

# Audio-Visual Generalized Zero-Shot Learning using Pre-Trained Large Multi-Modal Models

## Supplementary Material

### A. Additional Details about Textual Feature Extraction

#### A.1. CLIP Feature Extraction

To boost zero-shot classification performance, [1] calculate normalized CLIP text embeddings for an ensemble of text prompts to retrieve final textual embeddings. Then the mean is taken and the result is normalized again. Normalizing the individual CLIP text representations is necessary in order to obtain a meaningful averaged vector. The second normalization facilitates the calculation of cosine similarity scores. Note that image embeddings are normalized as well.

For UCF-GZSL<sup>cls</sup> and ActivityNet-GZSL<sup>cls</sup>, we use an ensemble of 48 different prompt templates for each class. UCF-GZSL<sup>cls</sup> and ActivityNet-GZSL<sup>cls</sup> have a similar context since both are action recognition datasets. Hence, we use the same text prompts for these two datasets (see listing 1). These templates are taken from the CLIP repository<sup>1</sup>.

```
1 CLIP_prompt_templates = [  
2     'a photo of a person {}.',  
3     'a video of a person {}.',  
4     'a example of a person {}.',  
5     'a demonstration of a person {}.',  
6     'a photo of the person {}.',  
7     'a video of the person {}.',  
8     'a example of the person {}.',  
9     'a demonstration of the person {}.',  
10    'a photo of a person using {}.',  
11    'a video of a person using {}.',  
12    'a example of a person using {}.',  
13    'a demonstration of a person using {}.',  
14    'a photo of the person using {}.',  
15    'a video of the person using {}.',  
16    'a example of the person using {}.',  
17    'a demonstration of the person using {}.',  
18    'a photo of a person doing {}.',  
19    'a video of a person doing {}.',  
20    'a example of a person doing {}.',  
21    'a demonstration of a person doing {}.',  
22    'a photo of the person doing {}.',  
23    'a video of the person doing {}.',  
24    'a example of the person doing {}.',  
25    'a demonstration of the person doing {}.',  
26    'a photo of a person during {}.',  
27    'a video of a person during {}.',  
28    'a example of a person during {}.',  
29    'a demonstration of a person during {}.',  
30    'a photo of the person during {}.',  
]
```

<sup>1</sup><https://github.com/openai/CLIP/blob/main/data/prompts.md#ucf101>

```
31     'a video of the person during {}.',  
32     'a example of the person during {}.',  
33     'a demonstration of the person during {}.',  
34     'a photo of a person performing {}.',  
35     'a video of a person performing {}.',  
36     'a example of a person performing {}.',  
37     'a demonstration of a person performing {}.',  
38     'a photo of the person performing {}.',  
39     'a video of the person performing {}.',  
40     'a example of the person performing {}.',  
41     'a demonstration of the person performing {}.',  
42     'a photo of a person practicing {}.',  
43     'a video of a person practicing {}.',  
44     'a example of a person practicing {}.',  
45     'a demonstration of a person practicing {}.',  
46     'a photo of the person practicing {}.',  
47     'a video of the person practicing {}.',  
48     'a example of the person practicing {}.',  
49     'a demonstration of the person practicing {}.',  
50 ]
```

Listing 1. Text prompt templates that were used to create CLIP label embeddings for UCF-GZSL<sup>cls</sup> and ActivityNet-GZSL<sup>cls</sup>.

VGGSound-GZSL<sup>cls</sup> contains videos of a variety of categories and hence more general prompts are required. The prompts that we used to create CLIP text embeddings for VGGSound-GZSL<sup>cls</sup> can be seen in listing 2.

```
1 VGGSound_CLIP_prompt_templates = [  
2     'a photo of {}.',  
3     'a video of {}.',  
4     'a example of {}.',  
5     'a demonstration of {}.',  
6     'a photo of the person {}.',  
7     'a video of the {}.',  
8     'a example of the {}.',  
9     'a demonstration of the {}.',  
10 ]
```

Listing 2. Text prompt templates that were used to create CLIP text embeddings for VGGSound-GZSL<sup>cls</sup>.

#### A.2. CLAP Feature Extraction

We use the same procedure as in A.1 to extract textual CLAP embeddings. For UCF-GZSL<sup>cls</sup> and ActivityNet-GZSL<sup>cls</sup> we use the prompts as in listing 3. For VGGSound-GZSL<sup>cls</sup>, we use the prompts given in listing 4.

```
1 CLAP_prompt_templates = [  
     'a person {} can be heard.',  
]
```

```

3 'a example of a person {} can be heard.',
4 'a demonstration of a person {} can be heard.'
5 ',
6 'the person {} can be heard.',
7 'a example of the person {} can be heard.',
8 'a demonstration of the person {} can be
9 heard.',
10 'a person using {} can be heard.',
11 'a example of a person using {} can be heard.'
12 ',
13 'a demonstration of a person using {} can be
14 heard.',
15 'a example of the person using {} can be
16 heard.',
17 'a demonstration of the person using {} can
18 be heard.',
19 'a person doing {} can be heard.',
20 'a example of a person doing {} can be heard.'
21 ',
22 'a demonstration of a person doing {} can be
23 heard.',
24 'a example of the person doing {} can be
25 heard.',
26 'a demonstration of the person doing {} can
27 be heard.',
28 'a example of a person during {} can be heard.'
29 ',
30 'a demonstration of a person during {} can be
31 heard.',
32 'a example of the person during {} can be
33 heard.',
34 'a demonstration of the person during {} can
35 be heard.',
36 'a person performing {} can be heard.',
37 'a example of a person performing {} can be
38 heard.',
39 'a demonstration of a person performing {}
40 can be heard.',
41 'a example of the person performing {} can be
42 heard.',
43 'a demonstration of the person performing {}
44 can be heard.',
45 'a person practicing {} can be heard.',
46 'a example of a person practicing {} can be
47 heard.',
48 'a demonstration of a person practicing {}
49 can be heard.',
50 'a example of the person practicing {} can be
51 heard.',
52 'a demonstration of the person practicing {}
53 can be heard.'
54 ]

```

Listing 3. Text prompt templates that were used to create CLAP label embeddings for UCF-GZSL<sup>cls</sup> and ActivityNet-GZSL<sup>cls</sup>.

```

1 VGGSound_CLAP_prompt_templates = [
2 'a {} can be heard.',
3 'a example of a {} can be heard.',
4 'a demonstration of a {} can be heard.',
5 'the {} can be heard.',
6 'a example of the {} can be heard.',
7 'a demonstration of the {} can be heard.',
8 '{} can be heard.',
9 'a example of {} can be heard.',
10 'a demonstration of {} can be heard.'

```

11 ]

Listing 4. Text prompt templates that were used to create CLAP text embeddings for VGGSound-GZSL<sup>cls</sup>.

## References

- [1] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *ICML*, 2021. 1