

✓ Homework 3

Instructions

- This homework focuses on understanding and applying DETR for object detection and attention visualization. It consists of **three questions** designed to assess both theoretical understanding and practical application.
- Please organize your answers and results for the questions below and submit this jupyter notebook as a **.pdf file**.
- **Deadline: 11/14 (Thur) 23:59**

Reference

- End-to-End Object Detection with Transformers (DETR): <https://github.com/facebookresearch/detr>

✓ Q1. Understanding DETR model

- Fill-in-the-blank exercise to test your understanding of critical parts of the DETR model workflow.

```

1 from torch import nn
2 class DETR(nn.Module):
3     def __init__(self, num_classes, hidden_dim=256, nheads=8,
4                 num_encoder_layers=6, num_decoder_layers=6, num_queries=100):
5         super().__init__()
6
7         # create ResNet-50 backbone
8         self.backbone = resnet50()
9         del self.backbone.fc
10
11        # create conversion layer
12        self.conv = nn.Conv2d(2048, hidden_dim, 1)
13
14        # create a default PyTorch transformer
15        self.transformer = nn.Transformer(
16            hidden_dim, nheads, num_encoder_layers, num_decoder_layers)
17
18        # prediction heads, one extra class for predicting non-empty slots
19        # note that in baseline DETR linear_bbox layer is 3-layer MLP
20        self.linear_class = nn.Linear(hidden_dim, num_classes + 1)
21        self.linear_bbox = nn.Linear(hidden_dim, 4)
22
23        # output positional encodings (object queries)
24        self.query_pos = nn.Parameter(torch.rand(100, hidden_dim))
25
26        # spatial positional encodings
27        # note that in baseline DETR we use sine positional encodings
28        self.row_embed = nn.Parameter(torch.rand(50, hidden_dim // 2))
29        self.col_embed = nn.Parameter(torch.rand(50, hidden_dim // 2))
30
31    def forward(self, inputs):
32        # propagate inputs through ResNet-50 up to avg-pool layer
33        x = self.backbone.conv1(inputs)
34        x = self.backbone.bn1(x)
35        x = self.backbone.relu(x)
36        x = self.backbone.maxpool(x)
37
38        x = self.backbone.layer1(x)
39        x = self.backbone.layer2(x)
40        x = self.backbone.layer3(x)
41        x = self.backbone.layer4(x)
42
43        # convert from 2048 to 256 feature planes for the transformer
44        h = self.conv(x)
45

```

```

46     # construct positional encodings
47     H, W = h.shape[-2:]
48     pos = torch.cat([
49         self.col_embed[:W].unsqueeze(0).repeat(H, 1, 1),
50         self.row_embed[:H].unsqueeze(1).repeat(1, W, 1),
51     ], dim=-1).flatten(0, 1).unsqueeze(1)
52
53     # propagate through the transformer
54     h = self.transformer(pos + 0.1 * h.flatten(2).permute(2, 0, 1),
55                         self.query_pos.unsqueeze(1)).transpose(0, 1)
56
57
58
59     # finally project transformer outputs to class labels and bounding boxes
60     pred_logits = self.linear_class(h)
61     pred_boxes = self.linear_bbox(h).sigmoid()
62
63     return {'pred_logits': pred_logits,
64           'pred_boxes': pred_boxes}

```

✓ Q2. Custom Image Detection and Attention Visualization

In this task, you will upload an **image of your choice** (different from the provided sample) and follow the steps below:

- Object Detection using DETR
 - Use the DETR model to detect objects in your uploaded image.
- Attention Visualization in Encoder
 - Visualize the regions of the image where the encoder focuses the most.
- Decoder Query Attention in Decoder
 - Visualize how the decoder's query attends to specific areas corresponding to the detected objects.

```

1 import math
2
3 from PIL import Image
4 import requests
5 import matplotlib.pyplot as plt
6 %config InlineBackend.figure_format = 'retina'
7
8 import ipywidgets as widgets
9 from IPython.display import display, clear_output
10
11 import torch
12 from torch import nn
13
14
15 from torchvision.models import resnet50
16 import torchvision.transforms as T
17 torch.set_grad_enabled(False);
18
19 # COCO classes
20 CLASSES = [
21     'N/A', 'person', 'bicycle', 'car', 'motorcycle', 'airplane', 'bus',
22     'train', 'truck', 'boat', 'traffic light', 'fire hydrant', 'N/A',
23     'stop sign', 'parking meter', 'bench', 'bird', 'cat', 'dog', 'horse',
24     'sheep', 'cow', 'elephant', 'bear', 'zebra', 'giraffe', 'N/A', 'backpack',
25     'umbrella', 'N/A', 'N/A', 'handbag', 'tie', 'suitcase', 'frisbee', 'skis',
26     'snowboard', 'sports ball', 'kite', 'baseball bat', 'baseball glove',
27     'skateboard', 'surfboard', 'tennis racket', 'bottle', 'N/A', 'wine glass',
28     'cup', 'fork', 'knife', 'spoon', 'bowl', 'banana', 'apple', 'sandwich',
29     'orange', 'broccoli', 'carrot', 'hot dog', 'pizza', 'donut', 'cake',
30     'chair', 'couch', 'potted plant', 'bed', 'N/A', 'dining table', 'N/A',

```

```

31     'N/A', 'toilet', 'N/A', 'tv', 'laptop', 'mouse', 'remote', 'keyboard',
32     'cell phone', 'microwave', 'oven', 'toaster', 'sink', 'refrigerator', 'N/A',
33     'book', 'clock', 'vase', 'scissors', 'teddy bear', 'hair drier',
34     'toothbrush'
35 ]
36
37 # colors for visualization
38 COLORS = [[0.000, 0.447, 0.741], [0.850, 0.325, 0.098], [0.929, 0.694, 0.125],
39           [0.494, 0.184, 0.556], [0.466, 0.674, 0.188], [0.301, 0.745, 0.933]]
40 # standard PyTorch mean-std input image normalization
41 transform = T.Compose([
42     T.Resize(800),
43     T.ToTensor(),
44     T.Normalize([0.485, 0.456, 0.406], [0.229, 0.224, 0.225])
45 ])
46
47 # for output bounding box post-processing
48 def box_cxcywh_to_xyxy(x):
49     x_c, y_c, w, h = x.unbind(1)
50     b = [(x_c - 0.5 * w), (y_c - 0.5 * h),
51          (x_c + 0.5 * w), (y_c + 0.5 * h)]
52     return torch.stack(b, dim=1)
53
54 def rescale_bboxes(out_bbox, size):
55     img_w, img_h = size
56     b = box_cxcywh_to_xyxy(out_bbox)
57     b = b * torch.tensor([img_w, img_h, img_w, img_h], dtype=torch.float32)
58     return b
59
60 def plot_results(pil_img, prob, boxes):
61     plt.figure(figsize=(16, 10))
62     plt.imshow(pil_img)
63     ax = plt.gca()
64     colors = COLORS * 100
65     for p, (xmin, ymin, xmax, ymax), c in zip(prob, boxes.tolist(), colors):
66         ax.add_patch(plt.Rectangle((xmin, ymin), xmax - xmin, ymax - ymin,
67                                   fill=False, color=c, linewidth=3))
68         cl = p.argmax()
69         text = f'{CLASSES[cl]}: {p[cl]:0.2f}'
70         ax.text(xmin, ymin, text, fontsize=15,
71               bbox=dict(facecolor='yellow', alpha=0.5))
72     plt.axis('off')
73     plt.show()
74
75

```

In this section, we show-case how to load a model from hub, run it on a custom image, and print the result. Here we load the simplest model (DETR-R50) for fast inference. You can swap it with any other model from the model zoo.

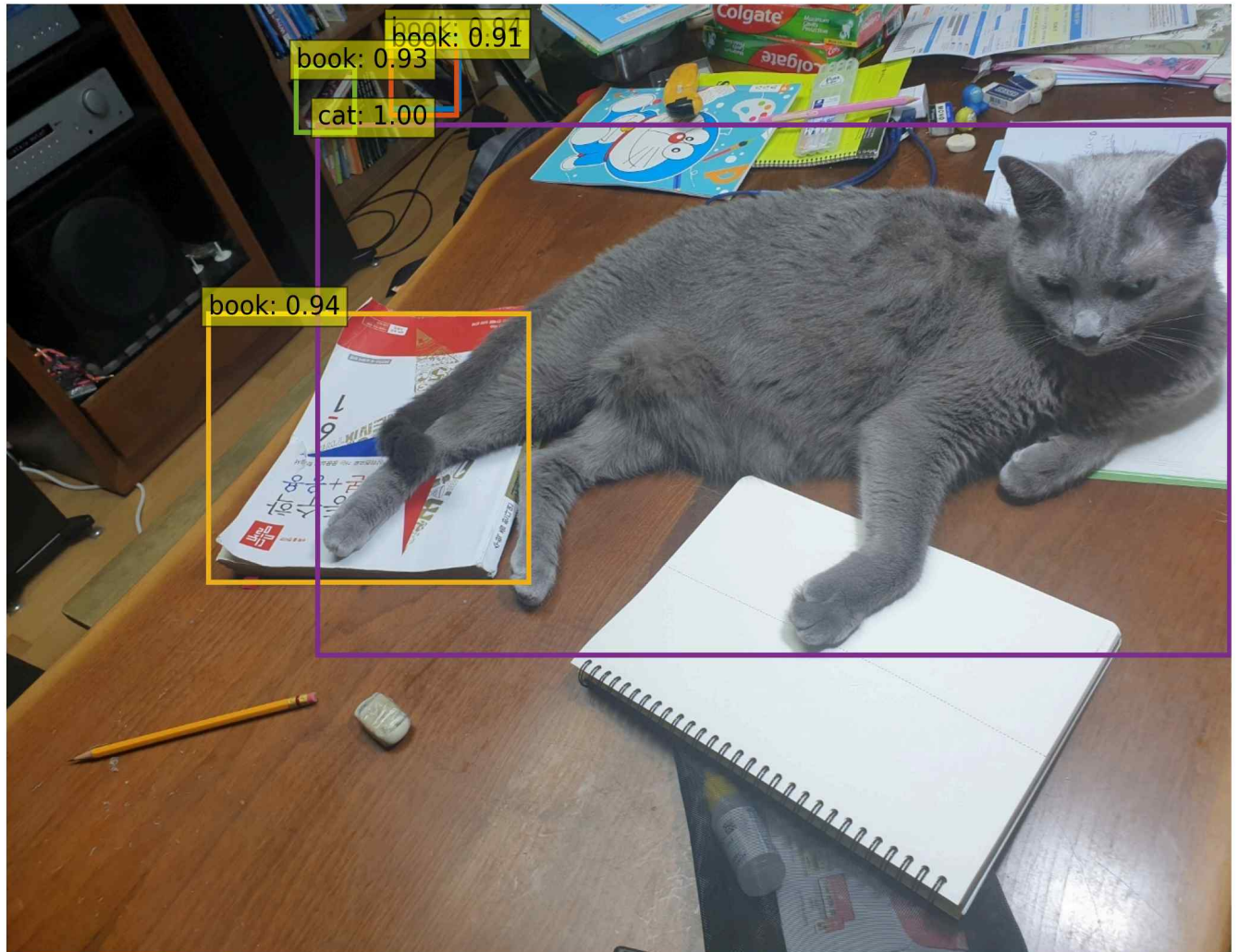
```

1 model = torch.hub.load('facebookresearch/detr', 'detr_resnet50', pretrained=True)
2 model.eval();
3
4 #url = 'http://images.cocodataset.org/val2017/000000039769.jpg'
5 #im = Image.open(requests.get(url, stream=True).raw) # put your own image
6 im = Image.open('/content/KakaoTalk_20241113_094243273.jpg')
7
8 # mean-std normalize the input image (batch-size: 1)
9 img = transform(im).unsqueeze(0)
10
11 # propagate through the model
12 outputs = model(img)
13

```

```
14 # keep only predictions with 0.7+ confidence
15 probas = outputs['pred_logits'].softmax(-1)[0, :, :-1]
16 keep = probas.max(-1).values > 0.9
17
18 # convert boxes from [0; 1] to image scales
19 bboxes_scaled = rescale_bboxes(outputs['pred_boxes'][0, keep], im.size)
20
21 # mean-std normalize the input image (batch-size: 1)
22 img = transform(im).unsqueeze(0)
23
24 # propagate through the model
25 outputs = model(img)
26
27 # keep only predictions with 0.7+ confidence
28 probas = outputs['pred_logits'].softmax(-1)[0, :, :-1]
29 keep = probas.max(-1).values > 0.9
30
31 # convert boxes from [0; 1] to image scales
32 bboxes_scaled = rescale_bboxes(outputs['pred_boxes'][0, keep], im.size)
33
34 # mean-std normalize the input image (batch-size: 1)
35 img = transform(im).unsqueeze(0)
36
37 # propagate through the model
38 outputs = model(img)
39
40 # keep only predictions with 0.7+ confidence
41 probas = outputs['pred_logits'].softmax(-1)[0, :, :-1]
42 keep = probas.max(-1).values > 0.9
43
44 # convert boxes from [0; 1] to image scales
45 bboxes_scaled = rescale_bboxes(outputs['pred_boxes'][0, keep], im.size)
46
47 plot_results(im, probas[keep], bboxes_scaled)
```

Using cache found in /root/.cache/torch/hub/facebookresearch_detr_main



Here we visualize attention weights of the last decoder layer. This corresponds to visualizing, for each detected objects, which part of the image the model was looking at to predict this specific bounding box and class.

```

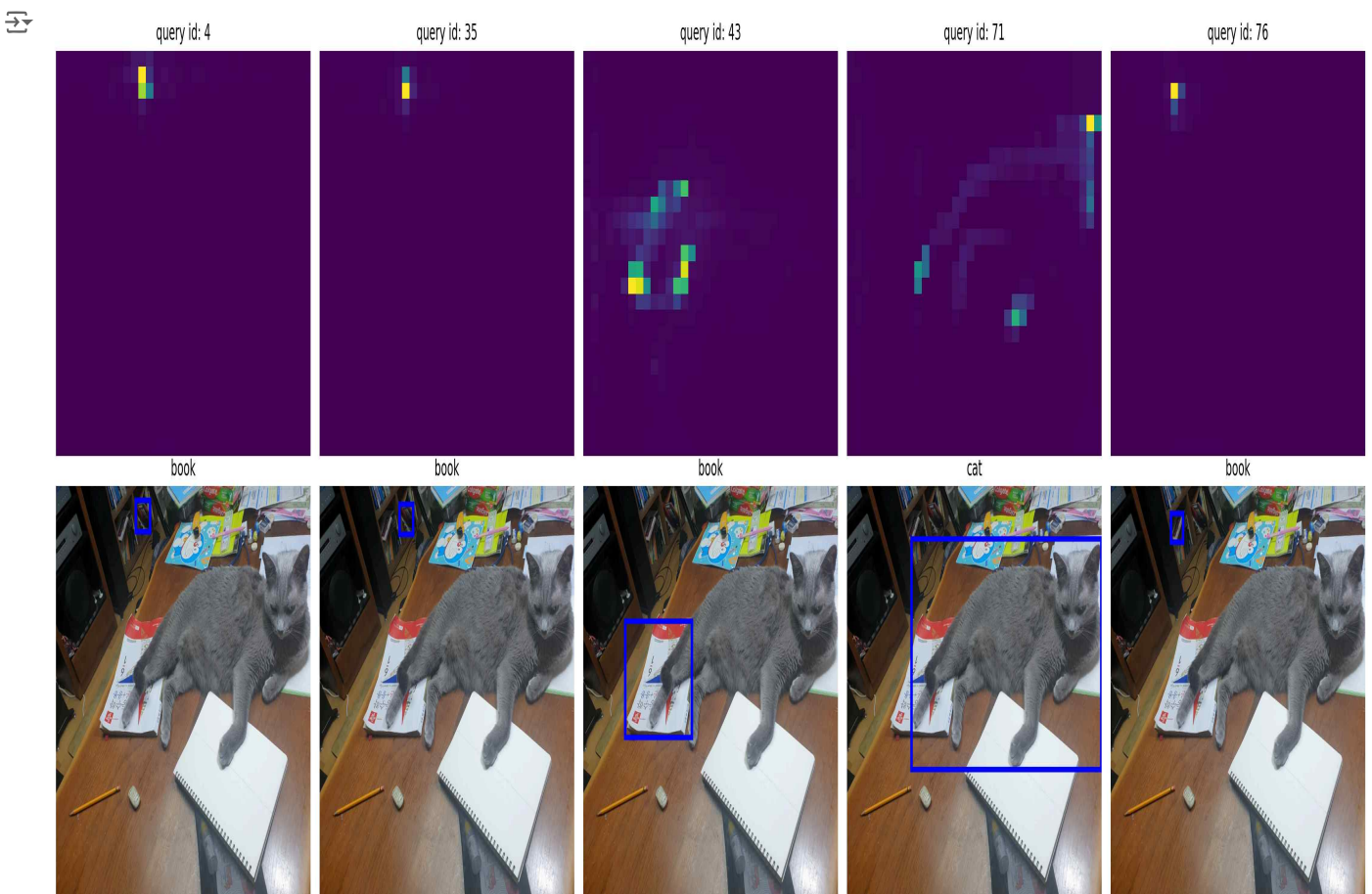
1 # use lists to store the outputs via up-values
2 conv_features, enc_attn_weights, dec_attn_weights = [], [], []
3
4 hooks = [
5     model.backbone[-2].register_forward_hook(
6         lambda self, input, output: conv_features.append(output)
7     ),
8     model.transformer.encoder.layers[-1].self_attn.register_forward_hook(
9         lambda self, input, output: enc_attn_weights.append(output[1])
10    ),
11    model.transformer.decoder.layers[-1].multihead_attn.register_forward_hook(
12        lambda self, input, output: dec_attn_weights.append(output[1])
13    ),
14 ]
15
16 # propagate through the model
17 outputs = model(img) # put your own image
18
19 for hook in hooks:
20     hook.remove()
21
22 # don't need the list anymore
23 conv_features = conv_features[0]
24 enc_attn_weights = enc_attn_weights[0]
25 dec_attn_weights = dec_attn_weights[0]

```

```

1 # get the feature map shape
2 h, w = conv_features['0'].tensors.shape[-2:]
3
4 fig, axs = plt.subplots(ncols=len(bboxes_scaled), nrows=2, figsize=(22, 7))
5 colors = COLORS * 100
6 for idx, ax_i, (xmin, ymin, xmax, ymax) in zip(keep.nonzero(), axs.T, bboxes_scaled):
7     ax = ax_i[0]
8     ax.imshow(dec_attn_weights[0, idx].view(h, w))
9     ax.axis('off')
10    ax.set_title(f'query id: {idx.item()}')
11
12    ax = ax_i[1]
13    ax.imshow(im)
14    ax.add_patch(plt.Rectangle((xmin, ymin), xmax - xmin, ymax - ymin,
15                             fill=False, color='blue', linewidth=3))
16    ax.axis('off')
17    ax.set_title(CLASSES[probas[idx].argmax()])
18 fig.tight_layout()

```



```

1 # output of the CNN
2 f_map = conv_features['0']
3 print("Encoder attention:      ", enc_attn_weights[0].shape)
4 print("Feature map:           ", f_map.tensors.shape)

```

```

Encoder attention:      torch.Size([850, 850])
Feature map:           torch.Size([1, 2048, 25, 34])

```

```

1 # get the HxW shape of the feature maps of the CNN
2 shape = f_map.tensors.shape[-2:]
3 # and reshape the self-attention to a more interpretable shape
4 sattn = enc_attn_weights[0].reshape(shape + shape)
5 print("Reshaped self-attention:", sattn.shape)

```

```

Reshaped self-attention: torch.Size([25, 34, 25, 34])

```



```

1 # downsampling factor for the CNN, is 32 for DETR and 16 for DETR DC5
2 fact = 32
3
4 # let's select 4 reference points for visualization
5 idxs = [(200, 200), (280, 400), (200, 600), (440, 800),]
6
7 # here we create the canvas
8 fig = plt.figure(constrained_layout=True, figsize=(25 * 0.7, 8.5 * 0.7))
9 # and we add one plot per reference point
10 gs = fig.add_gridspec(2, 4)
11 axs = [
12     fig.add_subplot(gs[0, 0]),
13     fig.add_subplot(gs[1, 0]),
14     fig.add_subplot(gs[0, -1]),
15     fig.add_subplot(gs[1, -1]),
16 ]
17
18 # for each one of the reference points, let's plot the self-attention
19 # for that point
20 for idx_o, ax in zip(idxs, axs):
21     idx = (idx_o[0] // fact, idx_o[1] // fact)
22     ax.imshow(sattn[..., idx[0], idx[1]], cmap='cividis', interpolation='nearest')
23     ax.axis('off')
24     ax.set_title(f'self-attention{idx_o}')
25
26 # and now let's add the central image, with the reference points as red circles
27 fcenter_ax = fig.add_subplot(gs[:, 1:-1])
28 fcenter_ax.imshow(im)
29 for (y, x) in idxs:
30     scale = im.height / img.shape[-2]
31     x = ((x // fact) + 0.5) * fact
32     y = ((y // fact) + 0.5) * fact
33     fcenter_ax.add_patch(plt.Circle((x * scale, y * scale), fact // 2, color='r'))
34     fcenter_ax.axis('off')

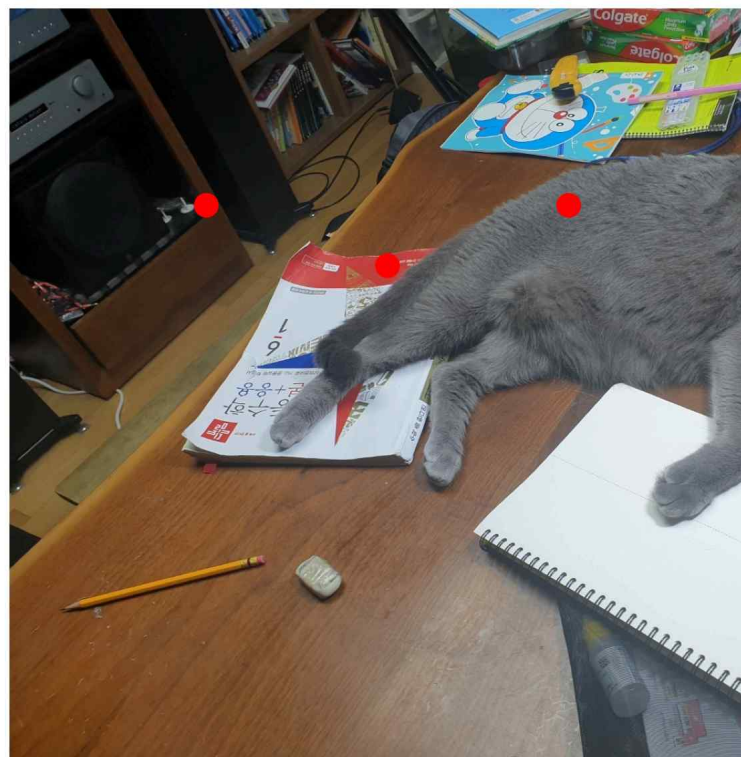
```



), 200)



), 400)



Q3. Understanding Attention Mechanisms

In this task, you focus on understanding the attention mechanisms present in the encoder and decoder of DETR.

- Briefly describe the types of attention used in the encoder and decoder, and explain the key differences between them.
- Based on the visualized results from Q2, provide an analysis of the distinct characteristics of each attention mechanism in the encoder and decoder. Feel free to express your insights.

✓ 1. Object Detection using DETR

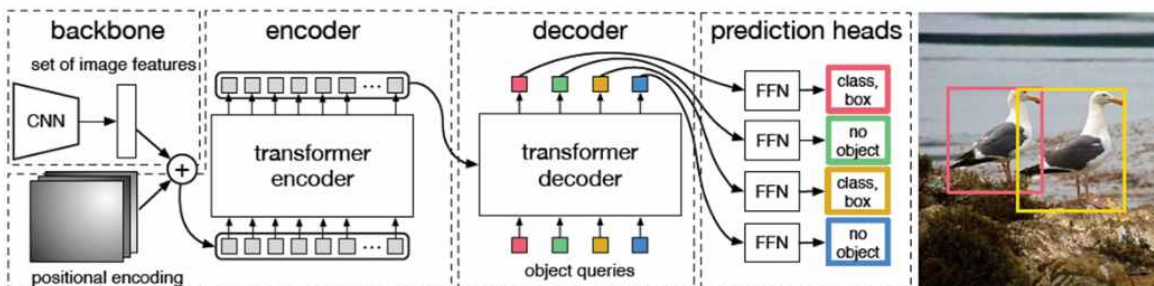
- **Input Processing:** We use ResNet(CNN Backbone) to split the image into feature maps(object detection), similar to the way transformers work with text by dividing sentences into tokens. This process provides a series of image features that can be input to the DETR model.
- **Encoder-Decoder Mechanism:** The DETR model processes the image features through an encoder-decoder transformer architecture.
 - **Encoder:** Processes the entire image to capture a global understanding of all regions.
 - **Decoder:** Uses learned queries to identify and localize specific objects within the image.
- **Output:** DETR produces bounding boxes around detected objects and assigns labels, highlighting what it identifies in the image.

2. About Attention

- **Self-Attention in Encoder:** Each patch of the image can attend to every other part, which allows the model to learn relationships between distant features (like an object's edges or parts). For example, if we see high attention around a specific area (like a cat or a book in the image), it indicates that the encoder recognizes these regions as significant for the entire image context.
- **Self-Attention in Decoder:** Each query specializes in detecting a particular object or feature in the image. Each query attends to other queries, allowing them to exchange information about what they are "seeing" in the image.
- **Cross-Attention (or Encoder-Decoder Attention):** The queries also attend to the encoder's output to focus on specific regions relevant to each query.

self-attention 400-800 is much more similar with original cat boundary. By visualization, i can get further explanation about using query **첨부 자료

DETR은 CNN Backbone + Transformer + FFN (Feed Forward Network) 로 구성되어있습니다. 매우 간단해서 PyTorch 로 구현한 코드도 굉장히 짧고 깔끔합니다.



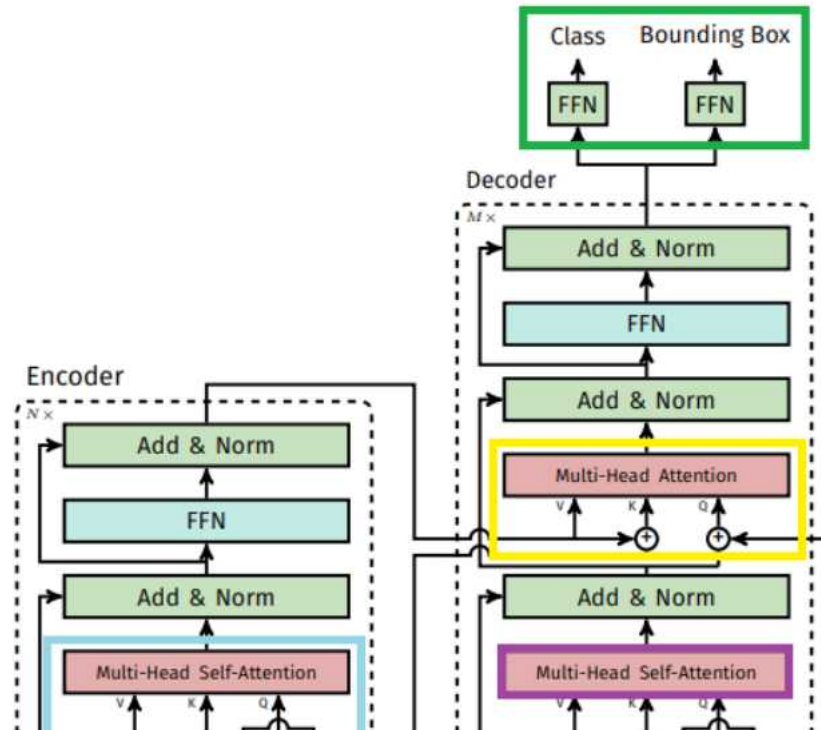
CNN Backbone

input image를 CNN Backbone을 통과시켜 feature map을 뽑아냅니다.

이 feature map이 transformer에 들어갈 수 있도록 처리를 해주어야 합니다.

- input image 크기는 $H_0 \times W_0$
- CNN을 통과하여 출력된 feature map은 $C \times H \times W$ (ResNet50을 사용하였기 때문에 $C=2048$, $H = H_0/32$, $W = W_0/32$)
- 1×1 convolution을 적용하여 $d \times H \times W$ 형태로 바꿈 ($C > d$)
- transformer에 들어가기 위해서는 2차원이어야 하므로, $d \times H \times W$ 의 3차원에서 $d \times HW$ 의 2차원으로 구조를 바꿈.

Transformer



Encoder

(파란색 박스) $d \times HW$ 의 feature matrix에 Positional encoding 정보를 더한 matrix를 multi-head self-attention에 통과시킵니다. Transformer의 특성 상 입력 matrix와 출력 matrix의 크기는 동일합니다.

Decoder

(분홍색 박스) N 개의 bounding box에 대해 N 개의 object query를 생성합니다. 초기 object query는 0으로 설정되어있습니다.

reCAPTCHA 서비스에 연결할 수 없습니다. 인터넷 연결을 확인한 후 페이지를 새로고침하여 reCAPTCHA 보안문자를 다시 로드하세요.