

```
import os
current_directory = os.getcwd()
```

```
import torch
import logging
import data
from models import imagebind_model
from models.imagebind_model import ModalityType, load_module
from models import lora as LoRA

logging.basicConfig(level=logging.INFO, force=True)

lora = False
linear_probing = False
device = "cuda:0" if torch.cuda.is_available() else "cpu"
load_head_post_proc_finetuned = True

assert not (linear_probing and lora), \
    "Linear probing is a subset of LoRA training procedure for ImageBind. " \
    "Cannot set both linear_probing=True and lora=True. "

if lora and not load_head_post_proc_finetuned:
    lora_factor = 12 / 0.07
else:
    lora_factor = 1
```

```

/usr/local/lib/python3.10/dist-packages/torchvision/transforms/_functional_video.py:10:
  warnings.warn(
/usr/local/lib/python3.10/dist-packages/torchvision/transforms/_transforms_video.py:10:
  warnings.warn(

```

```
image_paths = ["test_data/00088_2.jpg", "test_data/00095.jpg", "test_data/00158_4.jpg",
               "test_data/00169_3.jpg", "test_data/00220.jpg", "test_data/00222_3.jpg"]
audio_paths = ["test_data/00088_2.wav", "test_data/00095.wav", "test_data/00158_4.wav",
               "test_data/00169_3.wav", "test_data/00220.wav", "test_data/00222_3.wav"]
```

```
for i in range(len(image_paths)):
    image_paths[i] = current_directory + "/" + image_paths[i]
for i in range(len(audio_paths)):
    audio_paths[i] = current_directory + "/" + audio_paths[i]
```

```
model = imagebind_model.imagebind_huge(pretrained=True)
```

```
# if lora:
#     model.modality_trunks.update(
#         LoRA.apply_lora_modality_trunks(model.modality_trunks, rank=4,
#                                         modality_names=[ModalityType.TEXT, ModalityType.VISUAL])
#     )
#     LoRA.load_lora_modality_trunks(model.modality_trunks,
#                                     checkpoint_dir=".checkpoints/lora/", postfix="_lora")
#
#     if load_head_post_proc_finetuned:
#         load_module(model.modality_postprocessors, module_name="postprocessors",
#                     checkpoint_dir=".checkpoints/lora/", postfix="_last")
#         load_module(model.modality_heads, module_name="heads",
#                     checkpoint_dir=".checkpoints/lora/", postfix="_last")
# elif linear_probing:
#     load_module(model.modality_heads, module_name="heads",
#                 checkpoint_dir=".checkpoints/lora/", postfix="_last")
```

```
model.eval()  
model.to(device)
```

```

        (fc1): Linear(in_features=768, out_features=3072, bias=True)
        (act): GELU(approximate='none')
        (fc2): Linear(in_features=3072, out_features=768, bias=True)
        (drop): Dropout(p=0.0, inplace=False)
    )
    (norm_2): LayerNorm((768,), eps=1e-06, elementwise_affine=True)
)
(6): BlockWithMasking(
  (attn): MultiheadAttention(
    (out_proj): NonDynamicallyQuantizableLinear(in_features=768,
out_features=768, bias=True)
  )
  (drop_path): DropPath(drop_prob=0.055)
  (norm_1): LayerNorm((768,), eps=1e-06, elementwise_affine=True)
  (mlp): Mlp(
    (fc1): Linear(in_features=768, out_features=3072, bias=True)
    (act): GELU(approximate='none')
    (fc2): Linear(in_features=3072, out_features=768, bias=True)
    (drop): Dropout(p=0.0, inplace=False)
  )
  (norm_2): LayerNorm((768,), eps=1e-06, elementwise_affine=True)
)
(7): BlockWithMasking(
  (attn): MultiheadAttention(
    (out_proj): NonDynamicallyQuantizableLinear(in_features=768,
out_features=768, bias=True)
  )
  (drop_path): DropPath(drop_prob=0.064)
  (norm_1): LayerNorm((768,), eps=1e-06, elementwise_affine=True)
  (mlp): Mlp(
    (fc1): Linear(in_features=768, out_features=3072, bias=True)
    (act): GELU(approximate='none')
    (fc2): Linear(in_features=3072, out_features=768, bias=True)
    (drop): Dropout(p=0.0, inplace=False)
  )
  (norm_2): LayerNorm((768,), eps=1e-06, elementwise_affine=True)
)
(8): BlockWithMasking(
  (attn): MultiheadAttention(
    (out_proj): NonDynamicallyQuantizableLinear(in_features=768,
out_features=768, bias=True)
  )
  (drop_path): DropPath(drop_prob=0.073)
  (norm_1): LayerNorm((768,), eps=1e-06, elementwise_affine=True)
  (mlp): Mlp(

```

```
# Load data
#audio_query_paths = [current_directory + "/test_data/00222_3.wav"] # Example: Choos

inputs_audio = {
    ModalityType.AUDIO: data.load_and_transform_audio_data(audio_paths, device),
}

inputs_image = {
    ModalityType.VISION: data.load_and_transform_vision_data(image_paths, device),
}

# Generate embeddings for audio query
with torch.no_grad():
    embeddings_audio = model(inputs_audio)
with torch.no_grad():
    embeddings_image = model(inputs_image)

vision_embeddings = embeddings_image[ModalityType.VISION]
audio_query_embedding = embeddings_audio[ModalityType.AUDIO]

count_crt=0
for audio_index, audio_query_embedding in enumerate(embeddings_audio['audio']):

    similarity_scores = torch.softmax(vision_embeddings @ audio_query_embedding.T * (10

    most_relevant_image_index = torch.argmax(similarity_scores)

    most_relevant_image_path = image_paths[most_relevant_image_index]

    if most_relevant_image_index==audio_index:
        count_crt+=1

    print("Audio query:",audio_paths[audio_index].split('/')[ -1] )
    print("Most relevant image for audio query:", most_relevant_image_path.split('/')[ -1]
```

```
➞ Audio query: 00088_2.wav
Most relevant image for audio query: 00222_3.jpg
Audio query: 00095.wav
Most relevant image for audio query: 00095.jpg
Audio query: 00158_4.wav
Most relevant image for audio query: 00169_3.jpg
Audio query: 00169_3.wav
Most relevant image for audio query: 00222_3.jpg
Audio query: 00220.wav
Most relevant image for audio query: 00220.jpg
Audio query: 00222_3.wav
Most relevant image for audio query: 00220.jpg
```

```
count_crt
```

```
➞ 2
```

```
for index, audio_query_emb in enumerate(embeddings_audio['audio']):
    print(index, audio_query_emb)
```

```
➞ 0 tensor([ 0.0478, -0.2131,  0.2315, ...,  0.5279, -0.7604, -0.0132],
      device='cuda:0')
1 tensor([ 0.2381, -0.0772, -0.2055, ..., -0.4301,  0.3229, -0.4043],
      device='cuda:0')
2 tensor([ 0.0859, -0.0992,  0.6365, ..., -0.1347, -0.8802,  0.0288],
      device='cuda:0')
3 tensor([-0.3175, -0.6430, -0.4416, ...,  0.0071,  0.6742,  0.1453],
      device='cuda:0')
4 tensor([ 0.5646, -0.1296, -0.1822, ..., -0.1235,  0.5655,  0.2960],
      device='cuda:0')
5 tensor([-0.2691, -0.3681, -0.1953, ..., -0.4585,  0.3020,  0.2901],
      device='cuda:0')
```

```
vision_embeddings.shape
```

```
➞ torch.Size([6, 1024])
```

```
vision_embeddings
```

```
↳ tensor([[ 0.0016, -0.0102,  0.0057, ..., -0.0128,  0.0134,  0.0019],  
          [ 0.0013, -0.0126,  0.0278, ...,  0.0009,  0.0452, -0.0210],  
          [ 0.0209,  0.0147,  0.0075, ..., -0.0229,  0.0221, -0.0286],  
          [-0.0012,  0.0201, -0.0152, ...,  0.0071, -0.0067, -0.0100],  
          [-0.0206, -0.0143,  0.0086, ...,  0.0110,  0.0229, -0.0065],  
          [-0.0201,  0.0095,  0.0293, ...,  0.0029, -0.0007, -0.0067]],  
          device='cuda:0')
```