

## 3. Fine-Tuning Model Process

### 1. Introduction

In [\[1\]](#), We train the Image Segmentation and Classification model from `resnet-50` backbone, but we re-write a whole model for it to adapt to our dataset keys. And in that case, we re-write all heads and even a whole model. Which complicated the problem.

Refer to forward method In our model, we receive following parameters :

```
Python

1 def forward(self,
2     pixel_values: torch.Tensor,
3     pixel_mask: torch.Tensor,
4     class_ids: Optional[torch.Tensor] = None,
5     segmentation_mask: Optional[torch.Tensor] = None,
6     bboxes: Optional[torch.Tensor] = None,
7     **kwargs, # meta info (not used here)
8 ) -> FishSegmentationModelOutput:
```

#### (1) Model Construction

Actually, we can firstly simply modify `num_labels`,

```
Python

1 from transformers import DetrForObjectDetection
2
3 model = DetrForObjectDetection.from_pretrained(
4     "facebook/detr-resnet-50-panoptic",
5     num_labels=9, # number of fish species
6     ignore_mismatched_sizes=True, # avoid mismatched weighting size error
7 )
8 # or use following (use num_labels in config)
9 config = DetrConfig(num_labels=9, )
10 model = DetrForObjectDetection.from_pretrained(
11     "facebook/detr-resnet-50-panoptic", # number of fish species
12     ignore_mismatched_sizes=True,
13     config=config
14 )
```

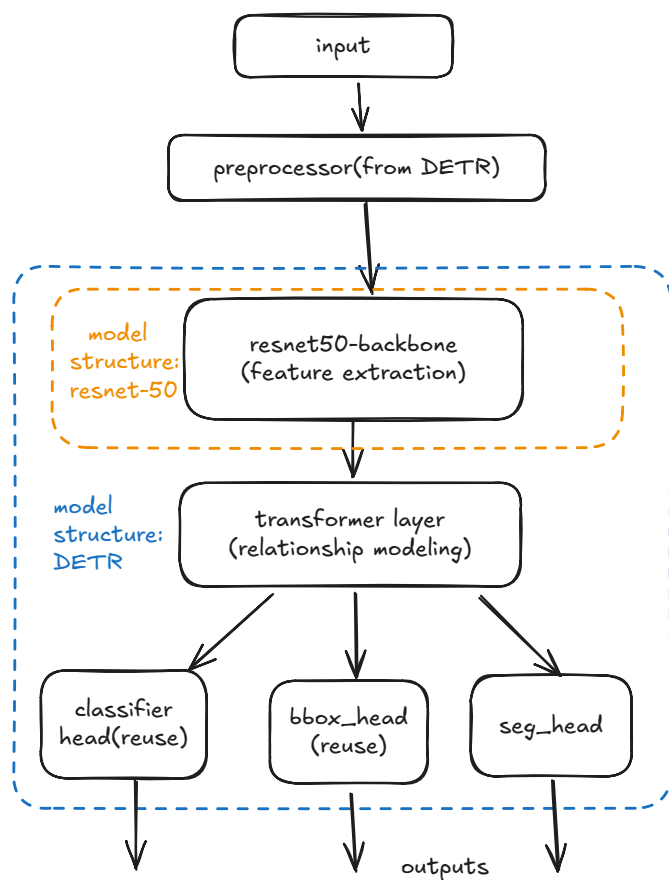
Then jump to the source code of `DetrForObjectDetection` to see what it needs :

```

1231 class DetrForObjectDetection(DetrPreTrainedModel):
1232 > def __init__(self, config: DetrConfig):...
1248
1249 @auto_docstring
1250 def forward(
1251     self,
1252     pixel_values: torch.FloatTensor,
1253     pixel_mask: Optional[torch.LongTensor] = None,
1254     decoder_attention_mask: Optional[torch.FloatTensor] = None,
1255     encoder_outputs: Optional[torch.FloatTensor] = None,
1256     inputs_embeds: Optional[torch.FloatTensor] = None,
1257     decoder_inputs_embeds: Optional[torch.FloatTensor] = None,
1258     labels: Optional[list[dict]] = None,
1259     output_attentions: Optional[bool] = None,
1260     output_hidden_states: Optional[bool] = None,
1261     return_dict: Optional[bool] = None,
1262 > ) → Union[tuple[torch.FloatTensor], DetrObjectDetectionOutput]:...

```

Compare to [ResNet50-Strucutre-Model](#), actual DETR model has a transformer layer



So if we just want to fine-tune, we **only need image as the input** :

```

Python
1 def test_model():
2     # just use pretrained DETR model and fine-tune it
3     model = DetrForObjectDetection.from_pretrained(
4         "facebook/detr-resnet-50-panoptic",
5         num_labels=9, # number of fish species
6         ignore_mismatched_sizes=True,
7     )
8     processor = DetrImageProcessor.from_pretrained("facebook/detr-resnet-50")
9
10    url = "http://images.cocodataset.org/val2017/000000039769.jpg"
11    image = Image.open(requests.get(url, stream=True).raw)
12    inputs = processor(images=[image, image], return_tensors="pt")
13    inputs.data.keys() # Out[25]: dict_keys(['pixel_values', 'pixel_mask'])
14    outputs = model(**inputs)
15    outputs.keys() # OrderedDictKeys(['logits', 'pred_boxes', 'last_hidden_state',
16    'encoder_last_hidden_state'])
17    outputs["logits"].shape # Out[22]: torch.Size([2, 100, 10])
18
19    # we can also control the output by `output_attention` parameter and
20    `output_hidden_states` parameter
21    outputs2 = model(**inputs, output_ attentions=True, output_hidden_states=True)
22    outputs2.keys() # Out[23]: OrderedDictKeys(['logits', 'pred_boxes', 'last_hidden_state',
23    'decoder_hidden_states', 'decoder_ attentions', 'cross_ attentions',
24    'encoder_last_hidden_state', 'encoder_hidden_states', 'encoder_ attentions'])
25
26    results = image_processor.post_process_object_detection(outputs,
27    target_sizes=torch.tensor([image.size[:2]]), threshold=0.3)

```

We note that **there's no labels in the inputs**, this is **because of what processor does**.

Then the output should be like this :

```

> {..} model = (DetrForObjectDetection) DetrForObjectDetection(\n (model): DetrModel(\n (backbor... View
> {..} output = {BatchFeature: 2} {'pixel_values': tensor([[[[ 0.2796, 0.3138, 0.3481, ..., -0.2856, -0.3... View
v {..} outputs = {DetrObjectDetectionOutput: 4} DetrObjectDetectionOutput(loss=None, loss_dict=Nor... View
1 auxiliary_outputs = (NoneType) None
1 cross_ attentions = (NoneType) None
1 decoder_ attentions = (NoneType) None
1 decoder_hidden_states = (NoneType) None
1 encoder_ attentions = (NoneType) None
1 encoder_hidden_states = (NoneType) None
> {..} encoder_last_hidden_state = {Tensor: (2, 850, 256)} tensor([[[[nan, nan, nan, ..., nan, 1...View as Array
> {..} last_hidden_state = {Tensor: (2, 100, 256)} tensor([[[[nan, nan, nan, ..., nan, nan, nan]...View as Array

```

## (2) Re-implement of dataloader function

Our new data function will return raw image (not convert RGB) and also, `mask`, `class_name` and `class_id`, It's relatively simpler than what we did before :



Python

```
1 class FishSegmentDataSet(Dataset):
2     def __init__(self,
3                 img_dirs: Union[str, List[str]],
4                 mask_dirs: Union[str, List[str]]):
5         super().__init__()
6         self.img_dirs: List[str] = [img_dirs] if isinstance(img_dirs, str) else img_dirs
7         self.mask_dirs: List[str] = [mask_dirs] if isinstance(mask_dirs, str) else
mask_dirs
8         if len(self.img_dirs) != len(self.mask_dirs):
9             raise ValueError("Number of image and mask directories must match")
10
11         self.items = []
12         for idx, (img_dir, mask_dir) in enumerate(zip(self.img_dirs, self.mask_dirs)):
13             img_paths = sorted(os.listdir(img_dir))
14             mask_paths = sorted(os.listdir(mask_dir))
15             class_name = os.path.basename(img_dir) # assuming img_dir is like
.../class_name
16             if len(img_paths) != len(mask_paths):
17                 raise ValueError(f"Number of images and masks must match in {img_dir} and
{mask_dir}")
18             for img_path, mask_path in zip(img_paths, mask_paths):
19                 img_id = os.path.splitext(img_path)[0]
20                 mask_id = os.path.splitext(mask_path)[0]
21                 if img_id != mask_id:
22                     raise ValueError(f"Image and mask file names must match: {img_id} vs
{mask_id}")
23                 self.items.append(dict(img_path=img_path,
24                                     mask_path=mask_path,
25                                     class_name=class_name,
26                                     class_id=idx))
27
28     def __len__(self):
29         return len(self.items)
30
31     def __getitem__(self, item_id):
32         record = self.items[item_id]
33         img = Image.open(os.path.join(self.img_dirs[record['class_id']],
record['img_path']))
34         mask = Image.open(os.path.join(self.mask_dirs[record['class_id']],
record['mask_path'])).convert("L")
35         # here we also return image_id for tracking
36         return dict(image_id = item_id,
37                     image=img,
38                     mask=torch.tensor(np.array(mask)).long(),
39                     class_id=record['class_id'],
40                     class_name=record['class_name'])
```

Then, we can use a simple `load_dataset()` method, just call:



Python

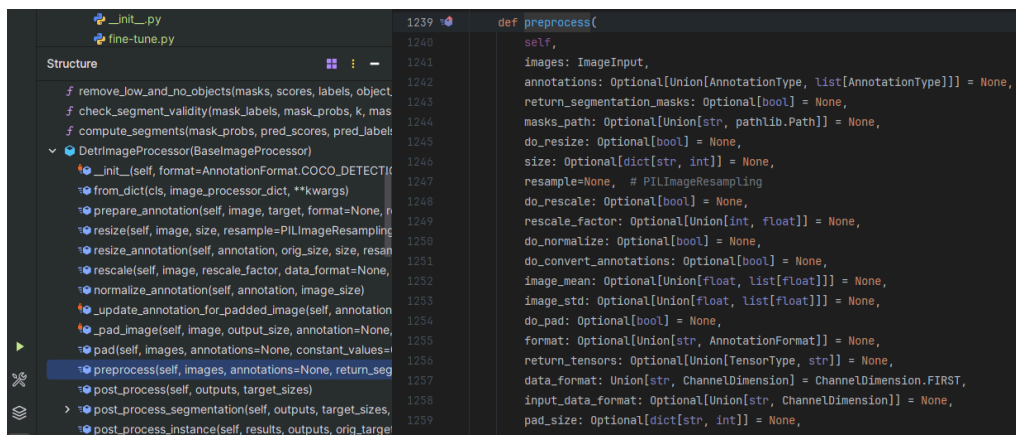
```

1  def load_dataset():
2      dataset_url = "crowww/a-large-scale-fish-dataset"
3      base_path = os.path.join(kagglehub.dataset_download(dataset_url), "Fish_Dataset",
                                "Fish_Dataset")
4      class_names = get_fish_classes()
5      img_dirs = [os.path.join(base_path, cls, cls) for cls in class_names]
6      mask_dirs = [os.path.join(base_path, cls, cls + " GT") for cls in class_names]
7      dataset = FishSegmentDataSet(img_dirs=img_dirs, mask_dirs=mask_dirs)
8      return dataset

```

### (3) Additional Parameters in Processor

For more precise parameters, we can refer to `DetrImageProcessor` <sup>[2]</sup>, the parameters are listed in `preprocess` method of `DetrImageProcessor`.



In the above picture, preprocessor receive those parameters, in previous chapter, we only use `images` :

```

Python
1     def preprocess(
2         self,
3         images: ImageInput,
4         annotations: Optional[Union[AnnotationType, list[AnnotationType]]] = None,
5         return_segmentation_masks: Optional[bool] = None,
6         masks_path: Optional[Union[str, pathlib.Path]] = None,
7         do_resize: Optional[bool] = None,
8         size: Optional[dict[str, int]] = None,
9         resample=None, # PILImageResampling
10        do_rescale: Optional[bool] = None,
11        rescale_factor: Optional[Union[int, float]] = None,
12        do_normalize: Optional[bool] = None,
13        do_convert_annotations: Optional[bool] = None,
14        image_mean: Optional[Union[float, list[float]]] = None,
15        image_std: Optional[Union[float, list[float]]] = None,
16        do_pad: Optional[bool] = None,
17        format: Optional[Union[str, AnnotationFormat]] = None,
18        return_tensors: Optional[Union[TensorType, str]] = None,
19        data_format: Union[str, ChannelDimension] = ChannelDimension.FIRST,
20        input_data_format: Optional[Union[str, ChannelDimension]] = None,
21        pad_size: Optional[dict[str, int]] = None,
22        **kwargs,
23    ) -> BatchFeatur

```

Now the `annotations` parameter is what we need. That's what should we input :

```

Python
1     """
2     annotations (`AnnotationType` or `list[AnnotationType]`, *optional*):
3         List of annotations associated with the image or batch of images. If annotation is
4         for object
5             detection, the annotations should be a dictionary with the following keys:
6             - "image_id" (`int`): The image id.
7             - "annotations" (`list[Dict]`): List of annotations for an image. Each annotation
8             should be a
9             dictionary. An image can have no annotations, in which case the list should be
10            empty.
11            If annotation is for segmentation, the annotations should be a dictionary with the
12            following keys:
13            - "image_id" (`int`): The image id.
14            - "segments_info" (`list[Dict]`): List of segments for an image. Each segment
15            should be a dictionary.
16            An image can have no segments, in which case the list should be empty.
17            - "file_name" (`str`): The file name of the image.
18    """

```

## 1) COCO-format annotations

The DETR model use **standard COCO dataset format** for annotations<sup>[3]</sup>, which is

randomly set to at least 480 and at most 800 pixels. At inference time, the shortest side is set to 800. One can use `DetrImageProcessor` to prepare images (and optional annotations in COCO format) for the model. Due to this resizing, images in a batch can have different sizes. DETR solves this by padding images up to

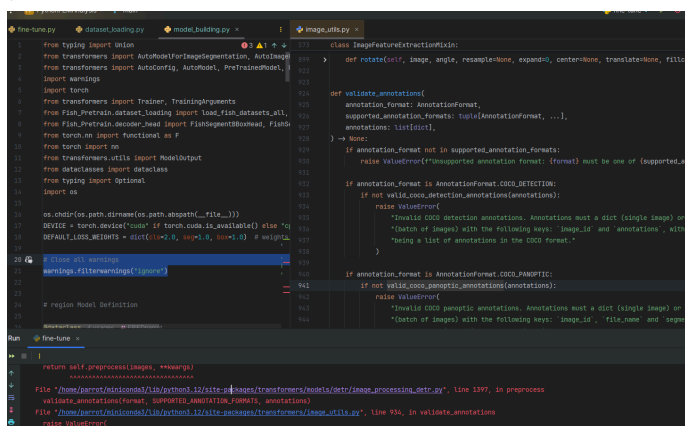
**Standard Coco dataset format annotations** are given as :

```
Python

1 # Construct annotations for object detection
2 "annotations": [
3     {
4         "id": 0,
5         "image_id": 0,
6         "category_id": 2,
7         "bbox": [
8             45,
9             2,
10            85,
11            85
12        ],
13        "area": 7225,
14        "segmentation": [],
15        "iscrowd": 0
16    },
17    {
18        "id": 1,
19        "image_id": 0,
20        "category_id": 2,
21        "bbox": [
22            324,
23            29,
24            72,
25            81
26        ],
27        "area": 5832,
28        "segmentation": [],
29        "iscrowd": 0
30    }
31 ]
```

But we have a easier way to find the things out, we can trace back and jump to `valid_coco_detection_annotations` and `valid_coco_panoptic_annotations` to check how it really is

:



The **source of checking coco-format functions** are pasted as follows :

```

Python
1  # **Detection-only model** → ignores segmentation
2  def is_valid_annotation_coco_detection(annotation: dict[str, Union[list, tuple]]) -> bool:
3      if (
4          isinstance(annotation, dict)
5          and "image_id" in annotation
6          and "annotations" in annotation
7          and isinstance(annotation["annotations"], (list, tuple))
8          and (
9              # an image can have no annotations
10             len(annotation["annotations"]) == 0 or isinstance(annotation["annotations"]
11             [0], dict)
12         ):
13         return True
14     return False
15
16 def is_valid_annotation_coco_panoptic(annotation: dict[str, Union[list, tuple]]) -> bool:
17     if (
18         isinstance(annotation, dict)
19         and "image_id" in annotation
20         and "segments_info" in annotation
21         and "file_name" in annotation
22         and isinstance(annotation["segments_info"], (list, tuple))
23         and (
24             # an image can have no segments
25             len(annotation["segments_info"]) == 0 or
26             isinstance(annotation["segments_info"][0], dict)
27         ):
28         return True
29     return False

```

## 2) How to use in our Model

Our goal is to predict both class, bounding box and binary segmentation mask, but in Processor **We use Object detection annotation format**, so **the processor information** includes **classification and bounding boxes**.

This is because Processor just define what data to keep, but **actually the Model defines what loss to calculate**, so for 1-object with mask prediction, we can use `detection processor`  
+ `segmentation model`:

### ✓ Model Inputs

- `DetrForObjectDetection`
  - Only looks at "annotations" with `bbox`, `category_id`, `area`, `iscrowd`.
  - Ignores any `segmentation` field.
  - Loss = classification + box (L1 + GloU).
- `DetrForSegmentation` ( `facebook/detr-resnet-50-panoptic` ):
  - Extends the object detection setup with a **mask head**.



- When you call the model with `labels` (from processor), it expects:
  - `labels["class_labels"]` (categories),
  - `labels["boxes"]` (bboxes),
  - `labels["mask_labels"]` **(segmentation masks)**.
- Loss = classification + box + segmentation loss.

We also need to set `return_segmentation_masks` to `True`

```

Python
1  """
2  return_segmentation_masks (`bool`, *optional*, defaults to
   self.return_segmentation_masks):
3      Whether to return segmentation masks.
4  """

```

Note there's another important parameter `format`, which is defined as :

```

Python
1  from transformers.image_utils import AnnotationFormat
2  class AnnotationFormat(ExplicitEnum):
3      COCO_DETECTION = "coco_detection"
4      COCO_PANOPTIC = "coco_panoptic"

```

Also, note the `segmentation` should be formatted into **COCO detection style (which is `[x1,y1, x2,y2,...]`)**. So when we use `numpy` to input, we got `Exception: input type is not supported`.

Firstly, we use a coco sample for test

```

Python
1  "annotations": [
2      {
3          "image_id": torch.tensor([f"class_id"]),
4          "category_id": torch.tensor([f"class_id"]),
5          "bbox": torch.tensor(bbox), # only gives bbox and area
6          "area": torch.tensor([area]),
7          # just make a square for test
8          "segmentation": [[bbox[0], bbox[1], bbox[0] + bbox[2], bbox[1], bbox[0] + bbox[2],
   bbox[1] + bbox[3], bbox[0], bbox[1] + bbox[3]]], # COCO polygon format
9          "iscrowd": torch.tensor([0]),
10     }
11 ]

```

The return result shape is :

```

Python
1  {
2      "pixel_values": Tensor[B, 3, H, W],      # normalized/resized images
3      "pixel_mask": Tensor[B, H, W],          # padding mask (1 = valid, 0 = padded)
4      "labels": [                             # list of dicts, one per image
5          {
6              "size": Tensor[2],               # processed (resized) image size
7              "orig_size": Tensor[2],          # original H, W
8              "image_id": Tensor[1],           # image id
9              "class_labels": Tensor[num_obj], # categories for each object
10             "boxes": Tensor[num_obj, 4],      # normalized boxes in cxcywh format
11             "area": Tensor[num_obj],          # object area
12             "iscrowd": Tensor[num_obj],       # crowd flags
13             "masks": Tensor[num_obj, H, W],   # binary masks aligned to pixel_values
14         },
15         ...
16     ]
17 }

```

We prepare to reuse `resize_mask` function in [\[4\]](#) add `masks` after processing to **simulate the calculation process**

So firstly, we test the shape of output :

```

Python
1  for batch in dataloader:
2      print(batch)
3      print(batch.keys())
4      print(batch['labels'][0].keys())
5      print(batch.get('labels')[0]["masks"].shape) # torch.Size([1, 800, 1060])

```

### 3) Full Implement of data loader

#### ⚠ sort directory name after

`os.listdir()` in Python **does not sort folders**

`os.listdir()` **in Python** does not sort folders or files automatically\*\*.

It returns the list of filenames in arbitrary order, which depends on the underlying filesystem and OS implementation.

So, the order is not guaranteed and should be considered "random"!

We get images from dataloader,



Python

```

1  class FishCollator: # (DataCollator)
2      def __init__(self, processor: DetrImageProcessor):
3          self.processor = processor
4
5      def __call__(self, features: List[Dict[str, Any]]) -> BatchFeature:
6          images = [f["image"] for f in features]
7          annotations = []
8
9          # we may have multiple objects in one image, but here in dataset we use only one
          # object per image
10         for f in features:
11             # prepare annotations from masks
12             bbox = _compute_bbox_from_mask(np.array(f["mask"]))
13             area = int(torch.prod(torch.tensor(bbox[2:])))
14             ann: AnnotationType = {
15                 "image_id": torch.tensor(f["image_id"]),
16                 "annotations": [
17                     {
18                         "image_id": torch.tensor(f["image_id"]), # not set to class_id,
19                         "category_id": torch.tensor(f["class_id"]),
20                         "bbox": torch.tensor(bbox), # only gives bbox and area
21                         "area": torch.tensor(area),
22                         # for COCO polygon format,
23                         # "segmentation": [[bbox[0], bbox[1], bbox[0] + bbox[2], bbox[1],
24                         #                  bbox[0] + bbox[2], bbox[1] + bbox[3], bbox[0], bbox[1] + bbox[3]]],
25                         "iscrowd": torch.tensor([0]),
26                     }
27                 ]
28             }
29             annotations.append(ann)
30
31             # it is a dict, but here we use batch feature for convenience
32             out = self.processor(images,
33                                 annotations=annotations,
34                                 return_segmentation_masks=False, # no mask info
35                                 format=AnnotationFormat.COCO_DETECTION,
36                                 return_tensors="pt")
37
38             # since out now have no masks, we need to add them back
39             masks = [f["mask"] for f in features]
40
41             for i, mask in enumerate(masks):
42                 # resize mask to the size of out['labels'][i]['size']
43                 target_size = tuple(out['labels'][i]['size'])
44                 mask_resize = resize_mask(mask, size=target_size)
45                 out['labels'][i]["masks"] = torch.tensor(mask_resize,
46                                                             dtype=torch.float32).unsqueeze(0) # add channel dim
47
48         return out

```

## Hint

At first I didn't discover that code completion and accidentally made `image_id` :

f["class\_id"] , This is a very critical issue. This will affects Hungarian Matching

Here is an reference :

### 1 How `image_id` affects Hungarian Matching

DETR's training loss is built around **set-based bipartite matching** — this is what makes DETR unique.

For every **image**, it:

1. Takes the **model predictions** (e.g., 100 queries per image → 100 predicted boxes).
2. Takes the **ground-truth annotations** for that same image.
3. Runs a **Hungarian algorithm** to find the *best matching* between predicted boxes and real boxes — minimizing a cost that combines classification, bounding box, and mask loss.

Then, it computes the losses **per image**.

So actually the image

In that case,

```
Python
1 print(batch.get('labels')[0]["masks"].shape) # it also returns torch.Size([1, 800, 1060])
```

Then we can test the output by using following :

```
Python
1 print(model(**batch)) # also, we can notice "pred_masks" is in the prediction, so the mask loss is correctly calculated
```

## 2. Train the fine-tune pretrained model

### (1) Train Code Implement

After implementing the above data collator code, the Process of fine-tuning model on Fish Dataset would become much simple than before, **all the backbones, transformers, heads, and loss functions** are reused :

#### Note

Be careful of `remove_unused_column` parameter, It should be set as `False` in order to keep the data transferred from the `collator` , or it will lose columns.



Python

```
1 def train_model():
2     # just use pretrained DETR model and fine-tune it
3     config = DetrConfig(
4         model_name="Fish_Segmentation_Model_Fine_Tuning_DETR",
5         num_labels=9
6     )
7     model = DetrForSegmentation.from_pretrained(
8         "facebook/detr-resnet-50-panoptic", # number of fish species
9         ignore_mismatched_sizes=True,
10        config=config,
11    )
12    processor = DetrImageProcessor.from_pretrained("facebook/detr-resnet-50")
13
14    train_dataset, val_dataset = load_dataset(train_size=0.8, seed=42)
15    collator = FishCollator(processor=processor)
16
17    training_args = TrainingArguments(
18        output_dir="./model",
19        per_device_train_batch_size=8,
20        per_device_eval_batch_size=8,
21        num_train_epochs=20,
22        logging_steps=10,
23        save_steps=50,
24        save_total_limit=2,
25        remove_unused_columns=False,
26        eval_strategy="steps",
27        learning_rate=5e-5,
28        weight_decay=0.01,
29        warmup_steps=100,
30    )
31
32    trainer = Trainer(
33        model=model,
34        args=training_args,
35        train_dataset=train_dataset,
36        eval_dataset=val_dataset,
37        data_collator=collator,
38    )
39
40    # trainer.train()
41    trainer.save_model('./models/final_model')
42    processor.save_pretrained('./models/final_model')
43    config.save_pretrained('./models/final_model')
44
45    repo_id = "FriedParrot/fish-segmentation-simple"
46    # push model to hub
47    backup_model_to_hub(repo_id, config, model, processor)
48
49    print("model training finished")
```

## Training Parameter Settings

Since this is a very large pretrained model, 1e-4 learning rate is too large for this model to converge. So we can use smaller learning rate (5e-5) and slightly increase the epoch num (to 20). And It's better for model training.

Here how to implement `backup_model_to_hub` has been given in [4-1].

## (2) Query-based

And example for how to evaluate model are also given in this chapter. But note for `detr` model, its return data includes :

```
Plain text
1 'loss', 'loss_dict', 'logits', 'pred_boxes', 'pred_masks', 'last_hidden_state',
  'encoder_last_hidden_state'
```

Since the DETR infers multiple predictions ( $N$  as is said in essay), Where the `logits` is of following shape (default `num_queries` is 100) :

```
Python
1 (batch_size, num_queries, num_classes + 1)
```

so actually, when we get `logits` from the network output, we got an output shape of `(batch_size, 100, 10)`

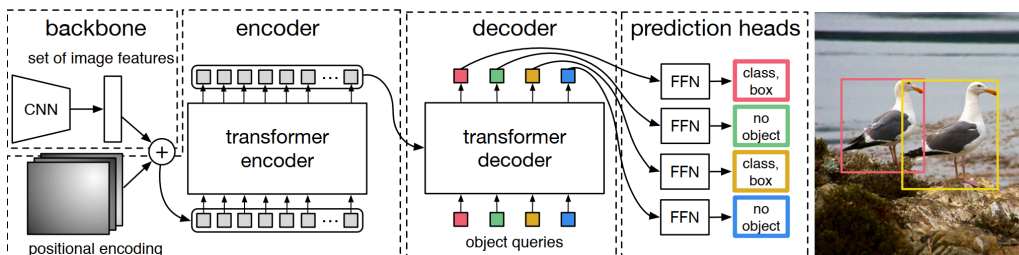
We need to know more about the output information of DERT. So for a brief view of the prediction principles of DETR network, we can refer to [5].

DETR use a direct set prediction approach to bypass the surrogate tasks. Which adopts the sequence-prediction effective transformers architecture. It simplifies the detection pipeline by dropping multiple hand-designed components that encode prior knowledge.

For multi-object detectors, DETR use Bipartite Matching Loss [6]. with following loss function :

$$\hat{\sigma} = \arg \min_{\sigma \in \Sigma} \sum_{i=1}^n \mathcal{L}_{match}(y_i, \hat{y}_{\sigma(i)}) \quad (2.2.1)$$

The ensemble structure of a DETR is as follows :



So the result is  $N$  predictions, and for `facebook/detr-resnet-50-panoptic` model, it also has its built-in mask head to predict the mask of each prediction.

Also, to predict the `background`, we needn't assign the background by ourselves. And when DETR predicts classes, it **automatically appends** a special **“no-object” (background)** class.

so we use following config :

```
1 config = DetrConfig(
2     model_name="Fish_Segmentation_Model_Fine_Tuning_DETR",
3     num_labels=9 # use number of Fish species (no need to +1 for background in DETR
4     segmentation)
```

Then DETR internally reserves one more class :

```
total classes = num_classes + 1
```

### (3) Problems in loading model

When we load model, we may encounter following problem :

```
Shell
```

```
1 File "/lib/python3.12/site-packages/safetensors/torch.py", line 381, in load_file with  
safe_open(filename, framework="pt", device=device) as f:  
^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^ safetensors_rust.SafetensorError: Error  
while deserializing header: header too small
```

This is because hugging face will cache the model weights, and it may load failed. So we can execute like following command to clear cache.

```
1  rm -rf ~/.cache/huggingface/hub
```

### 3. Evaluation

Evaluation Code are given as follows :

What to care :

- Since the input mask is



Python

```
1 def evaluate_model(
2     model, processor, batch_size: int = 1, max_vis: int = 2, save_dir: str = None
3 ):
4     _, test_set = load_dataset(train_size=0.95)
5     device = torch.device("cuda" if torch.cuda.is_available() else "cpu")
6     test_loader = DataLoader(test_set, batch_size=batch_size, shuffle=False,
7                             num_workers=2,
8                             collate_fn=FishCollator(processor), pin_memory=True if
9                             torch.cuda.is_available() else False)
10    classes = get_fish_classes()
11    model.eval()
12    model.to(device)
13
14    with torch.no_grad():
15        show_num = 0
16        correct = 0
17        for index, batch in enumerate(test_loader):
18            # since there is nested structure in inputs, only pin_memory for tensors is
19            # not enough
20            inputs = move_batch_to_device(batch, device)
21            outputs = model(**inputs)
22
23            images = inputs['pixel_values']
24            labels = [label for label in inputs['labels']]
25
26            # ===== predict result (written by GPT) =====
27            class_pred = outputs["logits"].argmax(-1) # (B, num_queries)
28            pred_masks = outputs["pred_masks"].sigmoid() # (B, num_queries, H, W)
29            pred_bboxes = outputs["pred_boxes"] # (B, num_queries, 4)
30            if show_num >= max_vis:
31                print("calculating accuracy ...", "batch: ", index, "/", len(test_loader))
32
33            batch_size = images.shape[0] # the bbox is predicted from mask
34            # calculate the classification accuracy -> gt : ground truth
35            for i in range(batch_size):
36                img = images[i] # shape [3, H, W]
37                label = labels[i] # predicted labels
38                # --- Handle predicted class, masks and bboxes (they should all be 1)---
39                class_ids = label['class_labels'] # input class ids
40                masks = label['masks'] # [B, 800, 1060]
41                bboxes = label['boxes'] # NOTE : the bbox is normalized
42                if len(class_ids) > 1:
43                    print(f"Warning: Multiple classes detected in sample {i}, using the
44                    first.")
45                if len(masks) > 1:
46                    print(f"Warning: Multiple masks detected in sample {i}, using the
47                    first.")
48                if len(bboxes) > 1:
49                    print(f"Warning: Multiple boxes detected in sample {i}, using the
50                    first.")
51
52            gt_label = class_ids[0] if len(class_ids) > 0 else None
53            gt_mask = masks[0] if len(masks) > 0 else None
```

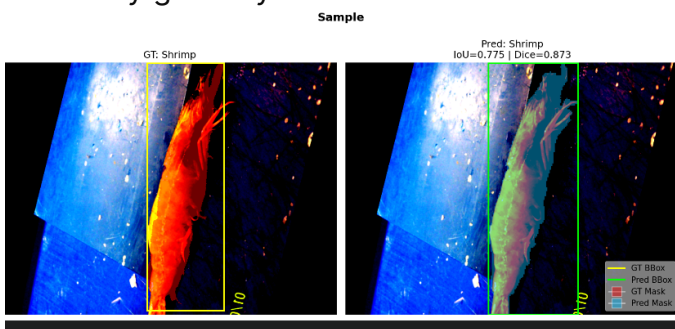


```

48         gt_box = bboxes[0] if len(bboxes) > 0 else None
49
50     # --- Prediction Extraction ---
51     valid_idx = class_pred[i] != model.config.num_labels
52     filtered_classes = class_pred[i][valid_idx]
53     filtered_masks = (pred_masks[i][valid_idx] > 0.5).long() # for sigmoid
function, we need to compare to 0.5
54     filtered_boxes = pred_bboxes[i][valid_idx]
55
56     # collate predicted results
57     pred_label = filtered_classes[0] if filtered_classes.nelement() > 0 else
None
58
59     if pred_label == gt_label:
60         correct += 1
61
62     if show_num >= max_vis:
63         continue
64     else:
65         show_num += 1
66     pred_mask = filtered_masks[0] if filtered_masks.nelement() > 0 else None
67     pred_box = filtered_boxes[0] if filtered_boxes.nelement() > 0 else None
68     # get label names
69     gt_label_name = classes[gt_label] if gt_label is not None else "N/A"
70     pred_label_name = classes[pred_label] if pred_label else "N/A"
71     visualize_sample_comparison(
72         img.cpu(),
73         gt_mask.cpu(), pred_mask.cpu(),
74         gt_label_name, pred_label_name,
75         gt_box.cpu(), pred_box.cpu(),
76         save_path=save_dir,
77         bbox_mode="center",
78     )
79     # compute accuracy
80     total = len(test_set)
81     accuracy = correct / total # total number of samples in the test set
82     print(f"Classification Accuracy on Test Set: {accuracy*100:.2f}%
({correct}/{total})")

```

We finally get very accurate result like following images :



1. Entrance of Hugging Face and pre-trained model ↩

2. [https://huggingface.co/docs/transformers/v4.56.2/en/model\\_doc/detr](https://huggingface.co/docs/transformers/v4.56.2/en/model_doc/detr) ↩

3. <https://www.v7labs.com/blog/coco-dataset-guide> ↩
4. 1. Entrance of Hugging Face and pre-trained model ↩ ↩
5. End-to-End Object Detection with Transformers ↩
6. Bipartite Matching Loss ↩