

CAVMA reading notes

1. Preparation

(1) audio and image pre-processing and tokenization

- dataset: 10s videos with parallel audios in AudioSet and VGGSound

- audio: AST

10s audio wave \longrightarrow 1024 (time) \times 128 (frequency) spectrogram
 \longrightarrow 512 16×16 square patches $\vec{a} = [a^1, \dots, a^{512}]$

- video: ViT

frame aggregation strategy: save computation resources

sample 10 RGB frames (1 FPS) $\begin{cases} \xrightarrow{\text{training}} \text{randomly select 1 RGB frame as input} \\ \xrightarrow{\text{inference}} \text{average each RGB frame prediction as the video prediction} \end{cases}$

each RGB frame resize + center crop \longrightarrow 224×224
 \longrightarrow 196 16×16 square patches $\vec{v} = [v_1, \dots, v_{196}]$

(2) transformer architecture

standard Transformer

a Transformer layer:

$$\left. \begin{aligned} x' &= \text{MSA}(\text{LN}_1(x)) + x \\ y &= \text{MLP}(\text{LN}_2(x')) + x' \end{aligned} \right\} y = \text{Transformer}(x; \text{MSA}, \text{LN}_1, \text{LN}_2, \text{MLP})$$

MSA: multi-headed self attention

LN: layer normalization

MLP: multilayer perceptron

(3) contrastive audio-visual learning (CAV)

$$\text{audio} \longrightarrow \vec{a}_i \longrightarrow c_i^a = \text{MeanPool}(E_a(\text{Proj}_a(a_i)))$$

$$\text{video} \longrightarrow \vec{v}_i \longrightarrow c_i^v = \text{MeanPool}(E_v(\text{Proj}_v(v_i)))$$

$$\text{Proj}_a: \vec{a}_i \rightarrow \mathbb{R}^{768}$$

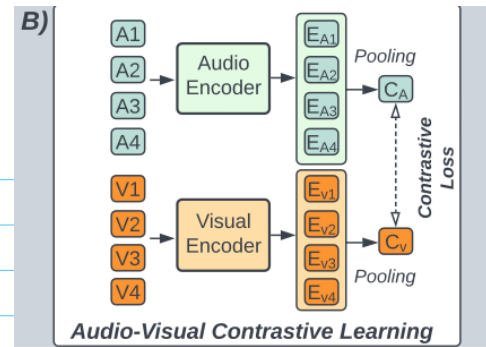
$$\text{Proj}_v: \vec{v}_i \rightarrow \mathbb{R}^{768}$$

contrastive loss:

$$L_c = -\frac{1}{N} \sum_{i=1}^N \log \left[\frac{\exp(s_{i,i}/\tau)}{\sum_{k \neq i} \exp(s_{i,k}/\tau) + \exp(s_{i,i}/\tau)} \right]$$

$$s_{i,j} = \|\mathbf{c}_i^v\|^T \|\mathbf{c}_j^a\|$$

τ is the temperature.



(4) single modality masked autoencoder (MAE)

$$\vec{x} = [x_1, x_2, \dots, x^n]$$

$$\vec{x} / x_{\text{mask}} \xrightarrow{\text{encoder-decoder model}} \hat{x}_{\text{mask}} \rightarrow \text{MSE}(x_{\text{mask}}, \hat{x}_{\text{mask}})$$

(5) vanilla audio-visual masked autoencoder (AV-MAE)

$$a' = \text{Proj}_a(a) + E_a + E_a^p$$

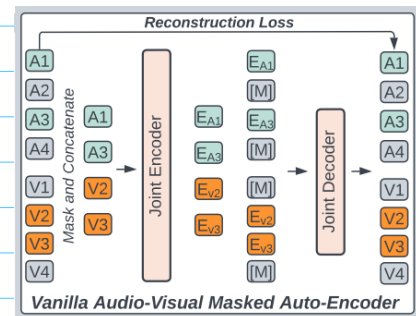
modality type embedding

2-D sinusoid positional embedding

$$v' = \text{Proj}_v(v) + E_v + E_v^p$$

$$x = [a', v']$$

mask a portion (75%) of x



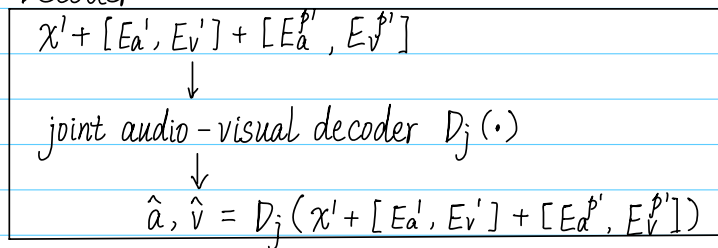
Encoder

$$x_{\text{unmask}} = x \setminus x_{\text{mask}}$$

audio-visual joint encoder $E_j(\cdot)$

$$x_{\text{unmask}} \xrightarrow{\text{pad with masked token}} x'$$

Decoder



minimize MSE \hat{a}, \hat{v} and normalized a, v

2. Contrastive audio-visual masked auto-encoder (CAV-MAE)

combine CAV and AV-MAE

mini-batch of N audio-visual pair samples

\downarrow preprocess

$\{a_i, v_i\}, i = 1, \dots, N$

$$a_i^{unmask} = \text{Mask}_{0.75}(\text{Proj}_a(a_i) + E_a + E_a^p)$$

$$v_i^{unmask} = \text{Mask}_{0.75}(\text{Proj}_v(v_i) + E_v + E_v^p)$$

$$a_i^{unmask} \xrightarrow{E_a} a_i'$$

$$v_i^{unmask} \xrightarrow{E_v} v_i'$$

multi-stream forward pass:

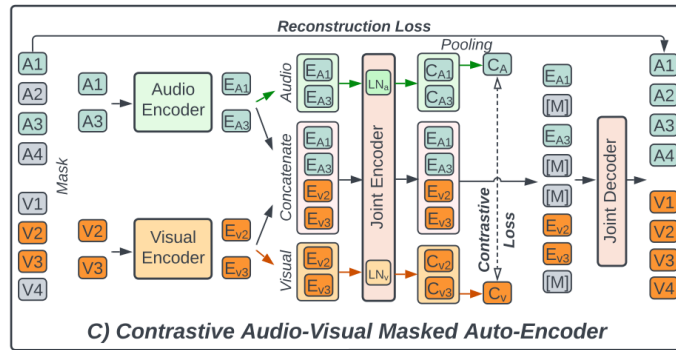
$$c_i^a = \text{MeanPool}(E_j(E_a(a_i^{unmask})); \text{LN}|_a, \text{LN}2_a)) \rightarrow \text{for CL}$$

$$c_i^v = \text{MeanPool}(E_j(E_v(v_i^{unmask})); \text{LN}|_v, \text{LN}2_v)) \rightarrow \text{for CL}$$

$$\chi_i = E_j([E_a(a_i^{unmask}), E_v(v_i^{unmask})]; \text{LN}|_{av}, \text{LN}2_{av}) \rightarrow \text{for RC}$$

* shared weights for E_j

* different LN layers $\text{LN}|_a, \text{LN}|_v, \text{LN}|_{av}$



$$\hat{a}, \hat{v} = D_j(\chi' + [E_a', E_v'] + [E_a^p, E_v^p])$$

$$\text{Loss: } L_c = -\frac{1}{N} \sum_{i=1}^N \log \left[\frac{\exp(s_{i,i}/\tau)}{\sum_{k \neq i} \exp(s_{i,k}/\tau) + \exp(s_{i,i}/\tau)} \right]$$

$$L_r = \frac{1}{N} \sum_{i=1}^n \left[\frac{\sum (\hat{a}_i^{\text{mask}} - \text{norm}(a_i^{\text{mask}}))^2}{|a_i^{\text{mask}}|} + \frac{\sum (\hat{v}_i^{\text{mask}} - \text{norm}(v_i^{\text{mask}}))^2}{|v_i^{\text{mask}}|} \right]$$

number of masked patches

$$L_{\text{CAV-MAE}} = L_r + \lambda_c L_c$$

* only keep the encoders for downstream tasks.

- ① single-modality stream output + multi-modal stream output
 - ② multi-modal stream output
- ① and ② perform similarly