

03/07-06-2023

Topic: One-Peace: Exploring one general representation model toward unlimited modalities

Abstract:

In this work, we explore a scalable way of building a general representation model toward unlimited modalities. We release One-Peace, a highly extensible model with 4B parameters that can seamlessly again 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 pre-train One-Peace, we develop two modality-agnostic pre-training tasks, cross-modal aligning 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 pre-training tasks, One-Peace has the potential to expand to unlimited modalities without using any vision or language pre-trained models.

for Initialization, One Piece 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, USTED), audio classification (ESC-50, FSD501, VGG-Sound), audio question answering (AVQA), image-text retrieval (MSCOCO, Flickr30K) and visual grounding (RefCOCO1+g).

Conclusion :

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 Piece, a general representation model that can seamlessly align and integrate representations across vision, audio and language modalities. We conduct a series of experiments across 8 modalities, 11 tasks and 16 datasets. The experimental results demonstrate that One Piece 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 Piece possesses a strong

Date: / /

emergent zero-shot retrieval capability, enabling it to align modalities that are not paired in the training data.