

# MultiModal learning

---

Kyungseon Lee

September 6, 2024

Seoul National University

# Outline

- 1 Introduction
- 2 AVDCNN: Audio-visual speech enhancement using multimodal deep convolutional neural networks - 2017 IEEE
- 3 Stable Diffusion: High-Resolution Image Synthesis with Latent Diffusion Models - 2022 CVPR
- 4 CLIP: Learning Transferable Visual Models From Natural Language Supervision - 2021 ICML
- 5 Conclusion

# Introduction

---

- ① What is multimodal learning?
  - ▶ Multimodal learning is a method that takes different types of data, such as text, images, and audio, as inputs and learns from them simultaneously.
- ② What is the goal of multimodal learning?
  - ▶ The goal is to maximize the joint likelihood of the data observed from different modalities.
  - ▶ To do this, we need to convert the data from each modality into a common latent space.

## How the data from each modality is combined?

### ① Early Fusion

$$\blacktriangleright \hat{y} = f(x_1 \oplus x_2 \oplus \dots \oplus x_n)$$

### ② Late Fusion

$$\blacktriangleright \hat{y}_i = f_i(x_i), \quad i = 1, 2, \dots, n$$

$$\blacktriangleright \hat{y}_{\text{final}} = \sum_{i=1}^n w_i \hat{y}_i$$

### ③ Joint Fusion

$$\blacktriangleright z_{\text{joint}} = f_1(x_1) \oplus \dots \oplus f_n(x_n); \quad f_i \text{ is updated only by the loss from } x_i$$

$$\blacktriangleright \hat{y} = g(z_{\text{joint}})$$

- $n$ : of modality,  $\oplus$ : concatenation symbol
- $x_i$ : embedding input vector for the  $i$ -th modality,  $\hat{y}$ : output
- $w_i$ : the ensemble weight
- $f(\cdot), g(\cdot), h(\cdot)$ : Neural network models

**AVDCNN: Audio-visual speech  
enhancement using multimodal  
deep convolutional neural  
networks - 2017 IEEE**

---

## Joint Fusion method

- Objective : Noise reduction in speech data.
- Solution : Use lip movement images as extra data to enhance the speech data.

**Solution :** Use lip movement images as extra data to enhance the speech data.

- ①  $Z_{\text{joint}} = f_A(X_A) \oplus f_V(X_V)$
- ②  $\hat{Y} = g_A \circ g_{\text{common}}(Z_{\text{joint}})$  ; enhanced speech data
- ③  $\hat{X}_V = g_V \circ g_{\text{common}}(Z_{\text{joint}})$  ; reconstructed image data

$$\mathcal{L} = \min_{\theta} \left( \frac{1}{K} \sum_{i=1}^K \left\| \hat{Y}_i - Y_i \right\|_2^2 + \mu \left\| \hat{X}_{V,i} - X_{V,i} \right\|_2^2 \right)$$

- $\oplus$  : Concatenation symbol,  $\mu$ : weight hyperparameter.
- $X_A, X_V$  : An embedding vector of noisy audio and visual data.
- $Y$ : The clean audio data.
- $f_A, f_V$ : The feature extraction encoder functions for audio and visual data.
- $g_A, g_V$ : The reconstruction decoder functions for audio and visual data.



# Stable Diffusion: High-Resolution Image Synthesis with Latent Diffusion Models - 2022 CVPR

---

## Conditional method

- Objective: Generate the desired image each time.
- Solution : Use a different modality as a condition for image generation.

# Diffusion Model

## ① Forward process

$$\blacktriangleright q(x_t | x_{t-1}) = \mathcal{N}(x_t; \sqrt{1 - \beta_t}x_{t-1}, \beta_t I_D)$$

## ② Reverse process

$$\blacktriangleright p_\theta(x_{t-1} | x_t) = \mathcal{N}(x_{t-1}; \mu_\theta(x_t, t), \sigma_t^2 I_D)$$

## ③ Objective function

$$\mathcal{L}_{DM} = \mathbb{E}_{x, \epsilon \sim \mathcal{N}(0,1), t} \left[ \|\epsilon - \epsilon_\theta(x_t, t)\|_2^2 \right]$$

with  $t$  uniformly sampled from  $\{1, \dots, T\}$ .

- $D$  : Dimension of data.
- $x_0$  : data  $\in \mathbb{R}^D$ ,  $x_t$  :  $x$  with noise added  $t$  times
- $q(x_0)$  : True density function of  $x_0$ .
- $\beta_t$  : small positive hyperparameter
- $\epsilon, \epsilon_\theta$  : noise and estimated noise

# Stable Diffusion

**Solution :** Use a different modality as a condition for image generation.

- ① input  $z_0 = f^E(x_0)$  ; Feature extraction function using an AE.
- ② In the denoising process of diffusion, we approximate true reverse process using the Transformer.

$$p_{\theta}(z_{t-1} | z_t, x_{other}) = \mathcal{N}(z_{t-1}; \mu_{\theta}(z_t, t, x_{other}), \sigma_t^2 |_D)$$

$$\mathcal{L}_{SD} = \mathbb{E}_{x, x_{other}, \epsilon \sim \mathcal{N}(0,1), t} \left[ \|\epsilon - \epsilon_{\theta}(x_t, x_{other}, t)\|_2^2 \right]$$

- $z_T$ :  $z$  after adding noise  $T$  times to become  $N(0, I_p)$  noise.
- $\tilde{z}$ : The output  $z$  generated conditionally.
- $x_{other}$ : the embedding vector of the other modality data as the condition.
- $\mu_{\theta}(\cdot)$ : A prediction function that calculates the average, using a Transformer.

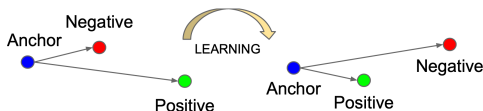
# CLIP: Learning Transferable Visual Models From Natural Language Supervision - 2021 ICML

---

## Self-supervised method

- Objective: Remove the data labeling process.
- Solution : Use the text attached to the image as input instead of a label.

## Contrastive learning



**Figure 1:** Basic contrastive learning process

Basic Contrastive loss:  $L(x_p, x_q, y) = y * ||f(x_p) - f(x_q)||^2$

- ①  $x_p, x_q$  : input pair,  $y$  : output,  $f(.)$  : contrastive learning model
- ② In the representation space, if an input pair is positive, it is mapped close together, and if it is negative, it is mapped farther apart.
- ③ This might look similar to clustering, but the goal of contrastive learning is to better learn representations through the input pairs.

**Solution :** Use the text attached to the image as input instead of a label.

- ①  $l_i = f_I(X_{I,i}), T_i = f_T(X_{T,i})$
- ②  $(l_i, T_i)$  ; the correct pair of the image and text data
- ③  $(l_i, T_j), i \neq j$  ; the uncorrect pair of the image and text data

$$\mathcal{L}_{contrast}(l, T) = \sum_i \left[ -\log\left(\frac{\exp(l_i \cdot T_i)}{\sum_j \exp(l_j \cdot T_i)}\right) - \log\left(\frac{\exp(l_i \cdot T_i)}{\sum_k \exp(l_i \cdot T_k)}\right) \right]$$

$$\hat{y}_i = \left[ \frac{\exp(l_i \cdot T_j)}{\sum_k \exp(l_i \cdot T_k)} \right], j = 1, 2, \dots, n$$

- $X_{I,i}, X_{T,i}$ : input of the image data and text data
- $l_i$ : The embedding vector of the i-th image
- $T_i$ : The embedding vector of the i-th text
- $f_I, f_T$ : The embedding functions for the image data and text data



## Conclusion

---

## Comments

- In multimodal learning, different papers use various methods to fuse different modalities. Among these, the Joint fusion method is effective at both extracting features from each modality and learning interactions between them.
- Multimodal learning also allows models to be trained without data labeling.
- Additionally, it can be applied in areas like harmful word filtering, human behavior recognition, data augmentation, and noise reduction.