# ONE-PEACE: Exploring one general Representation Model toward unlimited modalities

**Peng Wang**[1*]**, Shijie Wang**[1,2*]**, Junyang Lin**[1]**, Shuai Bai**[1]

**Xiaohuan Zhou**[1]**, Jingren Zhou**[1]**, Xinggang Wang**[2]**, Chang Zhou**[1†]

[1]DAMO Academy, Alibaba Group   [2]Huazhong University of Science and Technology

## ABSTRACT

In this work, we explore a scalable way for building a general representation model toward unlimited modalities. We release ONE-PEACE, a highly extensible model with 4B parameters that can seamlessly align and integrate representations across vision, audio, and language modalities. The architecture of ONE-PEACE comprises modality adapters, shared self-attention layers, and modality FFNs. This design allows for the easy extension of new modalities by adding adapters and FFNs, while also enabling multi-modal fusion through self-attention layers. To pretrain ONE-PEACE, we develop two modality-agnostic pretraining tasks, cross-modal aligning contrast and intra-modal denoising contrast, which align the semantic space of different modalities and capture fine-grained details within modalities concurrently. With the scaling-friendly architecture and pretraining tasks, ONE-PEACE has the potential to expand to unlimited modalities. Without using any vision or language pretrained model for initialization, ONE-PEACE achieves leading results on a wide range of uni-modal and multi-modal tasks, including image classification (ImageNet), semantic segmentation (ADE20K), audio-text retrieval (AudioCaps, Clotho), audio classification (ESC-50, FSD50K, VGGSound), audio question answering (AVQA), image-text retrieval (MSCOCO, Flickr30K), and visual grounding (RefCOCO/+/g).

Code is available at https://github.com/OFA-Sys/ONE-PEACE

## 1 Introduction

Representation models have received considerable attention in computer vision [1, 2, 3, 4, 5, 6, 7, 8, 9], speech processing [10, 11, 12, 13], natural language processing [14, 15, 16, 17, 18], etc. Learning from large amounts of data, representation models demonstrate strong generalization ability in a wide range of downstream tasks. Furthermore, the explosive growth of large-scale language models (LLMs) has sparked an escalating appetite for representation models. Until recently, representation models have shown their bedrock role to unleash LLMs to understand, perceive, and interact with other modalities (e.g., vision) [19, 20, 21, 22, 23, 24, 25].

Due to the distinct characteristics of different modalities, previous research mainly focuses on building uni-modal representation models with individual architectures and pretraining tasks. Despite achieving excellent results, uni-modal representation models face difficulties in effectively utilizing multi-modal data such as image-text pairs and audio-text pairs, which makes them challenging to extend to multi-modal tasks. With the development of unified architectures [26, 27, 28, 29] and efficient pretraining tasks [15, 1, 2, 30], recent works have achieved promising results in vision-language learning [31, 32, 33, 34, 35, 36] and audio-language learning [37, 38, 39, 40]. Nevertheless, there is still rare research on developing general models that can be applied to vision, audio, and language modalities. [34] utilize the Multiway Transformer to process both image and text modalities with a unified masked prediction task for pretraining. The masked prediction task requires a pretrained CLIP [30] model to discretize image data, which limits the scalability to other modalities such as audio. [41] proposes a general pretraining method that can be applied to

---

vision, audio, and language modalities without the need for third-party models (e.g., CLIP), but it doesn't extend the method to multi-modal data.

In this paper, we explore a scalable way to build a general representation model toward unlimited modalities. We advocate that a general representation model should meet the following conditions: 1. The model architecture must be flexible enough to accommodate various modalities and support multi-modal interaction. 2. Pretraining tasks should not only extract information within each modality but also ensure alignment across modalities. 3. Pretraining tasks should be general and straightforward, allowing them to be applied to different modalities.

Driven by these motivations, we propose ONE-PEACE, a model with 4B parameters that can seamlessly align and integrate representations across vision, audio, and language modalities. The architecture of ONE-PEACE consists of multiple modality adapters and a modality fusion encoder. Each modality is equipped with an adapter for converting the raw inputs into feature sequences. The modality fusion encoder operates on feature sequences with Transformer architecture. Each Transformer block contains a shared self-attention layer and multiple modality Feed Forward Networks (FFNs). The self-attention layer enables interaction between the multi-modal features through the attention mechanism, while the modality FFNs facilitate information extraction within modalities. With the clear division of labor in this architecture, extending new modalities only requires the injection of adapters and FFNs.

To pretrain ONE-PEACE, we design two modality-agnostic pretraining tasks. The first one is cross-modal contrastive learning, it contains both vision-language contrastive learning and audio-language contrastive learning, which effectively align the semantic spaces of vision, audio, and language modalities. The second one is intra-modal denoising contrastive learning, it can be viewed as a combination of masked prediction [15, 1, 2, 10] and contrastive learning [30, 7, 4], where we perform contrastive loss between the fine-grained masked features and visible features, such as image patches, language tokens, or audio waveform features. These tasks collaborate to enhance the model's fine-tuning performance while also maintaining cross-modal retrieval capability. Furthermore, they are universal for all modalities, which obviates the need for modality-specific designs. With the scaling-friendly model architecture and pretraining tasks, ONE-PEACE has the potential to expand to unlimited modalities.

We conduct comprehensive experiments on different tasks across various modalities, including vision, audio, vision-language, and audio-language tasks. Without using any vision or language pretrained model for initialization, ONE-PEACE achieves leading results in both uni-modal and multi-modal tasks, including image classification (89.8% accuracy on ImageNet w/o privately labeled data), semantic segmentation (63.0% mIoU on ADE20K), audio-text retrieval (outperforming previous SOTAs on AudioCaps and Clotho by a large margin), audio classification (91.8% zero-shot accuracy on ESC-50, 69.7% accuracy on FSD50K, 59.6% accuracy on VGGSound w/o visual information), audio question answering (86.2% accuracy on AVQA w/o visual information), image-text retrieval (84.1% I2T R@1 on MSCOCO and 97.6% I2T R@1 on Flickr30K w/o intermediate finetuning and ranking), and visual grounding (89.26%/83.23%/89.27% scores on RefCOCO/+/g test sets).

## 2   Related Work

**Vision-Language Pretraining.**   Recent years have witnessed the rapid development of vision-language pretraining. Early approaches [42, 43, 44, 45, 46, 47, 48, 49, 50, 51] relied heavily on region features extracted by object detectors, which is resource&time-consuming. With the increasing popularity of Vision Transformer [27], numerous works use Transformer to jointly learn vision-language data and demonstrate superior performance in downstream tasks [52, 53, 54, 55, 56, 57, 35, 21, 58]. To facilitate alignment between vision and language modalities, researchers propose various efficient pretraining tasks. Among them, contrastive learning is one of the most representative methods that has been widely adopted in a lot of works [30, 59, 60, 61, 7, 4, 53, 62]. There also emerge some works explore the unified frameworks to handle vision-language tasks [31, 63, 64, 65, 66, 33, 34, 36, 67, 68]. [31] adopts an encoder-decoder model to transform all vision-language tasks into generation tasks. [33] uses contrastive learning and text generation as the pretraining objectives, thus can be applied to image-text retrieval and vision-language generation tasks. [34] employs the Multiway Transformer to process vision-language data, and discretizes images into image tokens through CLIP [30] for joint learning with text tokens.

**Audio-Language Pretraining.**   There is currently a significant amount of research being conducted in audio-language pretraining. One category of these works focuses on speech-text joint pretraining. For instance, some studies propose to train a unified encoder for speech and text, which utilizes a large amount of unlabeled speech and text with masked prediction tasks and paired speech-text data to learn alignment [69, 70, 71, 72]. There are also some works proposed to jointly pretrain speech and text under the encoder-decoder framework [40, 73, 74, 75], which can be well applied to generation tasks, such as speech recognition and synthesis. Another category introduces cross-modal contrastive learning to audio-language pretraining [37, 39, 38, 76]. [39] uses CNN14 [77] and BERT [15] to extract audio and text
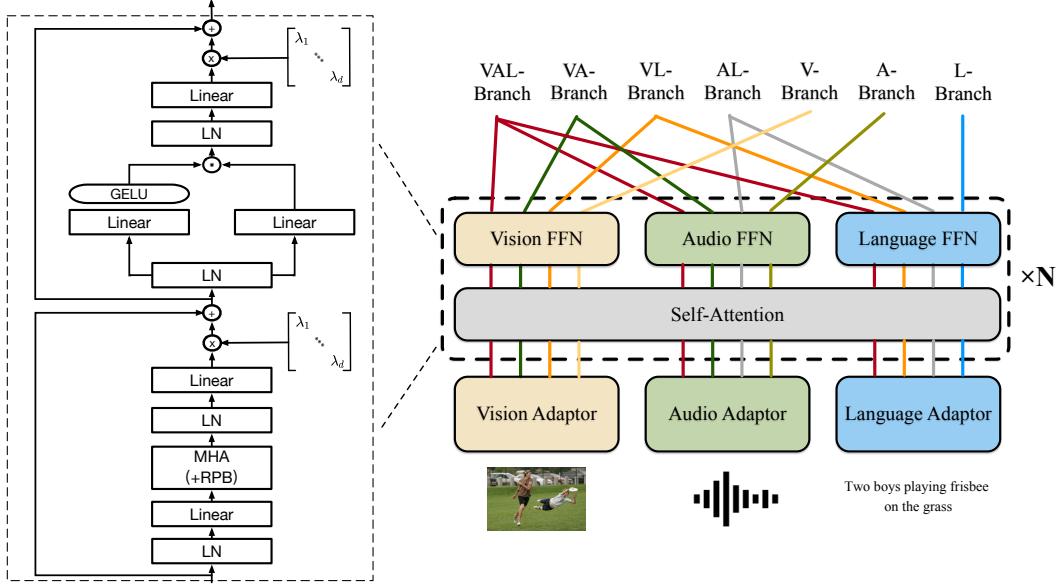
Figure 1: **The architecture of ONE-PEACE**. It consists of three modality adapters and a modality fusion encoder. ONE-PEACE can be disassembled into different branches to handle different tasks. For example, the vision adapter, self-attention layers, and vision FFNs can be combined into V-Branch to handle vision tasks.

information respectively, and conducts contrastive learning on environmental sound data. [76] further introduces more environmental sound data and trained with HTSAT [78] and RoBERTa [16]. It achieves state-of-the-art results in audio downstream tasks such as audio classification and audio-text retrieval.

**Vision-Audio-Language Pretraining.** Recently researchers have begun exploring joint learning of vision, audio, and language modalities. [79] employs a unified Transformer model to jointly learn video, text, and audio through cross-modal contrastive learning. [80, 81] utilize external models (e.g., VQ-VAE [82] and AV-HuBERT [83]) to discretize the video and audio data, and train the models with masked prediction objectives. [41] proposes a general self-supervised learning method that does not rely on external models. It successfully applies the method to vision, language, and audio modalities, but has not extended to multi-modal data.

Compared to previous works, ONE-PEACE has a flexible architecture that is compatible with multiple modalities. Furthermore, the pretraining tasks of ONE-PEACE are universally applicable without external models Therefore, ONE-PEACE can be easily extended to various modalities.

## 3 Method

### 3.1 Architecture

The model architecture of ONE-PEACE consists of three modality adapters and a modality fusion encoder. The overall architecture is shown in Figure 1.

**Modality Adapters.** We design modality adapters to convert different raw signals into unified features. Note that these adapters do not interact with each other, which affords us the flexibility to choose appropriate networks for them, such as Transformers [26, 27, 84], CNNs [85, 86], RNNs [87, 88], etc. We design three lightweight modality adapters for ONE-PEACE:

- **Vision Adapter (V-Adapter).** Given an image, we use a hierarchical MLP (hMLP) stem [89] to patchify the image by gradually increasing the patch size to $16 \times 16$. There is no interaction between different patches. Then the image patches are flattened into a sequence and prepended with a vision class embedding. By adding the absolute positional embeddings into the image embeddings, the image representation is $E^V = \langle e_{cls}^V, e_1^V, e_2^V, ..., e_M^V \rangle$, where $M$ denotes the total number of image patches.

3

- **Audio Adapter (A-Adapter).** Given an audio, we set the sample rate to 16kHz and normalize the raw audio waveform to zero mean and unit variance. Then the normalized waveform is processed by a convolutional feature extractor [11] to get the audio embeddings. Instead of using the absolute positional embeddings, we use a convolution layer to extract relative position information and add it to the audio embeddings [90]. With a prepended audio class embedding, we obtain the audio representation $E^A = \langle e^A_{cls}, e^A_1, e^A_2, ..., e^A_N \rangle$, where $N$ denotes the length of the audio representation.
- **Language Adapter (L-Adapter).** Given a text, we first apply byte-pair encoding (BPE) [91] to transform it to a subword sequence. Two special tokens [CLS] and [EOS] are inserted at the beginning and end of the sentence to indicate its start and end. Then an embedding layer is used to embed the subword sequence to the text embeddings. After summing the text embeddings with absolute positional embeddings, we obtain the text representation $E^L = \langle e^L_{cls}, e^L_1, e^L_2, ..., e^L_K, e^L_{eos} \rangle$, where $K$ denotes the text sequence length.

**Modality Fusion Encoder.** Following previous works [63, 31, 33, 34, 92], the modality fusion encoder is based on the Transformer architecture [26]. We set up a shared self-attention layer and three modality feed-forward networks (FFNs) in each Transformer block. The shared self-attention layer enables the interaction between different modalities through the attention mechanism. The three modality FFNs (V-FFN, A-FFN, and L-FFN) can further extract information within their respective modalities. To stabilize training and enhance model performance, we make the following improvements:

- **Sub-LayerNorm.** We incorporate Sub-LayerNorm [93] into each Transformer block to enhance training stability. Specifically, We insert layer normalization before the input projection and output projection of each self-attention layer and FFN layer. In our preliminary experiments, we find that Sub-LayerNorm can achieve better performance compared to the Pre-LayerNorm [94].
- **GeGLU Activation Function.** To further improve performance, we replace the activation function in FFN with GeGLU [95] activation function. The intermediate dimension of FFN is set to $4$ times the embedding dimension, which is consistent with the practice of PaLM [96].
- **Relative Position Bias (RPB).** For positional information, we introduce 1D relative position bias [97] for text and audio, and 2D relative position bias for image [98]. At the pretraining stage, the relative position bias of different self-attention layers is shared. At the fine-tuning stage, we decouple the relative position bias of each self-attention layer and let them inherit the weights of the pretrained relative bias.
- **LayerScale.** We use LayerScale [99] to dynamically adjust the output of each residual block. Specifically, before adding to the residual, we multiply the output of each layer (e.g., self-attention layer and FFN) by a learnable diagonal matrix, whose values will be initialized to $1e-6$. In our preliminary experiments, LayerScale is beneficial for stabilizing training and improving performance.

This "sharing-separated" architecture enables ONE-PEACE to disassemble into different branches that handle tasks for various modalities. For example, the vision adapter, self-attention layer, and vision FFNs can be combined into the vision branch (V-Branch) to process vision tasks. Similarly, we named other branches as audio branch (A-Branch), language branch (L-Branch), vision-audio branch (VA-Branch), vision-language branch (VL-Branch), audio-language branch (AL-Branch), and vision-audio-language branch (VAL-Branch).

### 3.2 Tasks

The pretraining tasks of ONE-PEACE include cross-modal contrastive learning and intra-modal denoising contrastive learning. Cross-modal contrastive learning endows the model with cross-modal retrieval capability, while intra-modal denoising contrastive learning enables the model to achieve superior fine-tuning performance in downstream tasks. An illustration of the pretraining tasks is shown in Figure 2.

**Cross-Modal Contrastive Learning.** Cross-modal contrastive learning is a widely-used pretraining task that effectively aligns the semantic spaces of different modalities. The key idea of this method is to maximize the similarity of related sample pairs across different modalities while minimizing the similarity of unrelated sample pairs. Given a sample pair $(S^1, S^2)$ of arbitrary modalities (e.g., image-text pair or audio-text pair), we extract their features using the corresponding branches of ONE-PEACE. The outputs of the special tokens (e.g., vision class token or language class token) are regarded as global representations. Followed by a linear projection and normalization, we obtain the final representations $s^1$ and $s^2$. The loss function is shown below:

$$\mathcal{L}_{CL} = -\frac{1}{2N} \sum_{i=1}^{N} (\log \frac{\exp(s^1_i s^2_i/\sigma)}{\sum_{j=1}^{N} \exp(s^1_i s^2_j/\sigma)} + \log \frac{\exp(s^1_i s^2_i/\sigma)}{\sum_{j=1}^{N} \exp(s^1_j s^2_i/\sigma)}), \tag{1}$$
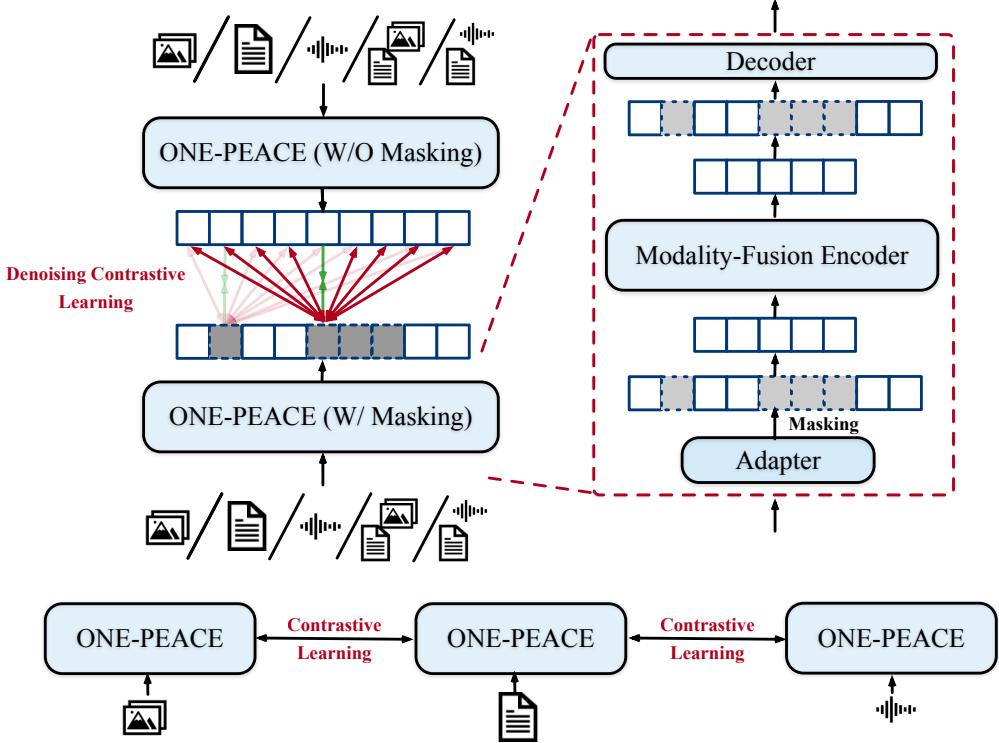
Figure 2: **The pretraining tasks of ONE-PEACE.** Intra-modal denoising contrastive learning encourages the features of the masked units (e.g., image patches or text tokens) close to the positive units (indicated by the green lines) and get away from the negative units (indicated by the red lines). Note that we compute the cross-modal contrastive loss by gathering negative features from all GPU devices, while the denoising contrastive loss is computed on the local batch.

where $N$ is the batch size, $i, j$ are indexes within the batch, and $\sigma$ is a learnable temperature parameter (initialized to 0.07). Following previous works [30, 59], the cross-modal contrastive loss is computed by gathering negative features from all GPU devices. We apply cross-modal contrastive learning to image-text pairs and audio-text pairs, denoted by $\mathcal{L}_{CL-VL}$ and $\mathcal{L}_{CL-AL}$ respectively.

**Intra-Modal Denoising Contrastive Learning.** Cross-modal contrastive learning mainly focuses on aligning features between different modalities. However, it lacks emphasis on the learning of fine-grained details within modalities, leading to suboptimal performance in downstream tasks [100]. To address this issue, we further introduce intra-modal denoising contrastive learning to train ONE-PEACE[2]. Intra-modal denoising contrastive learning can be viewed as a combination of masked prediction and contrastive learning, where we perform contrastive loss between the fine-grained masked features and visible features, such as image patches, text tokens, or audio waveform features.

Given a sample of arbitrary modalities, we first encode it into an embedding sequence through the corresponding modality adapter. Then, we randomly mask some units (e.g., text tokens or image patches) within the sequence. Following [2], we only input the unmasked units to the modality fusion encoder to reduce computation costs and save memory. The encoded unmasked features are concatenated with the learnable mask tokens and fed to a lightweight Transformer decoder, which generates the masked features. We also use the ONE-PEACE model to encode the raw input sample into target features without masking. Finally, we perform the contrastive loss between the masked features and target features, the loss function is shown below:

$$\mathcal{L}_{DCL} = -\frac{1}{N\hat{N}} \sum_{i=1}^{N} \sum_{j=1}^{\hat{N}} \log \frac{\exp(\hat{\boldsymbol{h}}_{ij} \cdot \mathrm{sg}(\boldsymbol{h}_{ij})/\tau)}{\sum_{m=1}^{N} \sum_{n=1}^{N} \exp(\hat{\boldsymbol{h}}_{ij} \cdot \mathrm{sg}(\boldsymbol{h}_{mn})/\tau)}, \tag{2}$$

---

[2]Intra-modal denoising contrastive learning is similar to [101], but extends to more modalities.

5

| #Layers | Hidden Size | Intermediate Size | Attention Size | #Parameters | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | | V-Adapter | A-Adapter | L-Adapter | V-FFN | A-FFN | L-FFN | Shared Attention | Total |
| 40 | 1536 | 6144 | 24 | 3.4M | 19M | 78M | 1.15B | 1.15B | 1.15B | 378M | 4B |

Table 1: **Detailed hyperparameters of ONE-PEACE model configuration.**

Where $\hat{\boldsymbol{h}}_{ij}$ is the representation of the masked unit, $\boldsymbol{h}_{ij}$ is the representation of the target unit, sg($\cdot$) is the stop gradient operation. $\hat{N}$ is the number of masked units within a sample, $N$ is the number of whole units within a sample. $\tau$ is a constant temperature value, we set it to $0.4$. This loss not only encourages the masked units close to the positive units but also gets away from the negative units. As a result, each unit acquires a unique semantic meaning, which makes the model better transfer to downstream tasks [100].

We apply intra-modal denoising contrastive learning to 5 types of data: image, audio, text, image-text pairs, and audio-text pairs. For image, we randomly mask $75\%$ patches, and the loss function used for this type of data is denoted by $\mathcal{L}_{DCL-V}$. For audio, we sample $p = 0.11$ of all time-steps to be starting indices and mask the subsequent 5 time-steps, and the loss function is denoted by $\mathcal{L}_{DCL-A}$. For text, we randomly mask $15\%$ tokens of the text sequence, and the loss function is denoted by $\mathcal{L}_{DCL-L}$. For image-text pairs, we randomly mask $68.75\%$ patches of the image and $40\%$ tokens of the text. The unmasked patches and tokens are concatenated together and encoded as masked features. The original image patches and text tokens are also concatenated together and encoded as target features. We then perform contrastive loss on the image patches and text tokens respectively, the average of these two losses is denoted by $\mathcal{L}_{DCL-VL}$. For audio-text pairs, we randomly mask $45\%$ time-steps of the audio waveform and $40\%$ tokens of the text. The loss is similar to the above one, we denote it by $\mathcal{L}_{DCL-AL}$.

## 3.3 Training

The overall pretraining process of ONE-PEACE is divided into two stages: vision-language pretraining and audio-language pretraining. At the vision-language pretraining stage, the model trains on image-text pairs and only updates parameters that are relevant to vision and language modalities. For each image-text pair, we not only utilize them to calculate $\mathcal{L}_{CL-VL}$ and $\mathcal{L}_{DCL-VL}$, but also separately using the image and text to calculate $\mathcal{L}_{DCL-V}$ and $\mathcal{L}_{DCL-L}$ respectively. The loss function at this stage is shown below:

$$\mathcal{L}_{VL} = \mathcal{L}_{CL-VL} + 1.0 * \mathcal{L}_{DCL-V} + 0.5 * \mathcal{L}_{DCL-L} + 1.0 * \mathcal{L}_{DCL-VL} \tag{3}$$

At the audio-language pretraining stage, the model trains on audio-text pairs, and we only update A-Adapter, A-FFNs, and other audio-related parameters. The remaining parameters including self-attention layers are totally frozen. Despite not training on image-audio pairs, the semantic space between vision and audio is still aligned by using language as the anchor. The loss function at the audio-language pretraining stage is shown below:

$$\mathcal{L}_{AL} = \mathcal{L}_{CL-AL} + 1.0 * \mathcal{L}_{DCL-A} + 1.0 * \mathcal{L}_{DCL-AL} \tag{4}$$

## 4 Pretraining Details

**Pretraining Datasets.** The pretraining datasets of ONE-PEACE are divided into two parts: image-text pairs and audio-text pairs. For image-text pairs, we use LAION-2B [102], a dataset obtained by web crawling. For audio-text pairs, we collect a large amount of open-source environmental sound datasets. To ensure reproducibility, all pretraining datasets are publicly available. We provide more details about the pretraining datasets in Appendix A.1.

**Pretraining Settings.** ONE-PEACE is a giant-size model with 4B parameters. We list the detailed hyper-parameters in Table 1. During pretraining, we introduce a lightweight Transformer decoder to recover the masked units from the visible units. The decoder is similar to the modality-fusion encoder, each block of it also consists of a shared self-attention layer and three modality FFNs. It has 2 layers with 768 hidden size, 2048 intermediate size, and 12 attention heads. The model weights of ONE-PEACE are randomly initialized at the beginning, except for the audio feature extractor of A-adapter, for which we use the weights of WavLM's feature extractor [12] for initialization. We find that incorporating WavLM's feature extractor significantly improves the model performance. More details about the pretraining settings are provided in Appendix A.2.

| Method | Enc. #Params | Patch Size | Image size | Top-1 acc |
|---|---|---|---|---|
| FD-SwinV2-G [103] | 3.0B | $16 \times 16$ | $336^2$ | 89.4 |
| InternImage [104] | 1.08B | $16 \times 16$ | $640^2$ | 89.2 |
| BEiT-3 [34] | 1.01B | $14 \times 14$ | $336^2$ | 89.6 |
| EVA [105] | 1.01B | $14 \times 14$ | $560^2$ | 89.7 |
| ONE-PEACE | 1.52B | $16 \times 16$ | $384^2$ | 89.6 |
| ONE-PEACE | 1.52B | $16 \times 16$ | $512^2$ | **89.8** |
| *methods using extra privately collected data:* | | | | |
| RevCol-H [106] | 2.16B | - | $640^2$ | 90.0 |
| ViT-G [107] | 1.84B | $14 \times 14$ | $518^2$ | 90.5 |
| Model Soups [107] | 1.84B | $14 \times 14$ | $500^2$ | 90.9 |
| CoCa [33] | 1.01B | $18 \times 18$ | $576^2$ | 91.0 |

Table 2: **System-level comparisons of image classification with the leading results on ImageNet-1k.** RevCol [106] is pretrained on a privately collected 168-million-image dataset. ViT-G [107] and Model Soups [107] are pretrained on JFT-3B [108] with supervision. CoCa [33] uses JFT-3B [108] and ALIGN [59] datasets for pretraining. ONE-PEACE achieves state-of-the-art results using the publicly available dataset with less token length.

| Method | Enc. #Params | Crop Size | mIoU$^{ss}$ | mIoU$^{ms}$ |
|---|---|---|---|---|
| RevCol-H [106] | 2.16B | $640^2$ | 60.4 | 61.0 |
| FD-SwinV2-G [103] | 3.00B | $896^2$ | - | 61.4 |
| ViT-Adapter [112] | 571M | $896^2$ | 61.2 | 61.5 |
| EVA [105] | 1.01B | $896^2$ | 61.5 | 62.3 |
| BEiT-3 [34] | 1.01B | $896^2$ | 62.0 | 62.8 |
| InternImage [104] | 1.08B | $896^2$ | **62.5** | 62.9 |
| ONE-PEACE | 1.52B | $896^2$ | 62.0 | **63.0** |

Table 3: **System-level comparisons of semantic segmentation with leading results on ADE20k.** mIoU$^{ss}$ means single-scale inference result while mIoU$^{ms}$ means multi-scale.

**Training Acceleration.** We introduce several model acceleration and memory optimization techniques to accelerate the training. Firstly, we use the memory-efficient attention technique [109, 110] implemented in the xformers library[3] to improve training speed. Secondly, we use the gradient checkpointing technique [111] to save memory, which allows us to train the model with a larger batch size. Furthermore, we replace the layer normalization with Fused LayerNorm implemented in the Flash Attention library[4], and leverage nvFuser[5] to fuse the operations of dropout, LayerScale, stochastic depth, and residual summing, which can bring additional speed improvements. To improve the training speed and prevent gradient overflow issues, we adopt Bfloat16 precision to train ONE-PEACE.

## 5 Experiments

### 5.1 Results on Vision Tasks

We transfer ONE-PEACE to various mainstream vision benchmarks, including image classification, semantic segmentation, object detection, instance segmentation and video action recognition. We provide the implementation details in Appendix B.1.

**Image Classification.** In our experiments, we assess the image classification transfer performance of ONE-PEACE using the ImageNet-1K [121] dataset, encompassing 1.28 million training images and 50,000 validation images distributed across 1,000 distinct categories. We also use intermediate fine-tuning on ImageNet-21k [122]. As demonstrated in Table 2, ONE-PEACE obtains **89.8** top-1 accuracy on ImageNet with less token length $(image\_size/patch\_size)^2$.

---

[3] https://github.com/facebookresearch/xformers
[4] https://github.com/HazyResearch/flash-attention
[5] https://pytorch.org/blog/introducing-nvfuser-a-deep-learning-compiler-for-pytorch

| Method | Detector | #Params | Image Size | $AP^{box}$ | $AP^{mask}$ |
|---|---|---|---|---|---|
| ViT-Adapter [112] | HTC++ | 401M | [400-1400, 1600] | 58.8 | 51.1 |
| ViTDet [113] | Cascade | 692M | $1280^2$ | 60.4 | 52.0 |
| RevCol-H [106] | HTC++ | 2.41B | [400-1400, 1600] | **61.1** | **53.0** |
| ONE-PEACE | Cascade | 1.59B | $1280^2$ | 60.4 | 52.9 |

Table 4: **System-level comparisons of object detection and instance segmentation on MSCOCO.** The reported results are obtained by directly fine-tuning the models on MSCOCO, without intermediate fine-tuning on Objects365.

| Method | Backbone | Input Size | Top-1 | Top-5 |
|---|---|---|---|---|
| VATT [79] | ViT-L | $32 \times 320^2$ | 82.1 | 95.5 |
| ViViT [114] | ViT-H | $32 \times 224^2$ | 84.9 | 95.8 |
| Florance [61] | Co-Swin-H | $N/A \times 384^2$ | 86.5 | 97.3 |
| SwinV2 [115] | Swin-G | $N/A \times 384^2$ | 86.8 | - |
| MAE-ST [116] | ViT-H | $16 \times 224^2$ | 86.8 | 97.2 |
| VideoMAE [117] | ViT-H | $32 \times 320^2$ | 87.4 | 97.6 |
| VideoMAE V2 [118] | ViT-H | $32 \times 320^2$ | 87.4 | 97.6 |
| MaskFeat [119] | MViTv2-L | $40 \times 352^2$ | 87.0 | 97.4 |
| CoCa (frozen) [33] | ViT-g | $16 \times 576^2$ | 88.0 | - |
| ViT-22B (frozen) [120] | ViT-22B | $128 \times 224^2$ | 88.0 | - |
| ONE-PEACE (frozen) | ViT-g | $16 \times 256^2$ | 88.0 | **97.8** |
| ONE-PEACE (frozen) | ViT-g | $32 \times 256^2$ | **88.1** | **97.8** |
| *methods using intermediate fine-tuning:* | | | | |
| EVA [105] | ViT-g | $16 \times 224^2$ | 89.7 | - |
| VideoMAE V2 [118] | ViT-g | $64 \times 266^2$ | 90.0 | 98 |

Table 5: **System-level comparisons of video action recognition with leading results on Kinetics-400.** Frozen means do not update pre-trained model parameters. EVA [105] and VideoMAE V2 [118] are intermediate fine-tuned on the merged Kinetics dataset (K400, K600, K700).

Note that FD-SwinV2-G, BEiT-3, and EVA all rely on the assistance of an external CLIP model for pretraining, while ONE-PEACE is trained from scratch without the help of external models. Even so, ONE-PEACE is able to achieve better results, which demonstrates its strong transferability.

**Semantic Segmentation.** We experiment on ADE20k [123] using ViT-Adapter [112] for task adaptation and Mask2Former [124] as the segmentation head. Following common practice, we first fine-tune the segmentation head on coco-stuff [125] then fine-tune on ADE20k. As demonstrated in Table 3, ONE-PEACE establishes a new state-of-the-art, achieving a mean Intersection over Union (mIoU) of **63.0**. This result indicates that ONE-PEACE exhibits exceptional transferring performance in the domain of dense prediction tasks.

**Object Detection and Instance Segmentation.** We perform fine-tuning experiments on the COCO 2017 [126] dataset. For the backbone, we employ the ONE-PEACE backbone and use the ViTDet [113] with Cascade Mask-RCNN architecture, which incorporates a straightforward feature pyramid and window attention for addressing object detection and instance segmentation tasks. The model is fine-tuned on the COCO dataset. Soft-NMS [127] is used during the inference stage. As illustrated in Table 4, the instance-level transfer capabilities of ONE-PEACE exhibit a performance that is on par with the current state-of-the-art methods.

**Video Action Recognition.** We benchmark ONE-PEACE on Kinetics 400 [128] dataset for video action recognition. Following AIM [129], we keep the whole model frozen and add several MLP adapters in each transformer layer. We use I3D [128] head as the classification layer. As demonstrated in Table 5, without fine-tuning the full encoder, ONE-PEACE could achieve **88.1** top-1 accuracy, even outperforming CoCa which is pre-trained on privately collected data, and ViT-22B with 14x more parameters.

| Method | AudioCaps | | | | | | Clotho | | | | | |
| | Text → Audio | | | Audio → Text | | | Text → Audio | | | Audio → Text | | |
| | R@1 | R@5 | R@10 | R@1 | R@5 | R@10 | R@1 | R@5 | R@10 | R@1 | R@5 | R@10 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| MMT [130] | 36.1 | 72.0 | 84.5 | 39.6 | 76.8 | 86.7 | 6.7 | 21.6 | 33.2 | 7.0 | 22.7 | 34.6 |
| ML-ACT [131] | 33.9 | 69.7 | 82.6 | 39.4 | 72.0 | 83.9 | 14.4 | 36.6 | 49.9 | 16.2 | 37.6 | 50.2 |
| CLAP-HTSAT [132] | 34.6 | 70.2 | 82.0 | 41.9 | 73.1 | 84.6 | 16.7 | 41.1 | 54.1 | 20.0 | 44.9 | 58.7 |
| TAP [133] | 36.1 | 72.0 | 85.2 | 41.3 | 75.5 | 86.1 | 16.2 | 39.2 | 50.8 | 17.6 | 39.6 | 51.4 |
| LAION-CLAP [76] | 35.1 | 71.5 | 83.6 | 45.8 | 80.9 | 91.6 | 18.2 | 42.5 | 54.4 | 25.7 | 51.5 | 63.4 |
| ONE-PEACE | **42.5** | **77.5** | **88.4** | **51.0** | **81.9** | **92.0** | **22.4** | **49.0** | **62.7** | **27.1** | **52.3** | **65.4** |

Table 6: **Experimental results on audio-text retrieval.** ONE-PEACE significantly outperforms baselines by a large margin.

| Method | ESC-50 | FSD50K | VGGSound (Audio Only) | AQA |
| | ZS | FT | FT | FT |
|---|---|---|---|---|
| Previous SOTA | 91.0 [76] | 65.6 [134] | 59.5 [135] | 83.5 [136] |
| Wav2CLIP [59] | 41.4 | 43.1 | 46.6 | - |
| AudioCLIP [137] | 69.4 | - | - | - |
| CLAP [39] | 82.6 | 58.6 | - | - |
| LAION-CLAP [76] | 91.0 | 46.2 | 55.1* | - |
| ONE-PEACE | **91.8** | **69.7** | **59.6** | **86.2** |

Table 7: **Experimental results on audio classification and audio question answering (AQA).** "ZS" is short for zero-shot results, "FT" is short for fine-tuning results. For the VGGSound dataset, we only use the audio data and discard the video data. *We use the official code of LAION-CLAP to reproduce the result on VGGSound.

## 5.2 Results on Audio(-Language) Tasks

We evaluate ONE-PEACE on various audio and audio-language tasks, including audio-text retrieval, audio classification, and audio question answering (AQA). The implementation details are provided in Appendix B.2.

**Audio-Text Retrieval.** Table 6 presents the performance of ONE-PEACE and baseline models in the audio-text retrieval task. As a general representation model, ONE-PEACE achieves SOTA results on both AudioCaps [138] and Clotho [139] datasets, outperforming the previous audio representation model by a large margin. On AudioCaps, ONE-PEACE achieves 21.1% improvement on R@1 in text-to-audio retrieval and 11.4% improvement on R@1 in audio-to-text retrieval. On Clotho, ONE-PEACE achieves 23.1% improvement on R@1 in text-to-audio retrieval and 5.4% on R@1 in audio-to-text retrieval.

**Audio Classification & Audio Question Answering.** Table 7 present the results of ONE-PEACE and baseline models in the audio classification and audio question answering (AQA) tasks. On ESC-50, ONE-PEACE achieves 91.8 zero-shot accuracy, outperforming LAION-CLAP by 0.8. On FSD50K, ONE-PEACE significantly outperforms the previous SOTA by 4.1. For the VGGSound dataset, which consists of both visual and audio information, we only utilized the audio information and disregarded the visual information. With this setting, ONE-PEACE achieves 59.6 score, surpassing the previous SOTA by 0.1. In the audio question answering task, ONE-PEACE outperforms the previous SOTA by 2.7. These results demonstrate the superior ability of ONE-PEACE on audio-related tasks.

## 5.3 Results on Vision-Language Tasks

We conduct experiments on various vision-language tasks, including image-text retrieval, visual grounding, visual question answering, and visual reasoning. the implementation details are provided in Appendix B.3.

**Image-Text Retrieval.** Table 8 presents the performance of ONE-PEACE and baseline models on the image-text retrieval task. Under the fine-tuning setting, ONE-PEACE achieves the best performance in both MSCOCO and Flickr30K test sets. This indicates that after combining both cross-modal contrastive learning and intra-modal denoising

| Method | COCO (5K test set) | | | | | | Flickr30K (1K test set) | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | Image → Text | | | Text → Image | | | Image → Text | | | Text → Image | | |
| | R@1 | R@5 | R@10 | R@1 | R@5 | R@10 | R@1 | R@5 | R@10 | R@1 | R@5 | R@10 |
| *Zero-shot Setting* | | | | | | | | | | | | |
| CLIP [30] | 58.4 | 81.5 | 88.1 | 37.8 | 62.4 | 72.2 | 88.0 | 98.7 | 99.4 | 68.7 | 90.6 | 95.2 |
| ALIGN [59] | 58.6 | 83.0 | 89.7 | 45.6 | 69.8 | 78.6 | 88.6 | 98.7 | 99.7 | 75.7 | 93.8 | 96.8 |
| FILIP [137] | 61.3 | 84.3 | 90.4 | 45.9 | 70.6 | 79.3 | 89.8 | 99.2 | 99.8 | 75.0 | 93.4 | 96.3 |
| Florence [61] | 64.7 | 85.9 | - | 47.2 | 71.4 | - | 90.9 | 99.1 | - | 76.7 | 93.6 | - |
| CoCa [33] | **66.3** | **86.2** | 91.8 | **51.2** | **74.2** | **82.0** | **92.5** | **99.5** | **99.9** | **80.4** | **95.7** | **97.7** |
| ONE-PEACE | 64.7 | 86.0 | **91.9** | 48.0 | 71.5 | 79.6 | 90.9 | 98.8 | 99.8 | 77.2 | 93.5 | 96.2 |
| *Fine-tuning Setting* | | | | | | | | | | | | |
| ALIGN [59] | 77.0 | 93.5 | 96.9 | 59.9 | 83.3 | 89.8 | 95.3 | 99.8 | 100.0 | 84.9 | 97.4 | 98.6 |
| FILIP [137] | 78.9 | 94.4 | 97.4 | 61.2 | 84.3 | 90.6 | 96.6 | 100.0 | 100.0 | 87.1 | 97.7 | 99.1 |
| Florence [61] | 81.8 | 95.2 | - | 63.2 | 85.7 | - | 97.2 | 99.9 | - | 87.9 | 98.1 | - |
| OmniVL [140] | 82.1 | 95.9 | 98.1 | 64.8 | 86.1 | 91.6 | 97.3 | 99.9 | **100.0** | 87.9 | 97.8 | 99.1 |
| BEiT-3 [34] | 82.7 | 96.0 | 98.2 | 65.1 | **86.6** | **92.3** | 97.5 | 99.9 | **100.0** | 89.1 | **98.6** | **99.3** |
| ONE-PEACE | **84.1** | **96.3** | **98.3** | **65.4** | 86.3 | 91.9 | **97.6** | **100.0** | **100.0** | **89.6** | 98.0 | 99.1 |

Table 8: **Experimental results on image-text retrieval.** We compare with baselines under both zero-shot and fine-tuning settings. For a fair comparison, the reported results of BEiT-3 are obtained by directly fine-tuning on downstream benchmarks without intermediate fine-tuning on pretraining data.

| Method | RefCOCO | | | RefCOCO+ | | | RefCOCOg | |
|---|---|---|---|---|---|---|---|---|
| | val | testA | testB | val | testA | testB | val-u | test-u |
| VL-T5 [50] | - | - | - | - | - | - | - | 71.3 |
| UNITER [45] | 81.41 | 87.04 | 74.17 | 75.90 | 81.45 | 66.70 | 74.86 | 75.77 |
| VILLA [141] | 82.39 | 87.48 | 74.84 | 76.17 | 81.54 | 66.84 | 76.18 | 76.71 |
| MDETR [142] | 86.75 | 89.58 | 81.41 | 79.52 | 84.09 | 70.62 | 81.64 | 80.89 |
| UNICORN [143] | 88.29 | 90.42 | 83.06 | 80.30 | 85.05 | 71.88 | 83.44 | 83.93 |
| X-VLM [144] | - | - | - | 84.51 | 89.00 | 76.91 | - | - |
| Grounding-DINO [145] | 90.56 | 93.19 | 88.24 | 82.75 | 88.95 | 75.92 | 86.13 | 87.02 |
| FIBER [56] | 90.68 | 92.59 | 87.26 | 85.74 | 90.13 | 79.38 | 87.11 | 87.32 |
| OFA [31] | 92.04 | 94.03 | 88.44 | 87.86 | 91.70 | 80.71 | 88.07 | 88.78 |
| ONE-PEACE | **92.58** | **94.18** | **89.26** | **88.77** | **92.21** | **83.23** | **89.22** | **89.27** |

Table 9: Experimental results on 3 visual grounding datasets: RefCOCO, RefCOCO+, RefCOCOg. ONE-PEACE achieves state-of-the-are results without using additional visual grounding datasets (e.g., Visual genome).

contrastive learning, ONE-PEACE can effectively transfer to downstream retrieval task. Under the zero-shot setting, ONE-PEACE can achieve better or competitive performance compared to previous dual-encoder models like CLIP and Florence. Notice that the results of ONE-PEACE are inferior to CoCa, which might be because ONE-PEACE only trained on 6.4 billion image-text pairs while CoCa trained on up to 32 billion image-text pairs.

**Visual Grounding.** To evaluate the capability of visual grounding, we conduct experiments on RefCOCO, Ref-COCO+, and RefCOCOg datasets [146, 147]. Table 9 presents the results of ONE-PEACE and baseline models. It is worth noting that previous SOTA OFA use additional visual grounding datasets for training (i.e., Visual Genome [148]). Without introducing additional visual grounding datasets, ONE-PEACE still achieves new SOTA results on the 3 datasets. We also compared the visual grounding ability of ONE-PEACE and OFA on an out-of-domain Pokémon picture.[6] As shown in Figure 3, given a specific description of a Pokémon, both ONE-PEACE and OFA can obtain the correct result. However, when we directly provide the name of the Pokémon, OFA fails to obtain the correct result while ONE-PEACE can give a correct answer.

---

[6]We use the Hugging Face spaces demo of OFA: `https://huggingface.co/spaces/OFA-Sys/OFA-Visual_Grounding`
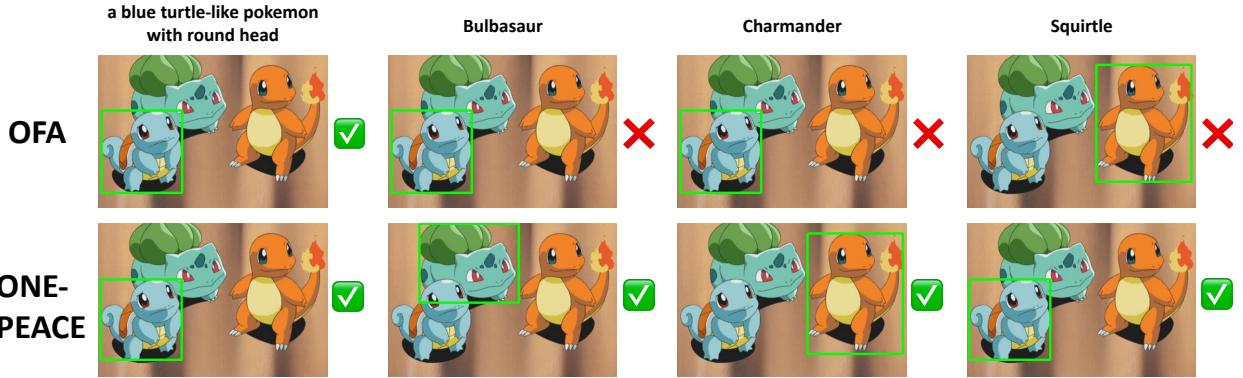
Figure 3: **Visualization of out-of-domain data in visual grounding task.** When given a specific description, both OFA and ONE-PEACE can give the correct result. However, OFA is unable to locate the correct region if we directly provide the name of the Pokémon, while ONE-PEACE can give a correct answer.

| Method | VQA | | NLVR-2 | |
|---|---|---|---|---|
| | test-dev | test-std | dev | test-P |
| ALBEF [53] | 75.8 | 76.0 | 82.55 | 83.14 |
| BLIP [35] | 78.25 | 78.32 | 82.15 | 82.24 |
| X-VLM [144] | 78.22 | 78.37 | 84.41 | 84.76 |
| SimVLM [63] | 80.0 | 80.3 | 84.5 | 85.2 |
| OFA [31] | 82.0 | 82.0 | - | - |
| Flamingo [32] | 82.0 | 82.1 | - | - |
| CoCa [33] | 82.3 | 82.3 | 86.1 | 87.0 |
| BLIP-2 [21] | 82.2 | 82.3 | - | - |
| BEiT-3 [34] | **84.2** | **84.0** | **91.5** | **92.6** |
| ONE-PEACE | 82.6 | 82.5 | 87.8 | 88.3 |

Table 10: **Results on vision-language understanding tasks.** Without initialized with language pretrained models or pretraining on pure text data, ONE-PEACE outperforms the strong baselines Flamingo and CoCa.

**Vision-Language Understanding.** Table 10 presents the results of ONE-PEACE and baselines on two popular multimodal understanding tasks: visual question answering (VQA [149]) and visual reasoning (NLVR-2 [150]). For the VQA task, ONE-PEACE achieves a score of 82.6 on the test-dev set and 82.5 on the test-std set, outperforming previous strong baselines like CoCa and BLIP-2. For the NLVR2 task, ONE-PEACE surpasses CoCa with gains of 1.7 and 1.3 on the dev set and test-P set respectively. Notice that our results on both tasks are lower than BEiT-3. This may be attributed to two reasons: Firstly, BEiT-3 is pretrained on in-domain datasets such as MSCOCO [126] and Visual Genome [148], which usually results in better downstream finetuning effects. Secondly, BEiT-3 incorporates pure text data for pretraining, which improves its language understanding ability and consequently enhances its multimodal understanding ability. In addition, OFA and BLIP-2 have shown that combined with language pretrained models can improve performance on multimodal understanding tasks. Therefore, we will explore the combination of ONE-PEACE and language pretrained models in the future.

### 5.4 Ablation Study

For the following ablation experiments, we utilize VIT-B/16 as the model backbone. The model is trained for 20 epochs with a batch size of 4096. We randomly selected 20 million image-text pairs from Laion-2B as the pretraining dataset.

**Ablation on Model Structures.** We first conduct ablation experiments to investigate the effects of sharing or separating different modules. As shown in Table 11, sharing both self-attention layers and FFN layers yields better results compared to not sharing. This suggests that utilizing a single Transformer can effectively align the semantic

| Structure | COCO zero-shot (5k test set) | | | | | | IN-1K |
| | Image → Text | | | Text → Image | | | |
| | R@1 | R@5 | R@10 | R@1 | R@5 | R@10 | ZS |
|---|---|---|---|---|---|---|---|
| Share ATTN & FFN | 33.96 | 60.92 | 72.24 | 22.51 | 46.13 | 57.63 | 42.21 |
| No Share | 29.08 | 54.30 | 65.56 | 18.98 | 40.34 | 51.44 | 37.73 |
| Share FFN | 27.36 | 52.66 | 64.36 | 18.08 | 38.68 | 49.84 | 33.98 |
| Share ATTN | **35.94** | **62.52** | **72.78** | **23.87** | **47.80** | **59.38** | **43.24** |

Table 11: **Ablation experiments on model structures.** "Share ATTN & FFN" means share both self-attention layers and FFN layers. "No Share" means separate both self-attention layers and FFN layers. "Share FFN" means separate self-attention layers and share FFN layers. "Share ATTN" means share self-attention layers and separate FFN layers, which is the default setting of ONE-PEACE. "ZS" is short for zero-shot accuracy.



(a) **Training loss**            (b) **Training accuracy**            (c) **Zero-shot accuracy**

Figure 4: **Training curves of different structures.** The model with shared self-attention layers and separated FFNs ("share attn") outperforms other structures, exhibiting the fastest convergence speed. (The curve appears as a staircase shape because we use exponential moving average to calculate relevant indicators in each epoch.)

space of vision and language. Furthermore, it is more beneficial to separate the FFN layer instead of sharing it. This implies that separating the FFN layer enhances the model's ability to extract modality-specific information, leading to more accurate representations. We also find that separating the self-attention layer and sharing the FFN layer yields the poorest results. We speculate that this is due to the self-attention layer playing a more significant role in aligning modalities compared to the FFN layer. Therefore, separating the self-attention layer lead to inferior performance. Figure 4 demonstrates the convergence performance of different architectures. Among all the architectures, the model with shared self-attention layers and separated FFNs exhibits the fastest convergence speed.

**Effects of Intra-modal Denoising Contrastive Learning.**    We examine the effects of intra-modal denoising contrastive learning (DCL). As shown in Table 12, applying DCL to language data (DCL-L) can enhance the model's performance in text retrieval tasks. Furthermore, applying DCL to vision data (DCL-V) can improve the model's cross-modal retrieval ability, as well as fine-tuning performance in image classification. By applying DCL to vision-language data (DCL-VL), ONE-PEACE achieves the best results in terms of all the evaluation metrics. These results demonstrate that intra-modal denoising contrastive learning can complement cross-modal contrastive learning. It not only enables ONE-PEACE to achieve excellent downstream fine-tuning performance but also enhances the model's capability for zero-shot cross-modal retrieval.

**Ablation on Different Denoising Losses.**    We conduct a systematic comparison of different denoising losses, including the smooth L1 loss used in [41, 151], the L2 loss used in [152, 153], the cosine loss used in [154, 105, 155], and the denoising contrastive loss used in this paper. As shown in Table 13, different types of denoising loss can improve the performance of the model in both cross-modal retrieval and image classification tasks. Among all the denoising losses, the denoising contrastive loss has the greatest improvement in terms of all the metrics compared to other losses. For example, it increased by +1.64 on COCO text retrieval R@1, increased by +1.47 on COCO image retrieval R@1, and increased by +0.6 on image classification. This indicates that denoising contrastive loss is more compatible with cross-modal contrastive loss than other denoising losses.

12

| CL | DCL-L | DCL-V | DCL-VL | COCO zero-shot (5k test set) | | | | | | IN-1K | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | | Image → Text | | | Text → Image | | | | |
| | | | | R@1 | R@5 | R@10 | R@1 | R@5 | R@10 | ZS | FT |
| ✓ | | | | 35.94 | 62.52 | 72.78 | 23.87 | 47.80 | 59.38 | 43.24 | 82.20 |
| ✓ | ✓ | | | 37.02 | 63.78 | 73.96 | 23.63 | 47.87 | 59.22 | 43.69 | 81.99 |
| ✓ | | ✓ | | 38.88 | 65.34 | 75.76 | 26.05 | 49.78 | 61.26 | 45.94 | 83.32 |
| ✓ | ✓ | ✓ | | 39.00 | 65.64 | 76.30 | 25.85 | 50.06 | 61.80 | 45.54 | 83.33 |
| ✓ | ✓ | ✓ | ✓ | **39.94** | **65.94** | **76.72** | **26.94** | **51.38** | **62.81** | **46.41** | **83.75** |

Table 12: **Ablation studies of intra-modal denoising contrastive learning.** "CL" is cross-modal contrastive learning. "DCL-L", "DCL-V", and "DCL-VL" means applying intra-modal denoising contrastive learning to language, vision, and vision-language data, respectively. "ZS" is short for zero-shot accuracy, "FT" is short for fine-tuning accuracy.

| Denoising Loss | COCO zero-shot (5k test set) | | | | | | IN-1K | |
|---|---|---|---|---|---|---|---|---|
| | Image → Text | | | Text → Image | | | | |
| | R@1 | R@5 | R@10 | R@1 | R@5 | R@10 | ZS | FT |
| None | 35.94 | 62.52 | 72.78 | 23.87 | 47.80 | 59.38 | 43.24 | 82.20 |
| Smooth L1 Loss | 36.72 | 63.86 | 74.30 | 24.37 | 48.17 | 59.86 | 44.36 | 82.50 |
| L2 Loss | 37.46 | 64.80 | 74.46 | 25.81 | 49.64 | 61.80 | 45.78 | 83.15 |
| Cosine Loss | 38.28 | 65.80 | 75.18 | 25.47 | 50.02 | 61.59 | 45.15 | 83.07 |
| Denoising Contrastive Loss | **39.94** | **65.94** | **76.72** | **26.94** | **51.38** | **62.81** | **46.41** | **83.75** |

Table 13: **Ablation studies of different denoising losses.** Among all the denoising losses, denoising contrastive loss shows the greatest performance improvement in both cross-modal retrieval and image classification tasks.

### 5.5 Emergent Zero-shot Retrieval

In our pretraining, we exclusively align other modalities with text which plays as an intermediary role. We assume that our model is able to align those modalities that are not paired in the training data. For example, ONE-PEACE should be able to align image and audio. Thus, we conduct experiments on the retrieval of those modalities to assess the emergent zero-shot capabilities [156].

To be more specific, we evaluate the audio-to-image, audio+image-to-image, and audio+text-to-image retrieval abilities and demonstrate case studies in Figure 5. The first two cases demonstrate the emergent capability of uni-modal retrieval, while the other cases show that of the retrieval of image based on multimodal inputs. Specifically, we find that ONE-PEACE is able to retrieve images that contain elements concerning inputs of different modalities, e.g., the model uses the text "snow" and the sound of bird chirping to retrieve the images of birds in the snow. These examples demonstrate that ONE-PEACE has strong potential in emergent zero-shot capabilities. This indicates that for a universal representation model, there is no need to learn all pairing relationships between modalities, but instead it is sufficient for modalities to be aligned to an intermediary one. We provide more quality examples in Appendix E.

## 6 Conclusion, Limitation and Future Work

In this work, we explore a scalable way for building a general representation model across different modalities. Based on the flexible architecture and modality-agnostic pretraining tasks, we release ONE-PEACE, a general representation model that can seamlessly align and integrate representations across vision, audio, and language modalities. We conduct a series of experiments across 3 modalities, 11 tasks, and 16 datasets. The experimental results demonstrate that ONE-PEACE achieves leading results in a wide range of tasks, including image classification, semantic segmentation, audio-text retrieval, audio classification, audio question answering, image-text retrieval, and visual grounding. Furthermore, we show that ONE-PEACE possesses a strong emergent zero-shot retrieval capability, enabling it to align modalities that are not paired in the training data.

**Limitation.** Although ONE-PEACE achieves leading results in a wide range of tasks, it falls short of achieving state-of-the-art results in zero-shot image-text retrieval and vision-language understanding tasks. There are two possible
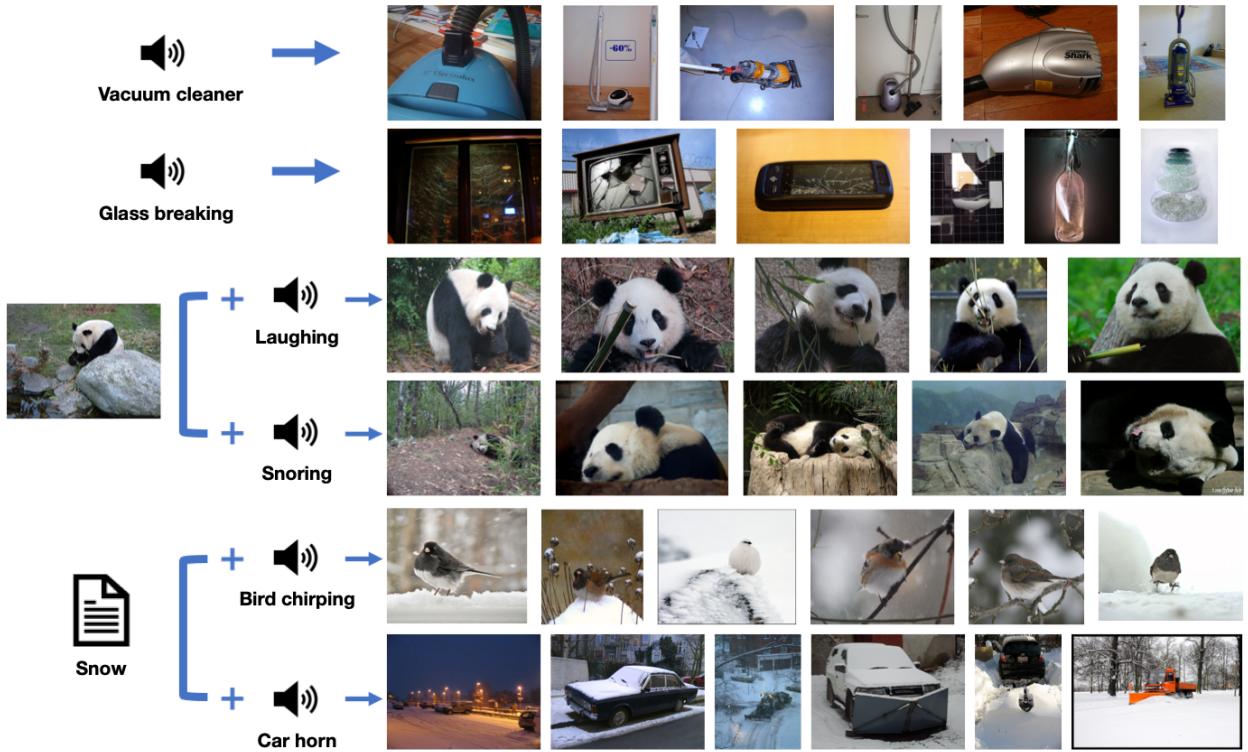
Figure 5: Examples of emergent zero-shot retrieval. ONE-PEACE is capable of aligning modalities and modality combinations, where there are no paired data of the modalities in the pretraining dataset. The images are retrieved from ImageNet-1K and MSCOCO.

reasons for this: 1). ONE-PEACE didn't see enough image-text pairs during pretraining. We only trained on 6.4 billion image-text pairs, while previous works [30, 59, 33] typically train on 12.8 billion image-text pairs or more. 2). ONE-PEACE didn't use language pretrained models for initialization or introduce any pure text data. Both the vision and language modules of ONE-PEACE are completely randomly initialized, while previous works [31, 21] show that introducing pure text data or initialized with the language pretrained models can greatly enhance the model's performance. In fact, as a highly extensible model, ONE-PEACE can combine with language pretrained models to achieve better results.

**Future Work.** In the future, we will test ONE-PEACE on more downstream tasks, such as vision-audio-language tasks and extend to more modalities for pretraining like video, 3D point cloud, etc. Also, we pursue an active interaction with large language models (LLMs) to continue influencing broader areas. This includes:

- With the help of LLMs, building a more powerful general representation model.
- By combining LLMs, creating a more general multimodal language model.

## Acknowledgments

## References

[1] Hangbo Bao, Li Dong, and Furu Wei. Beit: Bert pre-training of image transformers. In *ICLR*, 2021. 1, 2

[2] Kaiming He, Xinlei Chen, Saining Xie, Yanghao Li, Piotr Dollár, and Ross Girshick. Masked autoencoders are scalable vision learners. In *CVPR*, 2022. 1, 2, 5

[3] Jinghao Zhou, Chen Wei, Huiyu Wang, Wei Shen, Cihang Xie, Alan Yuille, and Tao Kong. ibot: Image bert pre-training with online tokenizer. In *ICLR*, 2022. 1

[4] Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross B. Girshick. Momentum contrast for unsupervised visual representation learning. In *CVPR*, 2020. 1, 2

[5] Xinlei Chen, Haoqi Fan, Ross Girshick, and Kaiming He. Improved baselines with momentum contrastive learning. *arXiv:2003.04297*, 2020. 1

[6] Xinlei Chen and Kaiming He. Exploring simple siamese representation learning. In *CVPR*, 2021. 1

[7] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for contrastive learning of visual representations. In *ICML*, 2020. 1, 2

[8] Mathilde Caron, Hugo Touvron, Ishan Misra, Hervé Jégou, Julien Mairal, Piotr Bojanowski, and Armand Joulin. Emerging properties in self-supervised vision transformers. *arXiv:2104.14294*, 2021. 1

[9] Maxime Oquab, Timoth'ee Darcet, Th'eo Moutakanni, Huy Q. Vo, Marc Szafraniec, Vasil Khalidov, Pierre Fernandez, Daniel Haziza, Francisco Massa, Alaaeldin El-Nouby, Mahmoud Assran, Nicolas Ballas, Wojciech Galuba, Russ Howes, Po-Yao Huang, Shang-Wen Li, Ishan Misra, Michael G. Rabbat, Vasu Sharma, Gabriel Synnaeve, Huijiao Xu, Hervé Jégou, Julien Mairal, Patrick Labatut, Armand Joulin, and Piotr Bojanowski. Dinov2: Learning robust visual features without supervision. *arXiv:2304.07193*, 2023. 1

[10] Steffen Schneider, Alexei Baevski, Ronan Collobert, and Michael Auli. wav2vec: Unsupervised pre-training for speech recognition. In *Interspeech*, 2019. 1, 2

[11] Alexei Baevski, Yuhao Zhou, Abdelrahman Mohamed, and Michael Auli. wav2vec 2.0: A framework for self-supervised learning of speech representations. In *NeurIPS*, 2020. 1, 4, 26

[12] Sanyuan Chen, Chengyi Wang, Zhengyang Chen, Yu Wu, Shujie Liu, Zhuo Chen, Jinyu Li, Naoyuki Kanda, Takuya Yoshioka, Xiong Xiao, Jian Wu, Long Zhou, Shuo Ren, Yanmin Qian, Yao Qian, Micheal Zeng, and Furu Wei. Wavlm: Large-scale self-supervised pre-training for full stack speech processing. In *JSTSP*, 2021. 1, 6, 26

[13] Wei-Ning Hsu, Benjamin Bolte, Yao-Hung Hubert Tsai, Kushal Lakhotia, Ruslan Salakhutdinov, and Abdelrahman Mohamed. Hubert: Self-supervised speech representation learning by masked prediction of hidden units. In *TASLP*, 2021. 1, 26

[14] Matthew E. Peters, Mark Neumann, Mohit Iyyer, Matt Gardner, Christopher Clark, Kenton Lee, and Luke Zettlemoyer. Deep contextualized word representations. In *NAACL*, 2018. 1

[15] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: pre-training of deep bidirectional transformers for language understanding. In *NAACL*, 2019. 1, 2

[16] Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. Roberta: A robustly optimized bert pretraining approach. *arXiv:1907.11692*, 2019. 1, 3

[17] Kevin Clark, Minh-Thang Luong, Quoc V. Le, and Christopher D. Manning. ELECTRA: pre-training text encoders as discriminators rather than generators. In *ICLR*, 2020. 1

[18] Pengcheng He, Xiaodong Liu, Jianfeng Gao, and Weizhu Chen. Deberta: decoding-enhanced bert with disentangled attention. In *ICLR*, 2021. 1

[19] OpenAI. Gpt-4 technical report. *arXiv:2303.08774*, 2023. 1

[20] Shaohan Huang, Li Dong, Wenhui Wang, Yaru Hao, Saksham Singhal, Shuming Ma, Tengchao Lv, Lei Cui, Owais Khan Mohammed, Qiang Liu, Kriti Aggarwal, Zewen Chi, Johan Bjorck, Vishrav Chaudhary, Subhojit Som, Xia Song, and Furu Wei. Language is not all you need: Aligning perception with language models. *arXiv:2302.14045*, 2023. 1

[21] Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. Blip-2: Bootstrapping language-image pre-training with frozen image encoders and large language models. *arXiv:2301.12597*, 2023. 1, 2, 11, 14

[22] Deyao Zhu, Jun Chen, Xiaoqian Shen, Xiang Li, and Mohamed Elhoseiny. Minigpt-4: Enhancing vision-language understanding with advanced large language models. *arXiv:2304.10592*, 2023. 1

[23] Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. *arXiv:2304.08485*, 2023. 1

[24] Qinghao Ye, Haiyang Xu, Guohai Xu, Jiabo Ye, Ming Yan, Yi Zhou, Junyan Wang, Anwen Hu, Pengcheng Shi, Yaya Shi, Chenliang Li, Yuanhong Xu, Hehong Chen, Junfeng Tian, Qiang Qi, Ji Chao Zhang, and Feiyan Huang. mplug-owl: Modularization empowers large language models with multimodality. *arXiv:2304.14178*, 2023. 1

[25] Danny Driess, F. Xia, Mehdi S. M. Sajjadi, Corey Lynch, Aakanksha Chowdhery, Brian Ichter, Ayzaan Wahid, Jonathan Tompson, Quan Ho Vuong, Tianhe Yu, Wenlong Huang, Yevgen Chebotar, Pierre Sermanet, Daniel Duckworth, Sergey Levine, Vincent Vanhoucke, Karol Hausman, Marc Toussaint, Klaus Greff, Andy Zeng, Igor Mordatch, and Peter R. Florence. Palm-e: An embodied multimodal language model. *arXiv:2303.03378*, 2023. 1

[26] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *NeurIPS*, 2017. 1, 3, 4

[27] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale. In *ICLR*, 2021. 1, 2, 3

[28] Andrew Jaegle, Felix Gimeno, Andrew Brock, Andrew Zisserman, Oriol Vinyals, and Joao Carreira. Perceiver: General perception with iterative attention. *arXiv:2103.03206*, 2021. 1

[29] Andrew Jaegle, Sebastian Borgeaud, Jean-Baptiste Alayrac, Carl Doersch, Catalin Ionescu, David Ding, Skanda Koppula, Daniel Zoran, Andrew Brock, Evan Shelhamer, et al. Perceiver io: A general architecture for structured inputs & outputs. *arXiv:2107.14795*, 2021. 1

[30] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. Learning transferable visual models from natural language supervision. In *ICML*, 2021. 1, 2, 5, 10, 14

[31] Peng Wang, An Yang, Rui Men, Junyang Lin, Shuai Bai, Zhikang Li, Jianxin Ma, Chang Zhou, Jingren Zhou, and Hongxia Yang. Ofa: Unifying architectures, tasks, and modalities through a simple sequence-to-sequence learning framework. In *ICML*, 2022. 1, 2, 4, 10, 11, 14, 26

[32] Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc, Antoine Miech, Iain Barr, Yana Hasson, Karel Lenc, Arthur Mensch, Katie Millican, Malcolm Reynolds, Roman Ring, Eliza Rutherford, Serkan Cabi, Tengda Han, Zhitao Gong, Sina Samangooei, Marianne Monteiro, Jacob Menick, Sebastian Borgeaud, Andy Brock, Aida Nematzadeh, Sahand Sharifzadeh, Mikolaj Binkowski, Ricardo Barreira, Oriol Vinyals, Andrew Zisserman, and Karen Simonyan. Flamingo: a visual language model for few-shot learning. *arXiv:2204.14198*, 2022. 1, 11

[33] Jiahui Yu, Zirui Wang, Vijay Vasudevan, Legg Yeung, Mojtaba Seyedhosseini, and Yonghui Wu. Coca: Contrastive captioners are image-text foundation models. In *TMLR*, 2022. 1, 2, 4, 7, 8, 10, 11, 14

[34] Wenhui Wang, Hangbo Bao, Li Dong, Johan Bjorck, Zhiliang Peng, Qiang Liu, Kriti Aggarwal, Owais Khan Mohammed, Saksham Singhal, Subhojit Som, and Furu Wei. Image as a foreign language: BEiT pretraining for vision and vision-language tasks. In *CVPR*, 2023. 1, 2, 4, 7, 10, 11

[35] Junnan Li, Dongxu Li, Caiming Xiong, and Steven C. H. Hoi. Blip: Bootstrapping language-image pre-training for unified vision-language understanding and generation. In *ICML*, 2022. 1, 2, 11, 26

[36] Xi Chen, Xiao Wang, Soravit Changpinyo, A. J. Piergiovanni, Piotr Padlewski, Daniel M. Salz, Sebastian Goodman, Adam Grycner, Basil Mustafa, Lucas Beyer, Alexander Kolesnikov, Joan Puigcerver, Nan Ding, Keran Rong, Hassan Akbari, Gaurav Mishra, Linting Xue, Ashish V. Thapliyal, James Bradbury, Weicheng Kuo, Mojtaba Seyedhosseini, Chao Jia, Burcu Karagol Ayan, Carlos Riquelme, Andreas Steiner, Anelia Angelova, Xiaohua Zhai, Neil Houlsby, and Radu Soricut. Pali: A jointly-scaled multilingual language-image model. *arXiv:2209.06794*, 2022. 1, 2

[37] Andrey Guzhov, Federico Raue, Jörn Hees, and Andreas R. Dengel. Audioclip: Extending clip to image, text and audio. In *ICASSP*, 2022. 1, 2

[38] Ho-Hsiang Wu, Prem Seetharaman, Kundan Kumar, and Juan Pablo Bello. Wav2clip: Learning robust audio representations from clip. In *ICASSP*, 2022. 1, 2

[39] Benjamin Elizalde, Soham Deshmukh, Mahmoud Al Ismail, and Huaming Wang. Clap: Learning audio concepts from natural language supervision. *arXiv:2206.04769*, 2022. 1, 2, 9

[40] Junyi Ao, Rui Wang, Long Zhou, Chengyi Wang, Shuo Ren, Yu Wu, Shujie Liu, Tom Ko, Qing Li, Yu Zhang, et al. Speecht5: Unified-modal encoder-decoder pre-training for spoken language processing. In *ACL*, 2022. 1, 2

[41] Alexei Baevski, Wei-Ning Hsu, Qiantong Xu, Arun Babu, Jiatao Gu, and Michael Auli. Data2vec: A general framework for self-supervised learning in speech, vision and language. In *ICML*, 2022. 1, 3, 12

[42] Weijie Su, Xizhou Zhu, Yue Cao, Bin Li, Lewei Lu, Furu Wei, and Jifeng Dai. Vl-bert: Pre-training of generic visual-linguistic representations. In *ICLR*, 2019. 2

[43] Jiasen Lu, Dhruv Batra, Devi Parikh, and Stefan Lee. Vilbert: Pretraining task-agnostic visiolinguistic representations for vision-and-language tasks. In *NeurIPS*, 2019. 2

[44] Liunian Harold Li, Mark Yatskar, Da Yin, Cho-Jui Hsieh, and Kai-Wei Chang. Visualbert: A simple and performant baseline for vision and language. *arXiv:1908.03557*, 2019. 2

[45] Yen-Chun Chen, Linjie Li, Licheng Yu, Ahmed El Kholy, Faisal Ahmed, Zhe Gan, Yu Cheng, and Jingjing Liu. Uniter: Universal image-text representation learning. In *ECCV*, 2020. 2, 10

[46] Junyang Lin, An Yang, Yichang Zhang, Jie Liu, Jingren Zhou, and Hongxia Yang. Interbert: Vision-and-language interaction for multi-modal pretraining. *arXiv:2003.13198*, 2020. 2

[47] Xiujun Li, Xi Yin, Chunyuan Li, Xiaowei Hu, Pengchuan Zhang, Lei Zhang, Lijuan Wang, Houdong Hu, Li Dong, Furu Wei, Yejin Choi, and Jianfeng Gao. Oscar: Object-semantics aligned pre-training for vision-language tasks. In *ECCV*, 2020. 2

[48] Pengchuan Zhang, Xiujun Li, Xiaowei Hu, Jianwei Yang, Lei Zhang, Lijuan Wang, Yejin Choi, and Jianfeng Gao. Vinvl: Revisiting visual representations in vision-language models. In *CVPR*, 2021. 2

[49] Wei Li, Can Gao, Guocheng Niu, Xinyan Xiao, Hao Liu, Jiachen Liu, Hua Wu, and Haifeng Wang. UNIMO: towards unified-modal understanding and generation via cross-modal contrastive learning. In *ACL*, 2021. 2

[50] Jaemin Cho, Jie Lei, Haochen Tan, and Mohit Bansal. Unifying vision-and-language tasks via text generation. In *ICML*, 2021. 2, 10

[51] Junyang Lin, Rui Men, An Yang, Chang Zhou, Ming Ding, Yichang Zhang, Peng Wang, Ang Wang, Le Jiang, Xianyan Jia, et al. M6: A chinese multimodal pretrainer. In *KDD*, 2021. 2

[52] Wonjae Kim, Bokyung Son, and Ildoo Kim. Vilt: Vision-and-language transformer without convolution or region supervision. In *ICML*, 2021. 2

[53] Junnan Li, Ramprasaath R Selvaraju, Akhilesh Deepak Gotmare, Shafiq Joty, Caiming Xiong, and Steven Hoi. Align before fuse: Vision and language representation learning with momentum distillation. In *NeurIPS*, 2021. 2, 11, 26

[54] Zi-Yi Dou, Yichong Xu, Zhe Gan, Jianfeng Wang, Shuohang Wang, Lijuan Wang, Chenguang Zhu, Nanyun Peng, Zicheng Liu, and Michael Zeng. An empirical study of training end-to-end vision-and-language transformers. In *CVPR*, 2022. 2

[55] Hangbo Bao, Wenhui Wang, Li Dong, Qiang Liu, Owais Khan Mohammed, Kriti Aggarwal, Subhojit Som, Songhao Piao, and Furu Wei. Vlmo: Unified vision-language pre-training with mixture-of-modality-experts. In *NeurIPS*, 2022. 2

[56] Zi-Yi* Dou, Aishwarya* Kamath, Zhe* Gan, Pengchuan Zhang, Jianfeng Wang, Linjie Li, Zicheng Liu, Ce Liu, Yann LeCun, Nanyun Peng, Jianfeng Gao, and Lijuan Wang. Coarse-to-fine vision-language pre-training with fusion in the backbone. In *NeurIPS*, 2022. 2, 10

[57] Xiaowei Hu, Zhe Gan, Jianfeng Wang, Zhengyuan Yang, Zicheng Liu, Yumao Lu, and Lijuan Wang. Scaling up vision-language pretraining for image captioning. In *CVPR*, 2022. 2

[58] Haiyang Xu, Qinghao Ye, Mingshi Yan, Yaya Shi, Jiabo Ye, Yuanhong Xu, Chenliang Li, Bin Bi, Qiuchen Qian, Wei Wang, Guohai Xu, Ji Zhang, Songfang Huang, Feiran Huang, and Jingren Zhou. mplug-2: A modularized multi-modal foundation model across text, image and video. *arXiv:2302.00402*, 2023. 2

[59] Chao Jia, Yinfei Yang, Ye Xia, Yi-Ting Chen, Zarana Parekh, Hieu Pham, Quoc V Le, Yunhsuan Sung, Zhen Li, and Tom Duerig. Scaling up visual and vision-language representation learning with noisy text supervision. *arXiv:2102.05918*, 2021. 2, 5, 7, 9, 10, 14

[60] Xiaohua Zhai, Xiao Wang, Basil Mustafa, Andreas Steiner, Daniel Keysers, Alexander Kolesnikov, and Lucas Beyer. Lit: Zero-shot transfer with locked-image text tuning. In *CVPR*, 2022. 2

[61] Lu Yuan, Dongdong Chen, Yi-Ling Chen, Noel C. F. Codella, Xiyang Dai, Jianfeng Gao, Houdong Hu, Xuedong Huang, Boxin Li, Chunyuan Li, Ce Liu, Mengchen Liu, Zicheng Liu, Yumao Lu, Yu Shi, Lijuan Wang, Jianfeng Wang, Bin Xiao, Zhen Xiao, Jianwei Yang, Michael Zeng, Luowei Zhou, and Pengchuan Zhang. Florence: A new foundation model for computer vision. *arXiv:2111.11432*, 2021. 2, 8, 10

[62] An Yang, Junshu Pan, Junyang Lin, Rui Men, Yichang Zhang, Jingren Zhou, and Chang Zhou. Chinese clip: Contrastive vision-language pretraining in chinese. *arXiv preprint arXiv:2211.01335*, 2022. 2

17

[63] Zirui Wang, Jiahui Yu, Adams Wei Yu, Zihang Dai, Yulia Tsvetkov, and Yuan Cao. Simvlm: Simple visual language model pretraining with weak supervision. In *ICLR*, 2022. 2, 4, 11

[64] Jianfeng Wang, Zhengyuan Yang, Xiaowei Hu, Linjie Li, Kevin Lin, Zhe Gan, Zicheng Liu, Ce Liu, and Lijuan Wang. Git: A generative image-to-text transformer for vision and language. *arXiv:2205.14100*, 2022. 2

[65] Xizhou Zhu, Jinguo Zhu, Hao Li, Xiaoshi Wu, Xiaogang Wang, Hongsheng Li, Xiaohua Wang, and Jifeng Dai. Uni-perceiver: Pre-training unified architecture for generic perception for zero-shot and few-shot tasks. In *CVPR*, 2022. 2

[66] Amanpreet Singh, Ronghang Hu, Vedanuj Goswami, Guillaume Couairon, Wojciech Galuba, Marcus Rohrbach, and Douwe Kiela. FLAVA: A foundational language and vision alignment model. In *CVPR*, 2022. 2

[67] Yaru Hao, Haoyu Song, Li Dong, Shaohan Huang, Zewen Chi, Wenhui Wang, Shuming Ma, and Furu Wei. Language models are general-purpose interfaces. *arXiv:2206.06336*, 2022. 2

[68] Jiasen Lu, Christopher Clark, Rowan Zellers, Roozbeh Mottaghi, and Aniruddha Kembhavi. Unified-io: A unified model for vision, language, and multi-modal tasks. In *ICLR*, 2023. 2

[69] Ankur Bapna, Yu-an Chung, Nan Wu, Anmol Gulati, Ye Jia, Jonathan H Clark, Melvin Johnson, Jason Riesa, Alexis Conneau, and Yu Zhang. Slam: A unified encoder for speech and language modeling via speech-text joint pre-training. *arXiv:2110.10329*, 2021. 2

[70] Ankur Bapna, Colin Cherry, Yu Zhang, Ye Jia, Melvin Johnson, Yong Cheng, Simran Khanuja, Jason Riesa, and Alexis Conneau. mslam: Massively multilingual joint pre-training for speech and text. *arXiv:2202.01374*, 2022. 2

[71] Zhehuai Chen, Yu Zhang, Andrew Rosenberg, Bhuvana Ramabhadran, Pedro J. Moreno, Ankur Bapna, and Heiga Zen. Maestro: Matched speech text representations through modality matching. In *Interspeech*, 2022. 2

[72] Ziqiang Zhang, Sanyuan Chen, Long Zhou, Yu Wu, Shuo Ren, Shujie Liu, Zhuoyuan Yao, Xun Gong, Lirong Dai, Jinyu Li, et al. Speechlm: Enhanced speech pre-training with unpaired textual data. *arXiv:2209.15329*, 2022. 2

[73] Yun Tang, Hongyu Gong, Ning Dong, Changhan Wang, Wei-Ning Hsu, Jiatao Gu, Alexei Baevski, Xian Li, Abdelrahman Mohamed, Michael Auli, et al. Unified speech-text pre-training for speech translation and recognition. In *ACL*, 2022. 2

[74] Ziqiang Zhang, Long Zhou, Junyi Ao, Shujie Liu, Lirong Dai, Jinyu Li, and Furu Wei. Speechut: Bridging speech and text with hidden-unit for encoder-decoder based speech-text pre-training. In *EMNLP*, 2022. 2

[75] Xiaohuan Zhou, Jiaming Wang, Zeyu Cui, Shiliang Zhang, Zhijie Yan, Jingren Zhou, and Chang Zhou. Mm-speech: Multi-modal multi-task encoder-decoder pre-training for speech recognition. *arXiv:2212.00500*, 2022. 2

[76] Yusong Wu, K. Chen, Tianyu Zhang, Yuchen Hui, Taylor Berg-Kirkpatrick, and Shlomo Dubnov. Large-scale contrastive language-audio pretraining with feature fusion and keyword-to-caption augmentation. In *ICASSP*, 2023. 2, 3, 9, 23

[77] Qiuqiang Kong, Yin Cao, Turab Iqbal, Yuxuan Wang, Wenwu Wang, and Mark D. Plumbley. Panns: Large-scale pretrained audio neural networks for audio pattern recognition. In *TASLP*, 2019. 2

[78] Ke Chen, Xingjian Du, Bilei Zhu, Zejun Ma, Taylor Berg-Kirkpatrick, and Shlomo Dubnov. Hts-at: A hierarchical token-semantic audio transformer for sound classification and detection. In *ICASSP*, 2022. 3

[79] Hassan Akbari, Linagzhe Yuan, Rui Qian, Wei-Hong Chuang, Shih-Fu Chang, Yin Cui, and Boqing Gong. VATT: Transformers for multimodal self-supervised learning from raw video, audio and text. In *NeurIPS*, 2021. 3, 8

[80] Ziyi Yang, Yuwei Fang, Chenguang Zhu, Reid Pryzant, Dongdong Chen, Yu Shi, Yichong Xu, Yao Qian, Mei Gao, Yi-Ling Chen, Liyang Lu, Yujia Xie, Robert Gmyr, Noel C. F. Codella, Naoyuki Kanda, Bin Xiao, Yuanxun Lu, Takuya Yoshioka, Michael Zeng, and Xuedong Huang. i-code: An integrative and composable multimodal learning framework. *arXiv:2205.01818*, 2022. 3

[81] Qiu shi Zhu, Long Zhou, Zi-Hua Zhang, Shujie Liu, Binxing Jiao, J. Zhang, Lirong Dai, Daxin Jiang, Jinyu Li, and Furu Wei. Vatlm: Visual-audio-text pre-training with unified masked prediction for speech representation learning. *arXiv:2211.11275*, 2022. 3

[82] Aäron van den Oord, Oriol Vinyals, and Koray Kavukcuoglu. Neural discrete representation learning. In *NeurIPS*, 2017. 3

[83] Bowen Shi, Wei-Ning Hsu, Kushal Lakhotia, and Abdel rahman Mohamed. Learning audio-visual speech representation by masked multimodal cluster prediction. In *ICLR*, 2022. 3

[84] Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo. Swin transformer: Hierarchical vision transformer using shifted windows. In *ICCV*, 2021. 3

[85] Yann LeCun, Léon Bottou, Yoshua Bengio, and Patrick Haffner. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 1998. 3

[86] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *CVPR*, 2016. 3

[87] Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural Computation*, 1997. 3

[88] Kyunghyun Cho, Bart van Merrienboer, Çaglar Gülçehre, Dzmitry Bahdanau, Fethi Bougares, Holger Schwenk, and Yoshua Bengio. Learning phrase representations using rnn encoder–decoder for statistical machine translation. In *EMNLP*, 2014. 3

[89] Hugo Touvron, Matthieu Cord, Alaaeldin El-Nouby, Jakob Verbeek, and Herve Jegou. Three things everyone should know about vision transformers. In *ECCV*, 2022. 3

[90] Abdel rahman Mohamed, Dmytro Okhonko, and Luke Zettlemoyer. Transformers with convolutional context for asr. *arXiv:1904.11660*, 2019. 4

[91] Rico Sennrich, Barry Haddow, and Alexandra Birch. Neural machine translation of rare words with subword units. In *ACL*, 2016. 4

[92] Jinze Bai, Rui Men, Han Yang, Xuancheng Ren, Kai fung Edward Dang, Yichang Zhang, Xiaohuan Zhou, Peng Wang, Sinan Tan, An Yang, Zeyu Cui, Yu Han, Shuai Bai, Wenhang Ge, Jianxin Ma, Junyang Lin, Jingren Zhou, and Chang Zhou. Ofasys: A multi-modal multi-task learning system for building generalist models. *ArXiv*, abs/2212.04408, 2022. 4

[93] Hongyu Wang, Shuming Ma, Shaohan Huang, Li Dong, Wenhui Wang, Zhiliang Peng, Yu Wu, Payal Bajaj, Saksham Singhal, Alon Benhaim, Barun Patra, Zhun Liu, Vishrav Chaudhary, Xia Song, and Furu Wei. Foundation transformers. *arXiv:2210.06423*, 2022. 4

[94] Tom B Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. *arXiv:2005.14165*, 2020. 4

[95] Noam M. Shazeer. Glu variants improve transformer. *arXiv:2002.05202*, 2020. 4

[96] Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, Parker Schuh, Kensen Shi, Sasha Tsvyashchenko, Joshua Maynez, Abhishek Rao, Parker Barnes, Yi Tay, Noam M. Shazeer, Vinodkumar Prabhakaran, Emily Reif, Nan Du, Benton C. Hutchinson, Reiner Pope, James Bradbury, Jacob Austin, Michael Isard, Guy Gur-Ari, Pengcheng Yin, Toju Duke, Anselm Levskaya, Sanjay Ghemawat, Sunipa Dev, Henryk Michalewski, Xavier García, Vedant Misra, Kevin Robinson, Liam Fedus, Denny Zhou, Daphne Ippolito, David Luan, Hyeontaek Lim, Barret Zoph, Alexander Spiridonov, Ryan Sepassi, David Dohan, Shivani Agrawal, Mark Omernick, Andrew M. Dai, Thanumalayan Sankaranarayana Pillai, Marie Pellat, Aitor Lewkowycz, Erica Moreira, Rewon Child, Oleksandr Polozov, Katherine Lee, Zongwei Zhou, Xuezhi Wang, Brennan Saeta, Mark Díaz, Orhan Firat, Michele Catasta, Jason Wei, Kathleen S. Meier-Hellstern, Douglas Eck, Jeff Dean, Slav Petrov, and Noah Fiedel. Palm: Scaling language modeling with pathways. *arXiv:2204.02311*, 2022. 4

[97] Colin Raffel, Noam M. Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. Exploring the limits of transfer learning with a unified text-to-text transformer. In *JMLR*, 2019. 4, 23

[98] Zihang Dai, Hanxiao Liu, Quoc V Le, and Mingxing Tan. Coatnet: Marrying convolution and attention for all data sizes. In *NeurIPS*, 2021. 4

[99] Hugo Touvron, Matthieu Cord, Alexandre Sablayrolles, Gabriel Synnaeve, and Herv'e J'egou. Going deeper with image transformers. In *ICCV*, 2021. 4

[100] Yixuan Wei, Han Hu, Zhenda Xie, Zheng Zhang, Yue Cao, Jianmin Bao, Dong Chen, and Baining Guo. Contrastive learning rivals masked image modeling in fine-tuning via feature distillation. *arXiv:2205.14141*, 2022. 5, 6

[101] Kun Yi, Yixiao Ge, Xiaotong Li, Shusheng Yang, Dian Li, Jianping Wu, Ying Shan, and Xiaohu Qie. Masked image modeling with denoising contrast. *arXiv:2205.09616*, 2022. 5

[102] Christoph Schuhmann, Romain Beaumont, Richard Vencu, Cade Gordon, Ross Wightman, Mehdi Cherti, Theo Coombes, Aarush Katta, Clayton Mullis, Mitchell Wortsman, Patrick Schramowski, Srivatsa Kundurthy, Katherine Crowson, Ludwig Schmidt, Robert Kaczmarczyk, and Jenia Jitsev. Laion-5b: An open large-scale dataset for training next generation image-text models. *arXiv:2210.08402*, 2022. 6, 23

[103] Yixuan Wei, Han Hu, Zhenda Xie, Zheng Zhang, Yue Cao, Jianmin Bao, Dong Chen, and Baining Guo. Contrastive learning rivals masked image modeling in fine-tuning via feature distillation. *arXiv:2205.14141*, 2022. 7

[104] Wenhai Wang, Jifeng Dai, Zhe Chen, Zhenhang Huang, Zhiqi Li, Xizhou Zhu, Xiaowei Hu, Tong Lu, Lewei Lu, Hongsheng Li, et al. Internimage: Exploring large-scale vision foundation models with deformable convolutions. In *CVPR*, 2023. 7

[105] Yuxin Fang, Wen Wang, Binhui Xie, Quan Sun, Ledell Wu, Xinggang Wang, Tiejun Huang, Xinlong Wang, and Yue Cao. Eva: Exploring the limits of masked visual representation learning at scale. In *CVPR*, 2023. 7, 8, 12

[106] Yuxuan Cai, Yizhuang Zhou, Qi Han, Jianjian Sun, Xiangwen Kong, Jun Li, and Xiangyu Zhang. Reversible column networks. In *ICLR*, 2023. 7, 8

[107] Mitchell Wortsman, Gabriel Ilharco, Samir Ya Gadre, Rebecca Roelofs, Raphael Gontijo-Lopes, Ari S Morcos, Hongseok Namkoong, Ali Farhadi, Yair Carmon, Simon Kornblith, et al. Model soups: averaging weights of multiple fine-tuned models improves accuracy without increasing inference time. In *ICML*, 2022. 7

[108] Xiaohua Zhai, Alexander Kolesnikov, Neil Houlsby, and Lucas Beyer. Scaling vision transformers. In *CVPR*, 2022. 7

[109] Markus N Rabe and Charles Staats. Self-attention does not need $O(n^2)$ memory. *arXiv:2112.05682*, 2021. 7

[110] Tri Dao, Daniel Y Fu, Stefano Ermon, Atri Rudra, and Christopher Ré. Flashattention: Fast and memory-efficient exact attention with io-awareness. In *NeurIPS*, 2022. 7

[111] Tianqi Chen, Bing Xu, Chiyuan Zhang, and Carlos Guestrin. Training deep nets with sublinear memory cost. *arXiv:1604.06174*, 2016. 7

[112] Zhe Chen, Yuchen Duan, Wenhai Wang, Junjun He, Tong Lu, Jifeng Dai, and Yu Qiao. Vision transformer adapter for dense predictions. In *ICLR*, 2023. 7, 8

[113] Yanghao Li, Hanzi Mao, Ross Girshick, and Kaiming He. Exploring plain vision transformer backbones for object detection. In *ECCV*, 2022. 8, 24

[114] Anurag Arnab, Mostafa Dehghani, Georg Heigold, Chen Sun, Mario Lučić, and Cordelia Schmid. Vivit: A video vision transformer. In *ICCV*, 2021. 8

[115] Ze Liu, Han Hu, Yutong Lin, Zhuliang Yao, Zhenda Xie, Yixuan Wei, Jia Ning, Yue Cao, Zheng Zhang, Li Dong, et al. Swin transformer v2: Scaling up capacity and resolution. In *CVPR*, 2022. 8

[116] Christoph Feichtenhofer, Yanghao Li, Kaiming He, et al. Masked autoencoders as spatiotemporal learners. In *NeurIPS*, 2022. 8

[117] Zhan Tong, Yibing Song, Jue Wang, and Limin Wang. Videomae: Masked autoencoders are data-efficient learners for self-supervised video pre-training. In *NeurIPS*, 2022. 8

[118] Limin Wang, Bingkun Huang, Zhiyu Zhao, Zhan Tong, Yinan He, Yi Wang, Yali Wang, and Yu Qiao. Videomae v2: Scaling video masked autoencoders with dual masking. In *CVPR*, 2023. 8

[119] Chen Wei, Haoqi Fan, Saining Xie, Chao-Yuan Wu, Alan Yuille, and Christoph Feichtenhofer. Masked feature prediction for self-supervised visual pre-training. In *CVPR*, 2022. 8

[120] Mostafa Dehghani, Josip Djolonga, Basil Mustafa, Piotr Padlewski, Jonathan Heek, Justin Gilmer, Andreas Steiner, Mathilde Caron, Robert Geirhos, Ibrahim Alabdulmohsin, et al. Scaling vision transformers to 22 billion parameters. *arXiv:2302.05442*, 2023. 8

[121] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, et al. Imagenet large scale visual recognition challenge. In *IJCV*, 2015. 7, 23

[122] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *CVPR*, 2009. 7, 24, 26

[123] Bolei Zhou, Hang Zhao, Xavier Puig, Sanja Fidler, Adela Barriuso, and Antonio Torralba. Semantic understanding of scenes through the ade20k dataset. In *IJCV*, 2016. 8, 23, 24

[124] Bowen Cheng, Ishan Misra, Alexander G. Schwing, Alexander Kirillov, and Rohit Girdhar. Masked-attention mask transformer for universal image segmentation. In *CVPR*, 2022. 8, 24

[125] Holger Caesar, Jasper Uijlings, and Vittorio Ferrari. Coco-stuff: Thing and stuff classes in context. In *CVPR*, 2018. 8, 24

[126] Tsung-Yi Lin, Michael Maire, Serge J. Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C. Lawrence Zitnick. Microsoft coco: Common objects in context. In *ECCV*, 2014. 8, 11, 23, 24, 25, 26

[127] Navaneeth Bodla, Bharat Singh, Rama Chellappa, and Larry S. Davis. Soft-nms — improving object detection with one line of code. In *ICCV*, 2017. 8

[128] Joao Carreira and Andrew Zisserman. Quo vadis, action recognition? a new model and the kinetics dataset. In *CVPR*, 2017. 8, 23, 24

[129] Taojiannan Yang, Yi Zhu, Yusheng Xie, Aston Zhang, Chen Chen, and Mu Li. Aim: Adapting image models for efficient video action recognition. In *ICLR*, 2023. 8, 24

[130] A. Sophia Koepke, Andreea-Maria Oncescu, João F. Henriques, Zeynep Akata, and Samuel Albanie. Audio retrieval with natural language queries: A benchmark study. *arXiv:2112.09418*, 2021. 9

[131] Xinhao Mei, Xubo Liu, Jianyuan Sun, MarkD . Plumbley, and Wenwu Wang. On metric learning for audio-text cross-modal retrieval. In *Interspeech*, 2022. 9

[132] Soham Deshmukh, Benjamin Elizalde, and Huaming Wang. Audio retrieval with wavtext5k and clap training. *arXiv:2209.14275*, 2022. 9, 23

[133] Yifei Xin, Dongchao Yang, and Yuexian Zou. Improving text-audio retrieval by text-aware attention pooling and prior matrix revised loss. *arXiv:2303.05681*, 2023. 9

[134] Khaled Koutini, Jan Schlüter, Hamid Eghbal-zadeh, and Gerhard Widmer. Efficient training of audio transformers with patchout. *arXiv:2110.05069*, 2021. 9

[135] Yuan Gong, Andrew Rouditchenko, Alexander H. Liu, David F. Harwath, Leonid Karlinsky, Hilde Kuehne, and James R. Glass. Contrastive audio-visual masked autoencoder. *arXiv:2210.07839*, 2022. 9

[136] Pinci Yang, Xin Wang, Xuguang Duan, Hong Chen, Runze Hou, Cong Jin, and Wenwu Zhu. Avqa: A dataset for audio-visual question answering on videos. In *ACMMM*, 2022. 9, 25

[137] Lewei Yao, Runhui Huang, Lu Hou, Guansong Lu, Minzhe Niu, Hang Xu, Xiaodan Liang, Zhenguo Li, Xin Jiang, and Chunjing Xu. Filip: Fine-grained interactive language-image pre-training. In *ICLR*, 2022. 9, 10

[138] Chris Dongjoo Kim, Byeongchang Kim, Hyunmin Lee, and Gunhee Kim. Audiocaps: Generating captions for audios in the wild. In *NAACL*, 2019. 9, 23, 24

[139] Konstantinos Drossos, Samuel Lipping, and Tuomas Virtanen. Clotho: an audio captioning dataset. In *ICASSP*, 2019. 9, 23, 24

[140] Junke Wang, Dongdong Chen, Zuxuan Wu, Chong Luo, Luowei Zhou, Yucheng Zhao, Yujia Xie, Ce Liu, Yu-Gang Jiang, and Lu Yuan. Omnivl: One foundation model for image-language and video-language tasks. *arXiv:2209.07526*, 2022. 10

[141] Zhe Gan, Yen-Chun Chen, Linjie Li, Chen Zhu, Yu Cheng, and Jingjing Liu. Large-scale adversarial training for vision-and-language representation learning. In *NeurIPS*, 2020. 10

[142] Aishwarya Kamath, Mannat Singh, Yann LeCun, Ishan Misra, Gabriel Synnaeve, and Nicolas Carion. Mdetr - modulated detection for end-to-end multi-modal understanding. In *ICCV*, 2021. 10

[143] Zhengyuan Yang, Zhe Gan, Jianfeng Wang, Xiaowei Hu, Faisal Ahmed, Zicheng Liu, Yumao Lu, and Lijuan Wang. Unitab: Unifying text and box outputs for grounded vision-language modeling. In *ECCV*, 2022. 10

[144] Yan Zeng, Xinsong Zhang, and Hang Li. Multi-grained vision language pre-training: Aligning texts with visual concepts. *arXiv:2111.08276*, 2021. 10, 11, 26

[145] Siyi Liu, Zhaoyang Zeng, Tianhe Ren, Feng Li, Hao Zhang, Jie Yang, Chun yue Li, Jianwei Yang, Hang Su, Jun-Juan Zhu, and Lei Zhang. Grounding dino: Marrying dino with grounded pre-training for open-set object detection. *arXiv:2303.05499*, 2023. 10

[146] Licheng Yu, Patrick Poirson, Shan Yang, Alexander C Berg, and Tamara L Berg. Modeling context in referring expressions. In *ECCV*, 2016. 10, 25

[147] Junhua Mao, Jonathan Huang, Alexander Toshev, Oana Camburu, Alan L Yuille, and Kevin Murphy. Generation and comprehension of unambiguous object descriptions. In *CVPR*, 2016. 10, 25

[148] Ranjay Krishna, Yuke Zhu, Oliver Groth, Justin Johnson, Kenji Hata, Joshua Kravitz, Stephanie Chen, Yannis Kalantidis, Li-Jia Li, David A. Shamma, Michael S. Bernstein, and Li Fei-Fei. Visual genome: Connecting language and vision using crowdsourced dense image annotations. In *IJCV*, 2017. 10, 11, 26

[149] Stanislaw Antol, Aishwarya Agrawal, Jiasen Lu, Margaret Mitchell, Dhruv Batra, C Lawrence Zitnick, and Devi Parikh. Vqa: Visual question answering. In *ICCV*, 2015. 11, 25

[150] Alane Suhr, Stephanie Zhou, Iris Zhang, Huajun Bai, and Yoav Artzi. A corpus for reasoning about natural language grounded in photographs. *arXiv:1811.00491*, 2018. 11, 25, 26

[151] Xiaoyi Dong, Yinglin Zheng, Jianmin Bao, Ting Zhang, Dongdong Chen, Hao Yang, Ming Zeng, Weiming Zhang, Lu Yuan, Dong Chen, Fang Wen, and Nenghai Yu. Maskclip: Masked self-distillation advances contrastive language-image pretraining. *arXiv:2208.12262*, 2022. 12

[152] Alexei Baevski, Arun Babu, Wei-Ning Hsu, and Michael Auli. Efficient self-supervised learning with contextualized target representations for vision, speech and language. *arXiv:2212.07525*, 2022. 12

[153] Xingbin Liu, Jinghao Zhou, Tao Kong, Xianming Lin, and Rongrong Ji. Exploring target representations for masked autoencoders. *arXiv:2209.03917*, 2022. 12

[154] Longhui Wei, Lingxi Xie, Wen gang Zhou, Houqiang Li, and Qi Tian. Mvp: Multimodality-guided visual pre-training. *arXiv:2203.05175*, 2022. 12

[155] Xinyu Zhang, Jiahui Chen, Junkun Yuan, Qiang Chen, Jian Wang, Xiaodi Wang, Shumin Han, Xiaokang Chen, Jimin Pi, Kun Yao, Junyu Han, Errui Ding, and Jingdong Wang. Cae v2: Context autoencoder with clip target. *arXiv:2211.09799*, 2022. 12

[156] Rohit Girdhar, Alaaeldin El-Nouby, Zhuang Liu, Mannat Singh, Kalyan Vasudev Alwala, Armand Joulin, and Ishan Misra. Imagebind: One embedding space to bind them all. 2023. 13

[157] Jort F. Gemmeke, Daniel P. W. Ellis, Dylan Freedman, Aren Jansen, Wade Lawrence, R. Channing Moore, Manoj Plakal, and Marvin Ritter. Audio set: An ontology and human-labeled dataset for audio events. In *ICASSP*, 2017. 23

[158] Frederic Font, Gerard Roma, and Xavier Serra. Freesound technical demo. In *ACMMM*, 2013. 23

[159] A. Sophia Koepke, Andreea-Maria Oncescu, João F. Henriques, Zeynep Akata, and Samuel Albanie. Audio retrieval with natural language queries: A benchmark study. In *TMM*, 2022. 23

[160] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. In *ICLR*, 2019. 23

[161] Gao Huang, Yu Sun, Zhuang Liu, Daniel Sedra, and Kilian Q. Weinberger. Deep networks with stochastic depth. In *ECCV*, 2016. 23

[162] Samyam Rajbhandari, Jeff Rasley, Olatunji Ruwase, and Yuxiong He. Zero: Memory optimizations toward training trillion parameter models. In *SC*, 2020. 24

[163] Golnaz Ghiasi, Yin Cui, Aravind Srinivas, Rui Qian, Tsung-Yi Lin, Ekin D Cubuk, Quoc V Le, and Barret Zoph. Simple copy-paste is a strong data augmentation method for instance segmentation. In *CVPR*, 2021. 24

[164] Karol J. Piczak. Esc: Dataset for environmental sound classification. In *ACMM*, 2015. 24, 26

[165] Eduardo Fonseca, Xavier Favory, Jordi Pons, Frederic Font, and Xavier Serra. Fsd50k: An open dataset of human-labeled sound events. In *TASLP*, 2020. 24

[166] Honglie Chen, Weidi Xie, Andrea Vedaldi, and Andrew Zisserman. Vggsound: A large-scale audio-visual dataset. In *ICASSP*, 2020. 24

[167] Peter Young, Alice Lai, Micah Hodosh, and Julia Hockenmaier. From image descriptions to visual denotations: New similarity metrics for semantic inference over event descriptions. In *TACL*, 2014. 25

[168] Andrej Karpathy and Li Fei-Fei. Deep visual-semantic alignments for generating image descriptions. In *CVPR*, 2015. 25

[169] Juho Lee, Yoonho Lee, Jungtaek Kim, Adam R. Kosiorek, Seungjin Choi, and Yee Whye Teh. Set transformer: A framework for attention-based permutation-invariant neural networks. In *ICML*, 2018. 25

[170] Yash Goyal, Tejas Khot, Douglas Summers-Stay, Dhruv Batra, and Devi Parikh. Making the v in vqa matter: Elevating the role of image understanding in visual question answering. In *CVPR*, 2017. 26

# A  Pretraining Details

## A.1  Pretraining Datasets

For image-text pairs, we use LAION-2B [102], a dataset obtained by web crawling that may contain some noisy pairs. To improve the data quality, we apply several pre-processing steps, including removing images with an aspect ratio greater than 3.5, removing images with the shortest side less than 128, and removing images with a CLIP score less than 0.3. We also remove texts containing non-English or emoji characters, as well as texts with lengths less than 3 or greater than 512. After these steps, we retain about 1.5 billion image-text pairs.

For audio-text pairs, we mainly use the environmental sound datasets processed by [76]. Specifically, for some datasets that only contain tags, [76] uses a pretrained language model T5 [97] to rewrite these tags into captions. We also perform simple cleaning on the data, which involves removing samples with text lengths less than 3 or greater than 512, as well as texts containing non-English or emoji characters. Ultimately, we obtain about 2.4 million audio-text pairs, with a total duration of around 8,000 hours. Table 14 presents the environmental sound datasets utilized by ONE-PEACE.

| Dataset | Number of Samples | Duration | T5 Augmentation |
|---|---|---|---|
| Epidemic Sound | 75618 | 220.40 hrs | Yes |
| AudioCaps [138] | 49494 | 135.56 hrs | No |
| AudioSet [157] | 1910918 | 5263.23 hrs | Yes |
| AudioStock | 9552 | 40.31 hrs | No |
| Clotho [139] | 3839 | 23.99 hrs | No |
| FreeSound [158] | 363618 | 2162.10 hrs | No |
| MACS | 3537 | 9.85 hrs | No |
| SoundDescs [159] | 10677 | 331.67 hrs | No |
| WavText5K [132] | 2248 | 17.23 hrs | No |

Table 14: **Statistics on the environmental sound datasets.** All datasets are publicly available.

## A.2  Pretraining Settings

As mentioned in Sec 3.3, the pretraining of ONE-PEACE is divided into two stages: vision-language pretraining and audio-language pretraining.

For vision-language pretraining, we pretrain ONE-PEACE for 200K steps with a batch size of 32768. We use the AdamW [160] optimizer with $(\beta_1, \beta_2) = (0.9, 0.98)$ and $\epsilon = 1e$-8. The peak learning rate is set to $5e - 4$, with a linear warmup of 3000 steps and a cosine decay scheduler. The image resolution is set to $256 \times 256$. The maximum text sequence length is set to 70. For regulation, we use weight decay with 0.05 and disable dropout. We employ drop path [161] with a 0.4 rate.

For audio-language pretraining, we keep the model parameters related to vision and language (e.g., self-attention layers) frozen and only update the parameters that pertain to audio, such as A-Adapter and A-FFN. In this stage, we pretrain ONE-PEACE for 10 epochs with a batch size of 3072. The peak learning rate is set to $2e - 4$, with a linear warmup of 1 epoch and cosine decay scheduler. The maximum audio duration is set to 15s. For audio with a duration of less than 1s, we first repeat the input and then truncate it to 1s. Other hyper-parameters remain the same as vision-language pretraining.

# B  Details of Downstream Tasks

## B.1  Vision Tasks

Here we describe the implementation details of different vision tasks, including image classification [121], semantic segmentation [123], object detection [126], and video action recognition [128]. All detailed hyperparameters are listed in Table 15.

**Image Classification**  We provide the fine-tuning results on ImageNet-1k [121]. Following recent studies in self-supervised learning for computer vision, we use global pooling of all image tokens excluding the class token, and

| Config | ImageNet-21k | ImageNet-1k | ADE20K | COCO | Kinetics 400 |
|---|---|---|---|---|---|
| Optimizer | | | AdamW | | |
| Optimizer momentum | | | $\beta_1, \beta_2 = 0.9, 0.999$ | | |
| Numerical precision | | | `fp16` | | |
| Peak learning rate | 1e-4 | 5e-5 | 1.5e-5 | 1e-4 | 3e-4 |
| Layer-wise lr decay | 0.85 | 0.9 | 0.95 | 0.9 | - |
| Weight decay | 0.05 | 0.05 | 0.05 | 0.1 | 0.05 |
| Batch size | 5120 | 1024 | 16 | 64 | 64 |
| Warmup ratio | 0.375 | 0.2 | 0.0375 | 0.003 | 0.1 |
| Training epochs | 40 | 15 | 30 | 50 | 30 |
| Drop path | 0.4 | 0.4 | 0.5 | 0.6 | 0.4 |
| Image resolution | $256^2$ | $512^2$ | $896^2$ | $1280^2$ | $256^2$ |

Table 15: **Fine-tuning setting for vision tasks.**

append a LayerNorm with a linear layer for classification. To further unleash the potential of ONE-PEACE, we perform intermediate fine-tuning on ImageNet-21k [122]. We set the label smoothing as 0.3 and do not use random erasing, mixup, and cutmix data augmentations. For fine-tuning on ImageNet-1k, we use exponential moving average (EMA) for model parameters and set the EMA decay rate as 0.9998. For intermediate fine-tuning on ImageNet-21k, we do not use EMA. We also use Zero Redundancy Optimizer [162] and set the stage as 1.

**Semantic Segmentation** We provide the fine-tuning results on ADE20k [123]. We use Mask2Former [124] as the segmentation head. We first intermediate fine-tune segmentation head on coco-stuff [125] dataset for 80k steps. The learning rate is set as 2e-5 and the rest hyperparameters are the same as ADE20K shown in Table 15. Then we fine-tune the model on ADE20K. Both experiments use the cosine learning rate decay scheduler.

**Object Detection** We provide the fine-tuning results on COCO [126] with ViTDet [113]. We use large-scale jitter [163] data augmentation and fine-tune for 50 epochs. We use the linear learning rate decay scheduler and decay the learning rate at 44 and 48 epochs respectively.

**Video Action Recognition** To perform video action recognition, following AIM [129], we freeze the parameters of the pre-trained model and add spatial and temporal MLP adapters in each transformer layer. We conduct experiments on Kinetics 400 [128] dataset. Due to the invalid video links, there are many different versions of the K400 dataset and we use the version released on AcademicTorrents. We use the cosine learning decay scheduler and set the backbone learning rate multiplier of 0.1.

## B.2 Audio-(language) Tasks

We describe the implementation details of audio-text retrieval, audio classification, and audio question answering here. All detailed hyperparameters are listed in Table 16.

**Audio-Text Retrieval** We evaluate ONE-PEACE on AudioCaps [138] and Clotho [139] datasets. To get better results, we merge the training set of AudioCaps [138], Clotho [139], and MACS as the fine-tuning dataset. Similar to image-text retrieval, we use A-Branch and L-Branch to extract the features of audio clips and texts respectively, and then calculate the cosine similarity between these features. The recall@k is employed as the evaluation metric.

**Audio Classification** We conduct experiments on three datasets: ESC-50 [164], FSD50K [165], and VG-GSound [166]. ESC-50 is an environmental sound dataset that contains 2000 environmental audio recordings and 50 labels. We directly use the pretrained ONE-PEACE model to perform zero-shot audio classification on ESC-50. Specifically, we use A-Branch to extract audio embeddings from the audio clips and use L-Branch to extract text embeddings from the label names. Then we determine the labels of the audio clips by calculating the similarity between the embeddings. For FSD50K and VGGSound, we input the original audio into the A-Branch and utilize multi-head attention pooling (MAP) to aggregate the features. FSD50K is a multi-label sound event dataset, for which we use BCELoss as the loss function and report the mean average precision on the test set. VGGSound is an audio-visual dataset, where each sample includes a video with audio. We extract the audio clips from the videos and excluded the visual information, using cross entropy as the loss function and reporting accuracy on the test set.

| Config | AudioCaps & Clotho | FSD50K | VGGSound | AQA |
|---|---|---|---|---|
| Optimizer | AdamW | | | |
| Optimizer momentum | $\beta_1, \beta_2 = 0.9, 0.999$ | | | |
| Weight decay | 0.05 | | | |
| Gradient clip | 0.0 | | | |
| Warmup ratio | 0.1 | | | |
| Learning rate schedule | cosine decay | | | |
| Numerical precision | bf16 | | | |
| Peak learning rate | 1.5e-4 | 1e-4 | 8e-5 | 7e-5 |
| Layer-wise lr decay | 0.95 | 0.9 | 0.95 | 0.9 |
| Batch size | 384 | 128 | 512 | 128 |
| Training epochs | 10 | 10 | 10 | 10 |
| Drop path | 0.9 | 0.5 | 0.6 | 0.5 |
| Max duration | 20s | 15s | 15s | 15s |

Table 16: **Fine-tuning setting for audio(-language) tasks.**

**Audio Question Answering** We conduct experiments on the AVQA dataset [136]. Each sample in this dataset consists of a video, a question, and four candidate answers. To perform the audio question answering task, we extract audio clips from the videos and excluded the visual information. During training, we concatenate each answer with the audio and question, and extracted the features through AL-Branch. We then minimize the pairwise hinge loss between the positive features and negative features.

## B.3 Vision-language tasks

Here we describe the implementation details of different vision-language tasks, including image-text retrieval [167, 126], visual grounding [146, 147], visual question answering [149], and visual reasoning [150]. All detailed hyperparameters are listed in Table 17.

| Config | MSCOCO | Flickr30K | RefCOCO/+/g | VQA | NLVR2 |
|---|---|---|---|---|---|
| Optimizer | AdamW | | | | |
| Optimizer momentum | $\beta_1, \beta_2 = 0.9, 0.999$ | | | | |
| Weight decay | 0.05 | | | | |
| Gradient clip | 0.0 | | | | |
| Warmup ratio | 0.1 | | | | |
| Learning rate schedule | cosine decay | | | | |
| Numerical precision | bf16 | | | | |
| Peak learning rate | 8e-5 | 7e-5 | 1.5e-4 | 3e-4 | 1e-4 |
| Layer-wise lr decay | 0.9 | 0.9 | 0.9 | 0.85 | 0.9 |
| Batch size | 3072 | 3072 | 256 | 512 | 128 |
| Training epochs | 15 | 20 | 30 | 10 | 25 |
| Drop path | 0.5 | 0.4 | 0.4 | 0.5 | 0.4 |
| Image resolution | 432 | 432 | 384 | 768 | 256 |

Table 17: **Fine-tuning setting for vision-language tasks.**

**Image-Text Retrieval** We evaluate ONE-PEACE on MSCOCO [126] and Flickr30K [167] datasets, and report the results on the widely used Karpathy test split [168]. We use V-Branch and L-Branch to extract the features of images and texts respectively, and then calculate the cosine similarity between these features. The recall@k is employed as the evaluation metric.

**Visual Grounding** This task requires the model to locate an image region based on a text description. We conduct experiments on RefCOCO, RefCOCO+, and RefCOCOg datasets [146, 147]. The image and text are fed to the VL-Branch simultaneously, then we use multi-head attention pooling (MAP) [169] to aggregate the features from all image

patches. The pooled output is used to predict the continuous corner coordinates of the bounding box $(x_1, y_1, x_2, y_2)$, where $x_1$ and $y_1$ denotes the normalized top left coordinates, $x_2$ and $y_2$ denotes the normalized bottom right coordinates. We report the standard metric Acc@0.5 on the validation and test sets.

**Visual Question Answering**    This task requires the model to answer the question based on an image. We perform experiments on the VQAv2 dataset [170]. Following previous works [31, 53, 35, 144], we use the training and validation set of VQAv2 for training, including additional question-answer pairs from Visual Genome [148]. The image and question are fed to the VL-Branch simultaneously, then we use MAP to aggregate the features from all text tokens. The pooled output is fed into a classifier to predict the answer from the 3,129 most frequent answers. We report the final score on the test-dev and test-std sets.

**Visual Reasoning**    Given a text and a pair of images, this task requires the model to distinguish whether the text truly describes the images. We conduct experiments on the NLVR2 dataset [150]. Following the common practice, We treat each sample as two image-text pairs, each containing a text and one image. Then we input these pairs into VL-branch respectively. The final pooled outputs are concatenated together and fed to a classifier to predict the label. We report accuracy on the dev and test-P sets.

## C    Effects of Pretrained Audio Feature Extractor

We conduct a systematic analysis of the impact of the pretrained audio feature extractor. We find that although the parameters of the feature extractor are only $4.6$M, accounting for only about $1\%$ of the total parameters, it has a significant impact on the model performance. As shown in Table 18, the feature extractor with random initialization only achieves $85.3$ accuracy on the ESC-50 dataset, while using pretrained feature extractors results in better performance. Notably, using the WavLM feature extractor can lead to the largest improvement (+6.2). We attribute this to the fact that WavLM is trained on a more diverse audio dataset compared to Hubert and Wav2Vec 2.0, making its feature extractor more suitable for environmental sound tasks.

| Feature Extractor | Random Init. | Init. with Hubert [13] | Init. with Wav2Vec 2.0 [11] | Init. with WavLM [12] |
|---|---|---|---|---|
| ESC-50 Acc. | 85.8 | 89.6 (+3.8) | 90.0 (+4.2) | 91.8 (+6.0) |

Table 18: **Ablation studies of pretrained audio feature extractors.** We report zero-shot accuracy on the ESC-50 dataset.

## D    Evaluate ONE-PEACE on *One Piece*

We further test the visual grounding ability of ONE-PEACE by using a more complex anime picture, *One Piece*. The model is fine-tuned on the RefCOCOg dataset. As shown in Figure 6, we ask ONE-PEACE to locate the characters based on their names. Although ONE-PEACE hasn't seen any anime pictures in the RefCOCOg dataset, it still achieves a recognition accuracy of 56.6%.

## E    More Examples of Emergent Zero-shot Retrieval

In this section, we provide more examples to demonstrate the emergent zero-shot abilities of ONE-PEACE, including audio-to-image, audio+image-to-image, and audio+text-to-image retrieval. The audios are selected from ESC-50 [164], and the images are retrieve from ImageNet-1K [122] and MSCOCO [126]. By reading this section, we hope that readers can better perceive ONE-PEACE.

Figure 6: Visualization of ONE-PEACE locating different characters of *One Piece*. Given the names of 9 members of the Straw Hat Pirates, ONE-PEACE correctly located 5 of them from the picture.

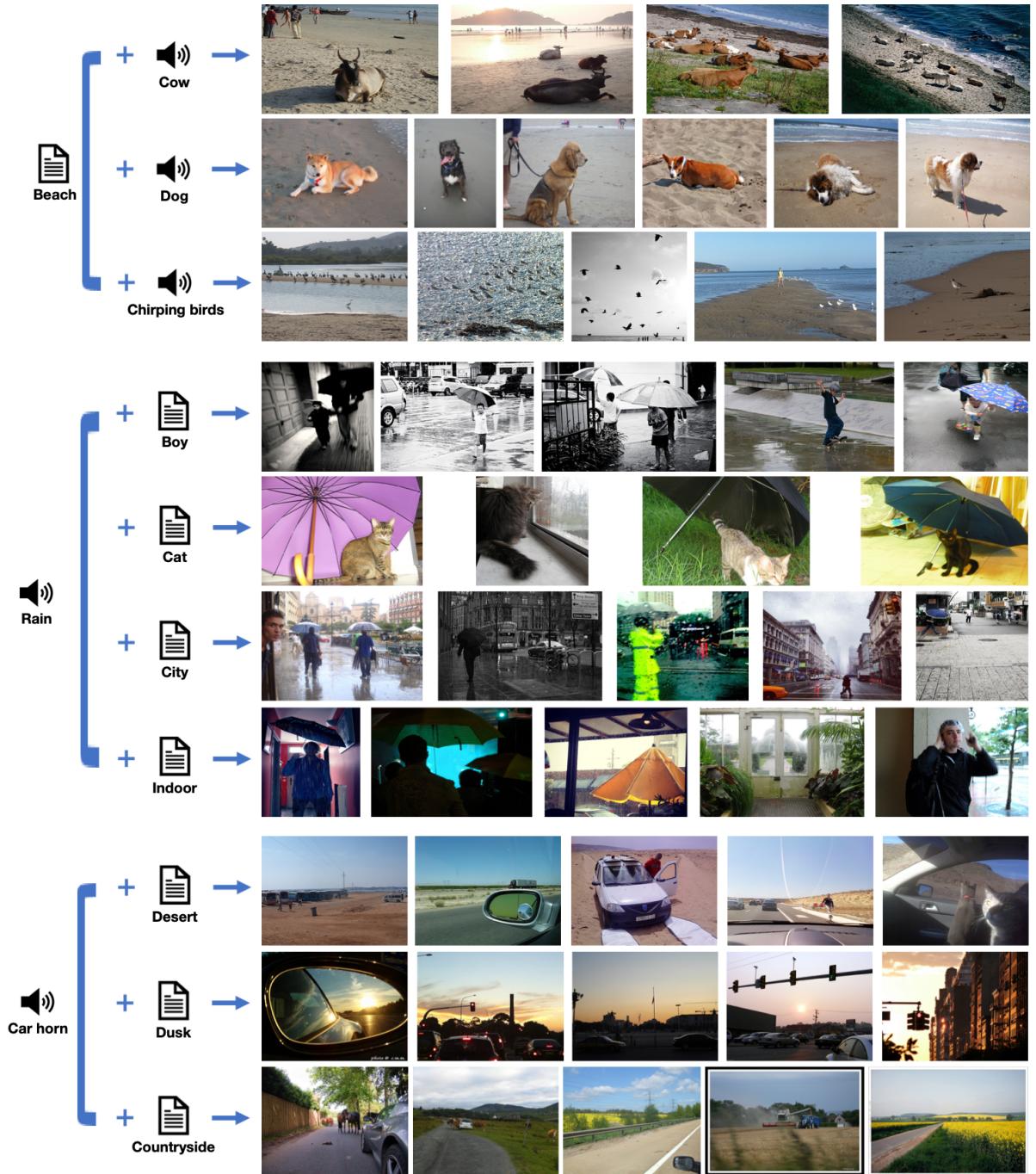Figure 7: **Audio-to-image retrieval.**

Figure 8: **Audio+image-to-image retrieval.**

Figure 9: **Audio+text-to-image retrieval.**