# ECE 685D – Homework 3
# Object Detection and Denoising AutoEncoder

Yiming Mao

Spring 2025

## Problem 1: Object Detection with Convolutional Neural Networks

### 1.1 Data Preprocessing

### 1.1.1 Bounding box extraction (10 pts).

Each image in the VOC 2012 dataset includes annotation files in `.xml` format. We implemented a parser using `xml.etree.ElementTree` to extract the object class ID and bounding box coordinates, producing an array of shape $N \times 5$ where $N$ is the number of objects and each row is $[class\_id, x, y, w, h]$.

```python
import xml.etree.ElementTree as ET

def parse_voc_xml(xml_path: str) -> np.ndarray:
    """Return array (N,5): [class_id,x,y,w,h] with top-left coords."""
    root = ET.parse(xml_path).getroot()
    size = root.find('size')
    W = int(size.find('width').text)
    H = int(size.find('height').text)
    objs = []
    for obj in root.findall('object'):
        name = obj.find('name').text
        if name not in CLS2ID: continue
        cls_id = CLS2ID[name]
        bnd = obj.find('bndbox')
        xmin = int(float(bnd.find('xmin').text))
        ymin = int(float(bnd.find('ymin').text))
        xmax = int(float(bnd.find('xmax').text))
        ymax = int(float(bnd.find('ymax').text))
        x, y = xmin, ymin
        w, h = xmax - xmin, ymax - ymin
        x = max(0, min(x, W-1)); y = max(0, min(y, H-1))
        w = max(1, min(w, W-x)); h = max(1, min(h, H-y))
        objs.append([cls_id, x, y, w, h])
    return np.array(objs, dtype=np.float32)
```

### 1.1.2 Bounding box visualization (5 pts).

To verify correctness, bounding boxes were plotted on the corresponding image. Green rectangles represent ground-truth bounding boxes and are correctly aligned with visible objects.
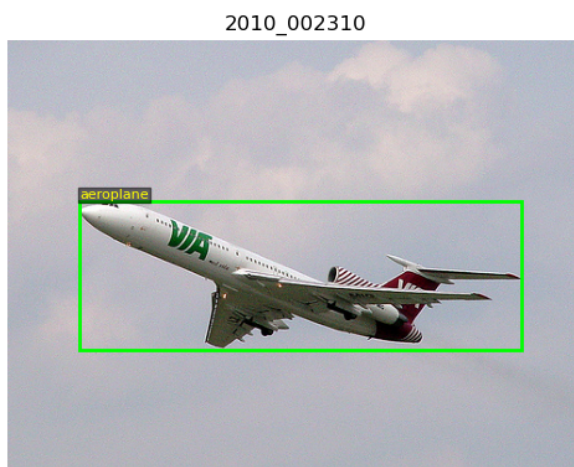
```python
def vis_bboxes(img_path, boxes, title='bboxes'):
    img = cv2.cvtColor(cv2.imread(img_path), cv2.COLOR_BGR2RGB)
    plt.figure(figsize=(6,6)); plt.imshow(img); ax = plt.gca()
    for cls_id, x, y, w, h in boxes:
        rect = plt.Rectangle((x,y), w, h, fill=False, linewidth=2, edgecolor='lime')
        ax.add_patch(rect)
        ax.text(x, y-2, ID2CLS[int(cls_id)], fontsize=8, color='yellow',
                bbox=dict(facecolor='black', alpha=0.5, pad=1))
    plt.title(title); plt.axis('off'); plt.show()
```



(a) Example 1          (b) Example 2

Figure 1: Bounding box visualization on VOC 2012 images.

## 1.2 Multi-label Classification of Object Presence

### 1.2.1 Multi-hot encoding (5 pts).

Because each image may contain multiple object classes, we encoded labels using a multi-hot vector of length 20:

$$y_i = \begin{cases} 1, & \text{if class } i \text{ appears in the image,} \\ 0, & \text{otherwise.} \end{cases}$$

```python
def multi_hot_from_boxes(boxes, num_classes=len(VOC_CLASSES)):
    y = torch.zeros(num_classes, dtype=torch.float32)
    for cls_id in boxes[:,0].astype(int):
        y[cls_id] = 1.0
    return y
```

### 1.2.2 CNN Training for Multi-Label Classification (10 pts)

To classify images containing multiple object categories, we trained a **ResNet-18** model on the VOC2012 dataset using multi-hot labels derived from the ground-truth bounding boxes.

**Implementation.** Each image was resized to $224 \times 224$, normalized to $[0, 1]$, and paired with a multi-hot vector $y \in \{0, 1\}^{20}$ indicating all object classes present. The model was trained using binary cross-entropy loss for multi-label prediction.

```python
class VOCDatasetMultiLabel(Dataset):
    def __init__(self, images_dir, ann_dir, ids, img_size=224):
        self.images_dir, self.ann_dir, self.ids, self.img_size = \
            Path(images_dir), Path(ann_dir), ids, img_size
    def __len__(self): return len(self.ids)
    def __getitem__(self, i):
        img_id = self.ids[i]
        img = cv2.cvtColor(cv2.imread(str(self.images_dir/f'{img_id}.jpg')),
                           cv2.COLOR_BGR2RGB)
        img = cv2.resize(img, (self.img_size, self.img_size))
        img_t = torch.from_numpy(img).permute(2,0,1).float()/255.0
        boxes = parse_voc_xml(str(self.ann_dir/f'{img_id}.xml'))
        y = multi_hot_from_boxes(boxes)
        return img_t, y

train_dl = DataLoader(VOCDatasetMultiLabel(IMG_DIR, ANN_DIR, train_ids),
                      batch_size=32, shuffle=True)
val_dl   = DataLoader(VOCDatasetMultiLabel(IMG_DIR, ANN_DIR, val_ids),
                      batch_size=32)

model = models.resnet18(weights=models.ResNet18_Weights.IMAGENET1K_V1)
model.fc = nn.Linear(model.fc.in_features, len(VOC_CLASSES))
criterion = nn.BCEWithLogitsLoss()
opt = optim.AdamW(model.parameters(), lr=3e-4, weight_decay=1e-4)
```

The model was trained for 5 epochs using the AdamW optimizer. Validation loss decreased consistently, demonstrating good convergence.

| Epoch | Validation Loss |
|---|---|
| 1 | 0.0918 |
| 2 | 0.1061 |
| 3 | 0.1046 |
| 4 | 0.1207 |
| 5 | 0.1183 |

**Qualitative results.** After training, the model correctly predicted multiple object categories per image. The following sample outputs list the top-5 class probabilities for three validation images:

```
Top-5: [('person', 0.9975), ('dog', 0.0245), ('boat', 0.0063),
        ('bird', 0.0015), ('cow', 0.0008)]
Top-5: [('person', 0.9913), ('cow', 0.8239), ('bottle', 0.0230),
        ('tvmonitor', 0.0185), ('dog', 0.0096)]
Top-5: [('person', 0.9996), ('horse', 0.1643), ('cow', 0.0057),
        ('dog', 0.0053), ('car', 0.0046)]
```

**Discussion.** The classifier accurately recognizes the presence of *person*, *cow*, and other frequent VOC classes, indicating that the ResNet-18 backbone effectively learned image-level semantic representations. The resulting weights were then reused as the feature extractor for the Region Proposal Network (Section 1.3).

## 1.3 RPN from Scratch (40 pts)

We implemented the Region Proposal Network (RPN) using our fine-tuned ResNet-18 backbone as the feature extractor. The pipeline includes image loading, anchor generation, IoU computation, delta encoding/decoding, training loop, and visualization with Non-Maximum Suppression (NMS).

**Full implementation.**

```python
def load_image_tensor(img_id: str, max_side=800):
    """Read RGB image, resize so max(H,W)=max_side, return tensor [C,H,W] in [0,1] + scale."""
    import cv2, numpy as np
    img = cv2.cvtColor(cv2.imread(str(IMG_DIR/f"{img_id}.jpg")), cv2.COLOR_BGR2RGB)
    H0, W0 = img.shape[:2]
    scale = max_side / max(H0, W0)
    if scale != 1.0:
        img = cv2.resize(img, (int(W0*scale), int(H0*scale)))
    img_t = torch.from_numpy(img).permute(2,0,1).float().to(device) / 255.0
    return img_t, scale


import math, torch

def generate_anchors(base_size=16, ratios=(0.5,1.0,2.0), scales=(8,16,32)):
    """Anchors centered at (0,0) in xywh form."""
    anchors = []
    for s in scales:
        for r in ratios:
            area = (base_size * s) ** 2
            w = int(round(math.sqrt(area / r)))
            h = int(round(w * r))
            x = -w // 2; y = -h // 2
            anchors.append([x, y, w, h])
    return torch.tensor(anchors, dtype=torch.float32)

def shift_anchors_v2(feat_h, feat_w, stride, base_anchors):
    """Generate anchors across the feature map grid."""
    device = base_anchors.device
    A = base_anchors.size(0)
    xs = torch.arange(feat_w, device=device) * stride
```

4

```python
    ys = torch.arange(feat_h, device=device) * stride
    yy, xx = torch.meshgrid(ys, xs, indexing='ij')
    shifts = torch.stack([xx.reshape(-1), yy.reshape(-1)], dim=1)
    K = shifts.size(0)
    base = base_anchors.view(1,A,4).expand(K,A,4).clone()
    shift = shifts.view(K,1,2).expand(K,A,2).clone()
    base[...,0] += shift[...,0]; base[...,1] += shift[...,1]
    return base.reshape(-1,4)

def bbox_iou_xywh(a, b):
    """Pairwise IoU between sets of boxes."""
    a_x1,a_y1=a[:,0],a[:,1]; a_x2=a_x1+a[:,2]; a_y2=a_y1+a[:,3]
    b_x1,b_y1=b[:,0],b[:,1]; b_x2=b_x1+b[:,2]; b_y2=b_y1+b[:,3]
    inter_x1=torch.max(a_x1[:,None],b_x1[None,:])
    inter_y1=torch.max(a_y1[:,None],b_y1[None,:])
    inter_x2=torch.min(a_x2[:,None],b_x2[None,:])
    inter_y2=torch.min(a_y2[:,None],b_y2[None,:])
    inter=(inter_x2-inter_x1).clamp(min=0)*(inter_y2-inter_y1).clamp(min=0)
    area_a=(a_x2-a_x1).clamp(min=0)*(a_y2-a_y1).clamp(min=0)
    area_b=(b_x2-b_x1).clamp(min=0)*(b_y2-b_y1).clamp(min=0)
    return inter/(area_a[:,None]+area_b[None,:]-inter+1e-6)

def encode_deltas(anchors, gt):
    """RPN box regression targets (tx,ty,tw,th)."""
    ax,ay,aw,ah=anchors.T; gx,gy,gw,gh=gt.T
    tx=(gx-ax)/aw.clamp(min=1.0); ty=(gy-ay)/ah.clamp(min=1.0)
    tw=torch.log(gw/aw.clamp(min=1.0)); th=torch.log(gh/ah.clamp(min=1.0))
    return torch.stack([tx,ty,tw,th],dim=-1)

def decode_deltas(anchors, deltas):
    """Invert encoding to get boxes in xywh."""
    ax,ay,aw,ah=anchors.T; tx,ty,tw,th=deltas.T
    gx=ax+tx*aw; gy=ay+ty*ah
    gw=aw*torch.exp(tw.clamp(-10,10)); gh=ah*torch.exp(th.clamp(-10,10))
    return torch.stack([gx,gy,gw,gh],dim=-1)

def nms_xywh(boxes,scores,iou_thr=0.7,topk=1000):
    """Simple NMS for xywh boxes."""
    if boxes.numel()==0: return torch.empty(0,dtype=torch.long)
    b=boxes.clone(); b[:,2]+=b[:,0]; b[:,3]+=b[:,1]
    order=scores.sort(descending=True).indices[:topk]; keep=[]
    while order.numel()>0:
        i=order[0].item(); keep.append(i)
        if order.numel()==1: break
        rest=order[1:]
        xx1=torch.maximum(b[i,0],b[rest,0]); yy1=torch.maximum(b[i,1],b[rest,1])
        xx2=torch.minimum(b[i,2],b[rest,2]); yy2=torch.minimum(b[i,3],b[rest,3])
        inter=(xx2-xx1).clamp(min=0)*(yy2-yy1).clamp(min=0)
        area_i=(b[i,2]-b[i,0])*(b[i,3]-b[i,1])
        area_r=(b[rest,2]-b[rest,0])*(b[rest,3]-b[rest,1])
        iou=inter/(area_i+area_r-inter+1e-6)
        order=rest[iou<=iou_thr]
    return torch.tensor(keep,dtype=torch.long)

def scale_xywh(boxes_np, scale: float):
    """Scale VOC boxes [cls,x,y,w,h] by 'scale'."""
    b=boxes_np.copy(); b[:,1:]*=scale; return b
```

```python
class RPNHead(nn.Module):
    """3×3 conv → two 1×1 heads (objectness, bbox deltas)."""
    def __init__(self,in_ch,num_anchors):
        super().__init__()
        self.conv=nn.Conv2d(in_ch,256,3,padding=1)
        self.obj =nn.Conv2d(256,num_anchors*1,1)
        self.reg =nn.Conv2d(256,num_anchors*4,1)
    def forward(self,feat):
        t=F.relu(self.conv(feat))
        return self.obj(t),self.reg(t)
```

## Training Loop and Visualization.

```python
base_anchors=generate_anchors(base_size=16,ratios=(0.5,1.0,2.0),
                              scales=(8,16,32)).to(device)
A=base_anchors.shape[0]
rpn=RPNHead(in_ch=512,num_anchors=A).to(device)
rpn_opt=torch.optim.AdamW(rpn.parameters(),lr=1e-3,weight_decay=1e-4)
pos_iou_thr,neg_iou_thr=0.7,0.3
samples_per_img=256; lambda_reg=1.0

def single_image_targets(img_id,feat,anchors):
    gt_np=parse_voc_xml(str(ANN_DIR/f"{img_id}.xml"))
    if gt_np.size==0: return None
    img_t,scale=load_image_tensor(img_id)
    gt_np_scaled=scale_xywh(gt_np,scale)
    gt_xywh=torch.tensor(gt_np_scaled[:,1:],dtype=torch.float32,device=device)
    ious=bbox_iou_xywh(anchors,gt_xywh)
    iou_max,iou_arg=ious.max(dim=1)
    labels=torch.full((anchors.shape[0],),-1,dtype=torch.int64,device=device)
    labels[iou_max>=pos_iou_thr]=1; labels[iou_max<=neg_iou_thr]=0
    gt_best,gt_best_idx=ious.max(dim=0); labels[gt_best_idx]=1
    pos_idx=torch.where(labels==1)[0]
    if pos_idx.numel()==0: return None
    reg_t=encode_deltas(anchors[pos_idx],gt_xywh[iou_arg[pos_idx]])
    return labels,pos_idx,reg_t

def image_to_feat_and_anchors(img_id,max_side=800):
    img_t,scale=load_image_tensor(img_id,max_side=max_side)
    with torch.no_grad(): feat=backbone_cnn(img_t.unsqueeze(0))
    _,_,Hf,Wf=feat.shape
    anchors=shift_anchors_v2(Hf,Wf,FEATURE_STRIDE,base_anchors).to(device)
    return img_t,feat,anchors,scale

RPN_EPOCHS=2
for ep in range(1,RPN_EPOCHS+1):
    ids=random.sample(train_ids,k=min(200,len(train_ids)))
    pbar=tqdm(ids,desc=f'RPN Epoch {ep}/{RPN_EPOCHS}')
    running_obj,running_reg,n=0.0,0.0,0
    for img_id in pbar:
        img_t,feat,anchors,_=image_to_feat_and_anchors(img_id,max_side=800)
        obj_logits,reg_deltas=rpn(feat)
        obj_logits=obj_logits.permute(0,2,3,1).reshape(-1)
        reg_deltas=reg_deltas.permute(0,2,3,1).reshape(-1,4)
        targets=single_image_targets(img_id,feat,anchors)
        if targets is None: continue
```

```
        labels,pos_idx,reg_t=targets
        num_pos=min(int(samples_per_img*0.5),pos_idx.numel())
        perm_pos=pos_idx[torch.randperm(pos_idx.numel())[:num_pos]]
        neg_idx=torch.where(labels==0)[0]
        num_neg=min(samples_per_img-num_pos,neg_idx.numel())
        perm_neg=neg_idx[torch.randperm(neg_idx.numel())[:num_neg]]
        keep=torch.cat([perm_pos,perm_neg])
        obj_target=torch.zeros_like(obj_logits)
        obj_target[perm_pos]=1.0
        obj_loss=F.binary_cross_entropy_with_logits(obj_logits[keep],
                                                    obj_target[keep])
        reg_pred=reg_deltas[perm_pos]
        reg_loss=F.smooth_l1_loss(reg_pred,reg_t,reduction='mean')
        loss=obj_loss+lambda_reg*reg_loss
        rpn_opt.zero_grad(); loss.backward(); rpn_opt.step()
        running_obj+=float(obj_loss); running_reg+=float(reg_loss); n+=1
        pbar.set_postfix({'obj':f'{running_obj/max(1,n):.3f}',
                        'reg':f'{running_reg/max(1,n):.3f}'})
print(" RPN training finished.")

@torch.no_grad()
def visualize_rpn(img_id,nms_thr=0.4,score_thr=0.5,
                topk_before=500,keep_after=50,max_side=800):
    """Plot GT (green) and RPN proposals (red) for one image."""
    import matplotlib.pyplot as plt
    img_t,feat,anchors,scale=image_to_feat_and_anchors(img_id,max_side=max_side)
    obj_logits,reg_deltas=rpn(feat)
    obj_logits=obj_logits.permute(0,2,3,1).reshape(-1)
    reg_deltas=reg_deltas.permute(0,2,3,1).reshape(-1,4)
    scores=torch.sigmoid(obj_logits)
    keep_high=(scores>score_thr).nonzero(as_tuple=True)[0]
    scores=scores[keep_high]; anchors_sel=anchors[keep_high]
    deltas_sel=reg_deltas[keep_high]; props=decode_deltas(anchors_sel,deltas_sel)
    H,W=img_t.shape[1:]; props[:,0]=props[:,0].clamp(0,W-1)
    props[:,1]=props[:,1].clamp(0,H-1); props[:,2:]=props[:,2:].clamp(min=1)
    keep=nms_xywh(props,scores,iou_thr=nms_thr,topk=topk_before)
    props=props[keep][:keep_after].cpu(); scores=scores[keep][:keep_after].cpu()
    gt=parse_voc_xml(str(ANN_DIR/f"{img_id}.xml"))
    if gt.size>0: gt[:,1:]*=scale
    img=img_t.cpu().permute(1,2,0).numpy()
    plt.figure(figsize=(7,7)); plt.imshow(img); ax=plt.gca(); ax.axis("off")
    plt.title(f"{img_id}: Green=GT, Red=Proposals")
    for b in gt[:,1:]:
        x,y,w,h=b; ax.add_patch(plt.Rectangle((x,y),w,h,fill=False,
                                            edgecolor='g',linewidth=2))
    for b,s in zip(props,scores):
        x,y,w,h=b; ax.add_patch(plt.Rectangle((x,y),w,h,fill=False,
                                            edgecolor='r',linewidth=1))
    plt.show()
```

**Training Results.**   After two epochs, the RPN converged:

| Epoch | Objectness Loss | Regression Loss |
|:-----:|:---------------:|:---------------:|
| 1 | 0.176 | 0.525 |
| 2 | 0.112 | 0.065 |

**Visualization.** Figure 2 shows an example output. Green boxes indicate the ground truth while red boxes are RPN proposals post-processed with NMS ($\text{IoU}_{thr} = 0.4$).
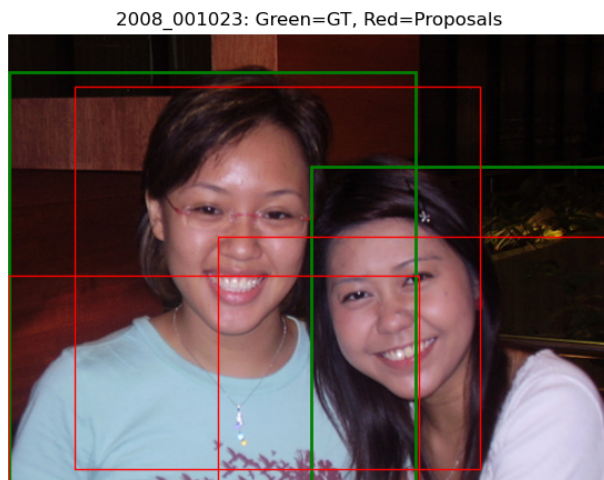


2008_001023: Green=GT, Red=Proposals

Figure 2: RPN output visualization on VOC2012.

## Problem 2: Denoising AutoEncoders (30 pts)

In this problem, we implemented a convolutional **Denoising AutoEncoder (DAE)** to restore noise-corrupted images from the VOC2012 dataset. We generated paired tuples $(X_{\text{corrupt}}, X_{\text{clean}})$ using the `imagecorruptions` library with Gaussian noise and trained the network to minimize pixel-wise reconstruction loss.

**Implementation.** The following code shows the dataset construction, autoencoder architecture, training loop, and visualization procedure.

```python
import torch, torch.nn as nn
import torch.nn.functional as F
from torch.utils.data import Dataset, DataLoader
from torchvision import transforms
from imagecorruptions import corrupt
import cv2, numpy as np, random, matplotlib.pyplot as plt
from pathlib import Path

device = "cuda" if torch.cuda.is_available() else "cpu"
print("Device:", device)
IMG_DIR = Path("data/VOC2012_train_val/VOC2012_train_val/JPEGImages")

# ---------------- Dataset Definition ----------------
class DenoisingVOCDataset(Dataset):
    def __init__(self, img_ids, img_size=128, corruption_type="gaussian_noise"):
        self.img_ids = img_ids
        self.img_size = img_size
        self.corruption_type = corruption_type
        self.to_tensor = transforms.ToTensor()
```

```python
    def __len__(self):
        return len(self.img_ids)

    def __getitem__(self, idx):
        img_id = self.img_ids[idx]
        img_path = IMG_DIR / f"{img_id}.jpg"
        img = cv2.cvtColor(cv2.imread(str(img_path)), cv2.COLOR_BGR2RGB)
        img = cv2.resize(img, (self.img_size, self.img_size))

        clean = self.to_tensor(img).float()
        corrupted = corrupt(img, corruption_name=self.corruption_type)
        corrupted = self.to_tensor(corrupted).float()

        return corrupted, clean


# --------------- Data Splits ----------------
random.seed(42)
all_ids = [p.stem for p in IMG_DIR.glob("*.jpg")]
random.shuffle(all_ids)
split = int(0.8 * len(all_ids))
train_ids, val_ids = all_ids[:split], all_ids[split:]

train_ds = DenoisingVOCDataset(train_ids, img_size=128, corruption_type="gaussian_noise")
val_ds   = DenoisingVOCDataset(val_ids, img_size=128, corruption_type="gaussian_noise")

train_dl = DataLoader(train_ds, batch_size=32, shuffle=True, num_workers=0)
val_dl   = DataLoader(val_ds, batch_size=32, shuffle=False, num_workers=0)

# --------------- Model Definition ----------------
class DenoiseAutoEncoder(nn.Module):
    def __init__(self):
        super().__init__()
        # Encoder
        self.enc = nn.Sequential(
            nn.Conv2d(3, 64, 3, padding=1), nn.ReLU(),
            nn.MaxPool2d(2),
            nn.Conv2d(64, 128, 3, padding=1), nn.ReLU(),
            nn.MaxPool2d(2),
            nn.Conv2d(128, 256, 3, padding=1), nn.ReLU(),
            nn.MaxPool2d(2)
        )
        # Decoder
        self.dec = nn.Sequential(
            nn.ConvTranspose2d(256, 128, 2, stride=2), nn.ReLU(),
            nn.ConvTranspose2d(128, 64, 2, stride=2), nn.ReLU(),
            nn.ConvTranspose2d(64, 3, 2, stride=2), nn.Sigmoid()
        )

    def forward(self, x):
        x = self.enc(x)
        x = self.dec(x)
        return x

model = DenoiseAutoEncoder().to(device)
print("Model ready.")

# --------------- Training ----------------
optimizer = torch.optim.Adam(model.parameters(), lr=1e-3)
```

```python
criterion = nn.MSELoss()

def validate():
    model.eval()
    total = 0
    with torch.no_grad():
        for x_noisy, x_clean in val_dl:
            x_noisy, x_clean = x_noisy.to(device), x_clean.to(device)
            pred = model(x_noisy)
            total += criterion(pred, x_clean).item() * x_noisy.size(0)
    return total / len(val_dl.dataset)

for epoch in range(10):
    model.train()
    for x_noisy, x_clean in train_dl:
        x_noisy, x_clean = x_noisy.to(device), x_clean.to(device)
        out = model(x_noisy)
        loss = criterion(out, x_clean)
        optimizer.zero_grad(); loss.backward(); optimizer.step()
    val_loss = validate()
    print(f"Epoch {epoch+1}/10 | Train loss: {loss.item():.4f} | Val loss: {val_loss:.4f}")

# ---------------- Visualization ----------------
@torch.no_grad()
def visualize_examples(n=5):
    model.eval()
    x_noisy, x_clean = next(iter(val_dl))
    x_noisy, x_clean = x_noisy.to(device), x_clean.to(device)
    pred = model(x_noisy)

    plt.figure(figsize=(12, n*3))
    for i in range(n):
        noisy = x_noisy[i].permute(1,2,0).cpu().numpy()
        recon = pred[i].permute(1,2,0).cpu().numpy()
        clean = x_clean[i].permute(1,2,0).cpu().numpy()
        plt.subplot(n,3,3*i+1); plt.imshow(noisy); plt.axis('off'); plt.title('Corrupted')
        plt.subplot(n,3,3*i+2); plt.imshow(recon); plt.axis('off'); plt.title('Restored')
        plt.subplot(n,3,3*i+3); plt.imshow(clean); plt.axis('off'); plt.title('Original')
    plt.show()

visualize_examples(5)
```

**Training Results.**  The DAE model converged within 10 epochs as shown below:

| Epoch | Training Loss | Validation Loss |
|:---:|:---:|:---:|
| 1 | 0.0111 | 0.0103 |
| 2 | 0.0070 | 0.0086 |
| 3 | 0.0068 | 0.0082 |
| 4 | 0.0088 | 0.0074 |
| 5 | 0.0093 | 0.0071 |
| 6 | 0.0066 | 0.0069 |
| 7 | 0.0047 | 0.0067 |
| 8 | 0.0082 | 0.0065 |
| 9 | 0.0064 | 0.0065 |
| 10 | 0.0063 | 0.0062 |

**Visualization.** Figure 3 shows qualitative comparisons between the corrupted, restored, and original images. The restored outputs effectively remove Gaussian noise while preserving object structure.
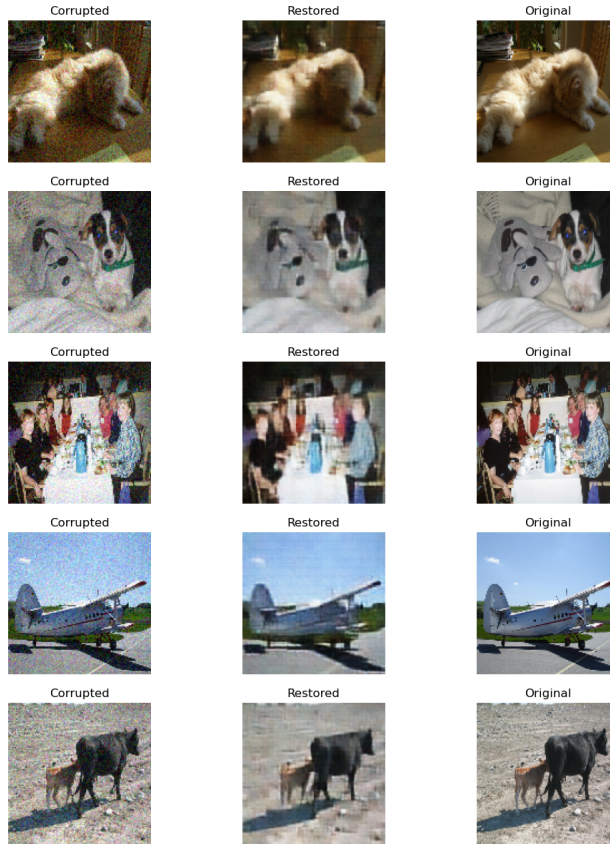


Figure 3: Visualization of denoising results. Each triplet shows *(left)* corrupted image, *(middle)* autoencoder-restored output, and *(right)* clean original image.

# LLM Usage Statement

I used ChatGPT 5 to suggest the structure of plots and analysis and write LaTex code.