels, inter_channels, kernel_size=1, bias=False)
41
42         self.W = nn.Conv2d(inter_channels, in_channels, kernel_size=1, bias=False)
43         # initialize W to zeros for residual learning stability
44         nn.init.constant_(self.W.weight, 0.0)
45
46 class HybridAttention(nn.Module):
47     def __init__(self, in_channels, reduction=16):
48         super().__init__()
49         self.cam = ChannelAttention(in_channels, reduction=reduction)
50         self.sam = NonLocalSpatialAttention(in_channels)
51
52         # small conv to smooth the fusion (optional but stabilizes)
53         self.fusion_conv = nn.Sequential(
54             nn.Conv2d(in_channels, in_channels, kernel_size=1, bias=False),
55             nn.BatchNorm2d(in_channels),
56             nn.ReLU(inplace=True)
57         )
```

```
DEBUG: Found 5122 groundtruth masks in GT folder.
[CASIA2] Loaded 4981 paired tampered samples.

Using device: cuda
C:\Users\UG\anaconda3\envs\tampering\lib\site-packages\torchvision\models\_utils.py:208: UserWarning: The parameter 'pretrained'
' instead.
  warnings.warn(
C:\Users\UG\anaconda3\envs\tampering\lib\site-packages\torchvision\models\_utils.py:223: UserWarning: Arguments other than a we
n the future. The current behavior is equivalent to passing `weights=ResNet50_Weights.IMAGENET1K_V1`. You can also use `weights
  warnings.warn(msg)
c5 shape: torch.Size([1, 2048, 16, 16])
c5_hat shape: torch.Size([1, 2048, 16, 16])
HAM test complete.

(tampering) C:\Users\UG\ImageManipulationLocalization\src\test>.
```


4. Boundary-Aware Decoder (BAM + Multi-Scale Refinement)

Integrates BAM at each scale to enhance feature refinement using boundary-focused attention. Combines deep and shallow features to progressively rebuild a precise tampering mask. Outputs both mask predictions and boundary cues for improved localization.

```
23
24 class BAM(nn.Module):
25     def __init__(self, enc_channels, dec_channels, topk=32, embed_dim=256):
26         super().__init__()
27         self.enc_channels = enc_channels
28         self.dec_channels = dec_channels
29         self.topk = topk
30         self.embed_dim = embed_dim
31         # for computing high-frequency mask: small conv on pooled features
32         self.hf_pool_conv = nn.Sequential(
33             nn.AdaptiveAvgPool2d(1),
34             nn.Conv2d(enc_channels + dec_channels, enc_channels + dec_channels, kernel_size=1),
35             nn.ReLU(inplace=True),
36             nn.Conv2d(enc_channels + dec_channels, enc_channels + dec_channels, kernel_size=1),
37             nn.Sigmoid()
38         )
39         # boundary predictor
40         self.boundary_pred = nn.Sequential(
41             nn.Conv2d(enc_channels + dec_channels, enc_channels // 4, kernel_size=3, padding=1),
42             nn.ReLU(inplace=True),
43             nn.Conv2d(enc_channels // 4, 1, kernel_size=1),
44             nn.Sigmoid()
45         )
46         # embeddings for attention (project to lower embed_dim)
47         self.enc_proj = nn.Conv2d(enc_channels, embed_dim, kernel_size=1, bias=False)
48         self.dec_proj = nn.Conv2d(dec_channels, embed_dim, kernel_size=1, bias=False)
49         # project back to enc channels
50         self.out_proj = nn.Conv2d(embed_dim, enc_channels, kernel_size=1, bias=False)
51         # Fusion conv after concat(updated_enc, upsampled_decoder)
52         self.fuse_conv = nn.Sequential(
53             nn.Conv2d(enc_channels + dec_channels, enc_channels, kernel_size=3, padding=1, bias=False),
54             nn.BatchNorm2d(enc_channels),
55             nn.ReLU(inplace=True)
56         )
```

```
Using device: cuda
C:\Users\UG\anaconda3\envs\tampering\lib\site-packages\torchvision\models\_utils.py:208: UserWarning:
' instead.
  warnings.warn(
C:\Users\UG\anaconda3\envs\tampering\lib\site-packages\torchvision\models\_utils.py:223: UserWarning:
n the future. The current behavior is equivalent to passing `weights=ResNet50_Weights.IMAGENET1
  warnings.warn(msg)
Shapes: c2 (1, 256, 128, 128) c3 (1, 512, 64, 64) c4 (1, 1024, 32, 32) c5 (1, 2048, 16, 16)
c5_hat shape: (1, 2048, 16, 16)
BAM outputs shapes:
X_hat i: (1, 1024, 32, 32)
boundary_up: (1, 1, 32, 32)
BAM test complete.

(tampering) C:\Users\UG\ImageManipulationLocalization\src\test>
```

```

8
9 import torch
10 import torch.nn as nn
11 import torch.nn.functional as F
12
13 from .bam import BAM
14
15 class SimpleDecoder(nn.Module):
16     def __init__(self, channels_c2=256, channels_c3=512, channels_c4=1024, channels_c5=2048,
17                 topk=32, embed_dim=256):
18         super().__init__()
19         self.bam_c4 = BAM(enc_channels=channels_c4, dec_channels=channels_c5, topk=topk, embed_dim=embed_dim)
20         self.bam_c3 = BAM(enc_channels=channels_c3, dec_channels=channels_c4, topk=topk, embed_dim=embed_dim)
21         self.bam_c2 = BAM(enc_channels=channels_c2, dec_channels=channels_c3, topk=topk, embed_dim=embed_dim)
22         self.smooth4 = nn.Sequential(
23             nn.Conv2d(channels_c4, channels_c4, kernel_size=3, padding=1, bias=False),
24             nn.BatchNorm2d(channels_c4),
25             nn.ReLU(inplace=True)
26         )
27         self.smooth3 = nn.Sequential(
28             nn.Conv2d(channels_c3, channels_c3, kernel_size=3, padding=1, bias=False),
29             nn.BatchNorm2d(channels_c3),
30             nn.ReLU(inplace=True)
31         )
32         self.smooth2 = nn.Sequential(
33             nn.Conv2d(channels_c2, channels_c2, kernel_size=3, padding=1, bias=False),
34             nn.BatchNorm2d(channels_c2),
35             nn.ReLU(inplace=True)
36         )
37
38         # Final mask head: project channels -> 1 and upsample to input size later
39         self.mask_head = nn.Sequential(
40             nn.Conv2d(channels_c2, channels_c2//2, kernel_size=3, padding=1, bias=False),
41             nn.BatchNorm2d(channels_c2//2),
42             nn.ReLU(inplace=True),
43             nn.Conv2d(channels_c2//2, 1, kernel_size=1)
44         )
45

```

Using device: cuda

C:\Users\UG\anaconda3\envs\tampering\lib\site-packages\torchvision\models_utils.py:208: UserWarning: The parameter 'p' instead.

warnings.warn(

C:\Users\UG\anaconda3\envs\tampering\lib\site-packages\torchvision\models_utils.py:223: UserWarning: Arguments other than 'weights' in the future. The current behavior is equivalent to passing 'weights=ResNet50_Weights.IMAGENET1K_V1'. You can also use

warnings.warn(msg)

Backbone shapes: c2 (1, 256, 128, 128) c3 (1, 512, 64, 64) c4 (1, 1024, 32, 32) c5 (1, 2048, 16, 16)

c5_hat: (1, 2048, 16, 16)

Decoder outputs:

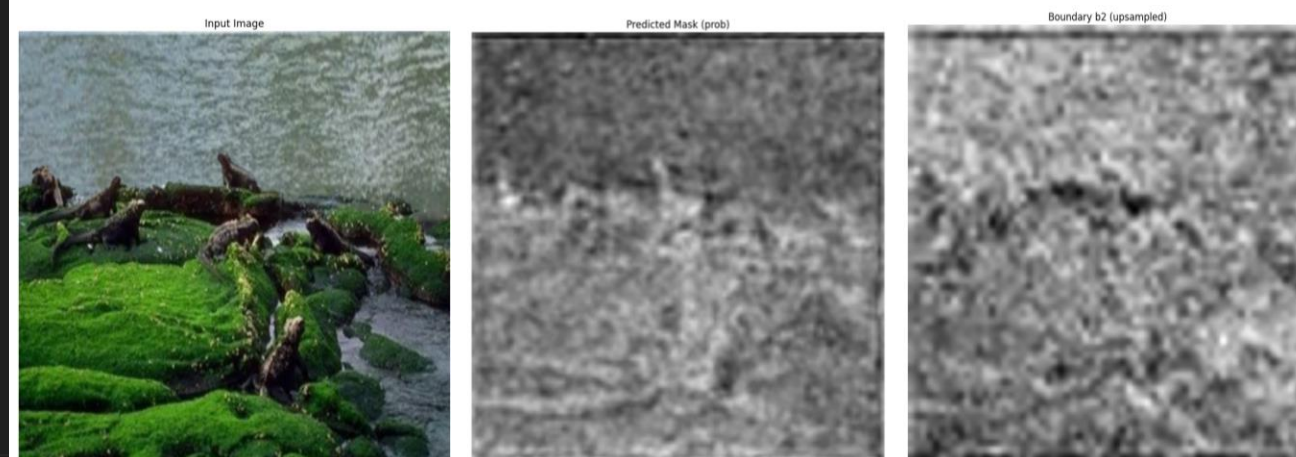
mask_logits_up: (1, 1, 512, 512)

boundary b4: (1, 1, 32, 32)

boundary b3: (1, 1, 64, 64)

boundary b2: (1, 1, 128, 128)

final_feat: (1, 256, 128, 128)



PERFORMANCE METRICS

1.	Intersection over Union (IoU)	Measures the overlap between the predicted tampered region and the ground truth. Higher IoU means better localization accuracy.	$\text{IoU} = \text{TP} / (\text{TP} + \text{FP} + \text{FN})$
2.	Boundary F1-Score (BF1)	Evaluates how accurately the model predicts the boundaries of tampered regions. Useful because your project includes Boundary Refinement.	$\text{BF1} = (2 * \text{Precision_b} * \text{Recall_b}) / (\text{Precision_b} + \text{Recall_b})$
3.	Mean Absolute Error (MAE)	Measures the average pixel-wise difference between the predicted probability map and ground-truth mask.	$\text{MAE} = (1/N) * \sum P_i - G_i $
4.	Precision–Recall AUC (PR-AUC)	Area under the Precision–Recall curve. More meaningful than ROC-AUC for imbalanced tampering masks (few tampered pixels).	$\text{PR-AUC} = \int \text{Precision}(\text{Recall}) d(\text{Recall})$

PERFORMANCE METRICS

5.	Matthews Correlation Coefficient (MCC)	A balanced performance metric that works even when classes are imbalanced.	$\text{MCC} = (\text{TP} * \text{TN} - \text{FP} * \text{FN}) / \sqrt{(\text{TP} + \text{FP}) * (\text{TP} + \text{FN}) * (\text{TN} + \text{FP}) * (\text{TN} + \text{FN})}$
6.	Boundary IoU (BIOU)	Measures IoU only along the boundary region (dilated ground-truth boundary). Good when your model uses boundary-aware refinement.	$\text{BIOU} = (\text{Boundary}(\text{P}) \cap \text{Boundary}(\text{G})) / (\text{Boundary}(\text{P}) \cup \text{Boundary}(\text{G}))$
7.	F2-Score (Recall-Focused)	A variant of F1 that gives higher weight to recall (important in forensics — missing tampered pixels is worse).	$\text{F2} = (5 * \text{Precision} * \text{Recall}) / (4 * \text{Precision} + \text{Recall})$
8.	Pixel Accuracy (Acc)	Measures overall correctness at pixel level.	$\text{Accuracy} = (\text{TP} + \text{TN}) / (\text{TP} + \text{TN} + \text{FP} + \text{FN})$

DELIVERABLES TO BE DEMONSTRATED

1. **Proposed Image Manipulation Localization Model** implemented using ResNet-50 backbone
2. **Dual-Branch Boundary Refinement Module** distinguishing object vs tampered boundaries
3. **Image-Level Adaptation Module** for authentic/manipulated classification
4. **Training and Evaluation** on benchmark datasets (CASIA v2.0, NIST16, Columbia, Coverage)
5. **Visualization of Results:**
 1. Tampering masks and boundary maps
 2. Image-level classification output
6. **Performance Metrics:** F1-Score, AUC, IoU, Boundary F1, Accuracy
7. **Comparative Analysis** with base paper performance

CONCLUSION

- Proposed a Boundary-Guided Tampering Detection Model with Dual-Branch Learning to distinguish tampered from natural edges.
- Enhanced localization using Boundary-Aware Attention (BAM) and Multi-Scale Contrastive Learning.
- Integrated an Image-Level Adaptation Module for global authenticity classification.
- Achieved state-of-the-art performance across benchmark datasets in both pixel-level and image-level detection.

REFERENCES

- [1] Wenxi Liu, Hao Zhang, Xinyang Lin, Qing Zhang, Qi Li, Xiaoxiang Liu, and Ying Cao, “Attentive and Contrastive Image Manipulation Localization With Boundary Guidance”, in IEEE Transactions on Information Forensics and Security, vol. 19, pp. 6764–6778, 2024.
- [2] Lusen Dong; Sichen Li; Jin Zheng, “Degraded Image Semantic Segmentation Using Intra-image and Inter-image Contrastive Learning”, in IEEE 2023 China Automation Congress (CAC), March 2024.
- [3] Yuyuan Zeng; Bowen Zhao; Shanzhao Qiu; Tao Dai; Shu-Tao Xia, “Toward Effective Image Manipulation Detection With Proposal Contrastive Learning”, in IEEE Transactions on Circuits and Systems for Video Technology, vol 33, no:9, pp 4703 - 4714, February 2023.
- [4] Muhammad Zubair Khan; Mohan Kumar Gajendran; Yugyung Lee; Muazzam A. Khan, “Deep Neural Architectures for Medical Image Semantic Segmentation: Review”, in IEEE Access, vol 9, June 2021.
- [5] Shervin Minaee; Yuri Boykov; Fatih Porikli; Antonio Plaza; Nasser Kehtarnavaz; Demetri Terzopoulos, “Image Segmentation Using Deep Learning: A Survey”, in IEEE Transactions on Pattern Analysis and Machine Intelligence, vol 44, no:7, pp 3523 - 3542, February 2021.