# VLM 끄적끄적 2

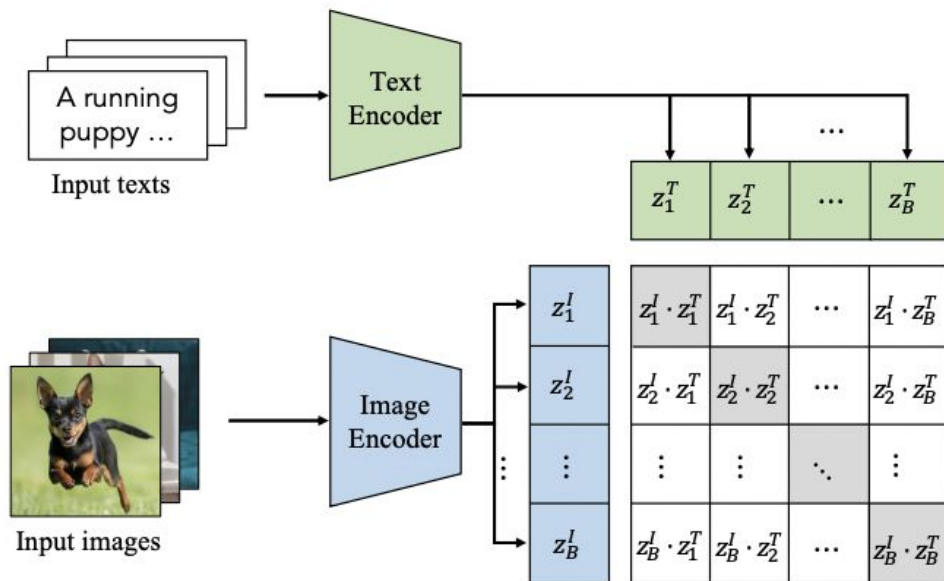Pretraining Task (1): Contrastive Learning

# 1. CLIP: Image-Text CL



Fig. 6: Illustration of the image-text contrastive learning in CLIP [10]. Figure is reproduced from [10].

# 2. CLIP 발전시키는 방향 (Image-Text CL)

1. Large-scale datasets: "데이터셋을 더 키우자"

2. Small-scale datasets: "데이터셋을 (알짜배기만) 효율적으로 사용하자"

3. Various semantic level: "not only 거시적(coarse) but also 미시적 (dense)"

4. Augmentation

5. Unifying vision & language encoder

# 2. CLIP 발전시키는 방향 (Image-Text CL)

1. Large-scale datasets
- ALIGN (ICML 2021),

2. Small-scale datasets
- DeCLIP (ICLR 2022), OTTER (ICLR 2022), ZeroVL (ECCV 2022)

3. Various semantic levels
- FILIP (ICLR 2022), PyramidCLIP (NeurIPS 2022)

4. Augmentation
- LaCLIP (NeurIPS 2023), ALIP (ICCV 2023), RA-CLIP (CVPR 2023)

5. Unifying vision & language encoder
- CLIPPO (CVPR 2023), OneR (AAAI 2023)

# 2. CLIP 발전시키는 방향 (Image-Text CL)

## (1) Large-scale datasets

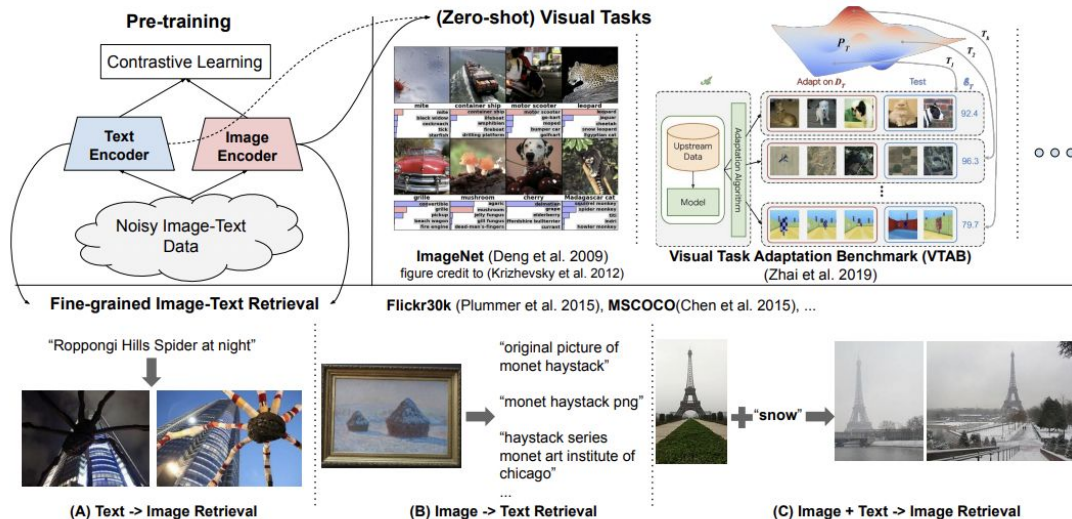- ALIGN (ICML 2021)
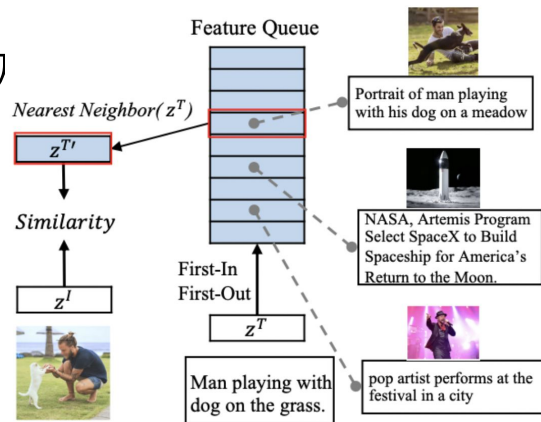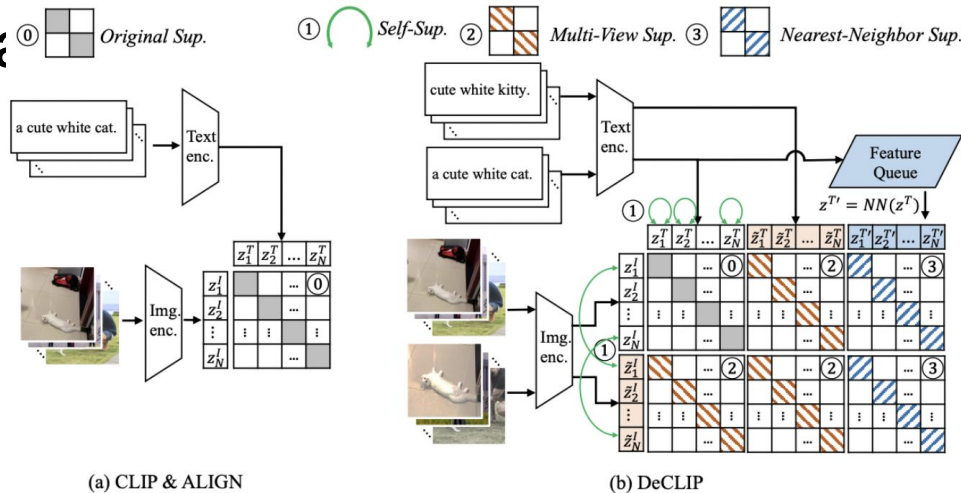
  - Noisy pair도 활용

  - 다양한 task



Figure 1. A summary of our method, ALIGN. Visual and language representations are jointly learned from noisy image alt-text data. The representations can be used for vision-only or vision-language task transfer. Without any fine-tuning, ALIGN powers zero-shot visual classification and cross-modal search including image-to-text search, text-to-image search and even search with joint image+text queries.

# 2. CLIP 발전시키는 방향 (Ima

## (2) Small-scale dataset

- DeCLIP (ICLR 2022)

    - CLIP 한계점: Data-hungry

    - 해결책: Data-efficient CLIP!

    - How? Image-text pair의 supervision을 활용하ㄱ

    - Details:
        - (1) Self-supervision **within** each modality
        - (2) Multi-view supervision **across** modalities
        - (3) NN supervision from other similar pairs



(a) CLIP & ALIGN

(b) DeCLIP

# 2. CLIP 발전시키는 방향 (Image-Text CL)

## (2) Small-scale dataset

- OTTER (ICLR 2022)

    - CLIP 한계점: Data-hungry (이유: Noisy pairs)
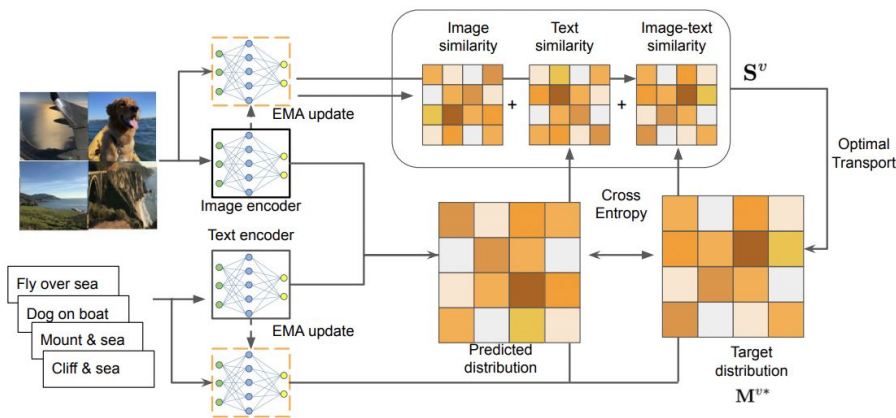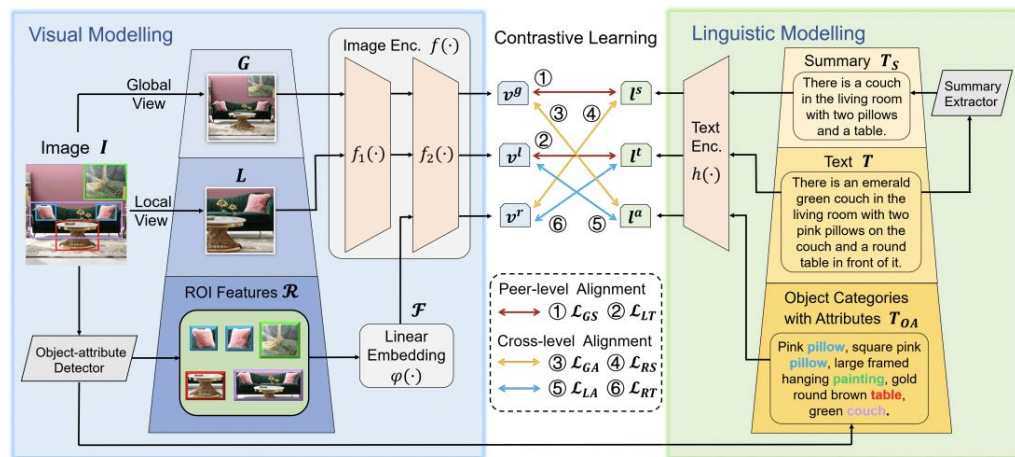
    - 해결책: 3M장만으로 만들기

- ZeroVL (ECCV 2022)



Figure 2: Architecture of OTTER. We use image and text embeddings to compute similarity matrices $S^v$ (and $S^t$), which is then used to solve for matching probabilities $M^{v*}$ (and $M^{t*}$) as targets.

# 2. CLIP 발전시키는 방향 (Image-Text CL)

## (3) Various Semantic Level

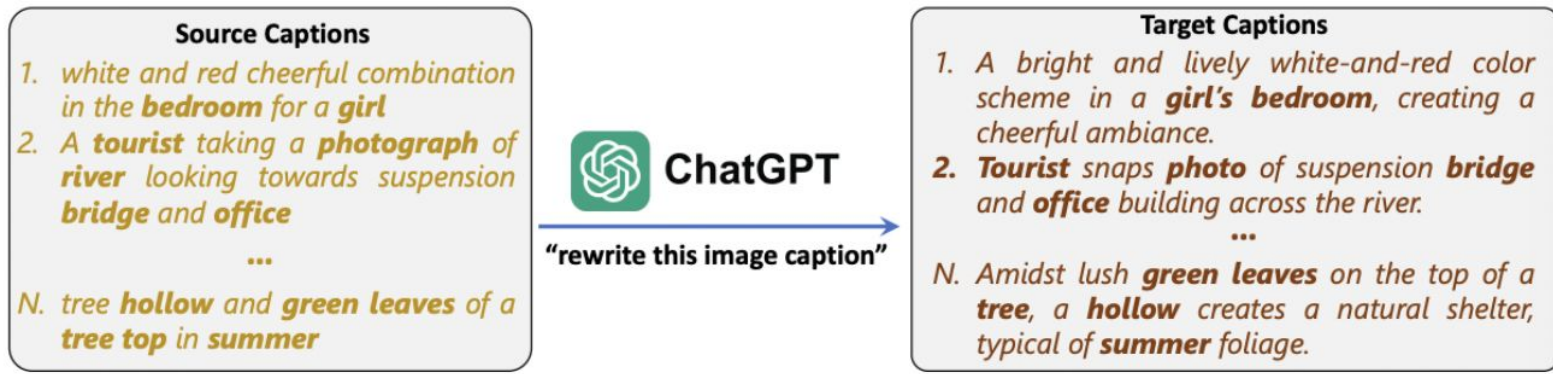- PyramidCLIP (NeurIPS 2022)

  - (1) Cross-level

  - (2) Peer-level



**Figure 2: Overall architecture of the proposed PyramidCLIP which is a dual-stream network.** The input elements of visual modelling and linguistic modelling both have three-level semantics. The elements of the two modalities interact through peer-level semantics alignment and cross-level relation alignment.

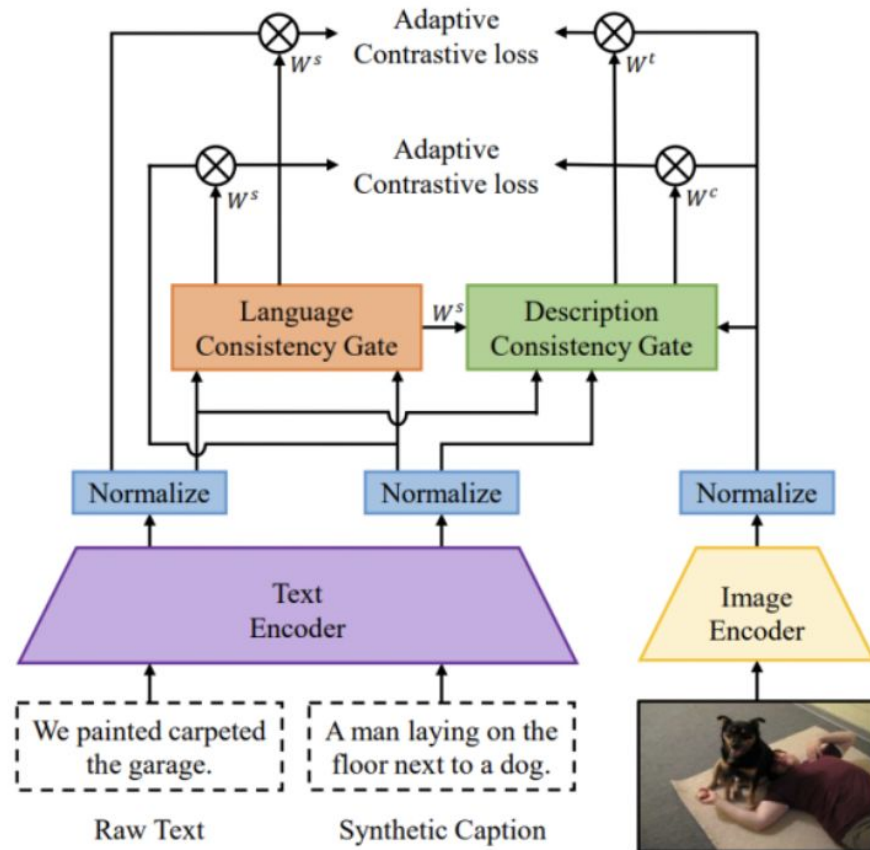# 2. CLIP 발전시키는 방향 (Image-Text CL)

(4) Augmentation

- LaCLIP (NeurIPS 2023)

    - La = Language augmented

    - Solution: Synthetic caption via LLM => Exhibit diversity

# 2. CLIP 발전시키는 방향 (Image-Text CL)

## (4) Augmentation

- ALIP (ICCV 2023)

  - ALIP = Adaptive LIP

  - CLIP 한계점: Noisy data

    ( = Unmatched image & text )

  - Solution: Use both….

    - (1) Raw text

    - (2) Synthetic caption

# 2. CLIP 발전시키는 방향 (Image-Text CL)

(4) Augmentation

- RA-CLIP (CVPR 2023)

    - RA = Retrieval Augmented

    - CLIP 한계점: Data-hungry

    - Solution: Augment embeddings by "online" retrieval

        - Step 1) Use "hold-out" reference set

        - Step 2) Given input, retrieve relevant pairs from reference set

# 2. CLIP 발전시키는 방향 (Image-Text CL)

## (4) Augmentation
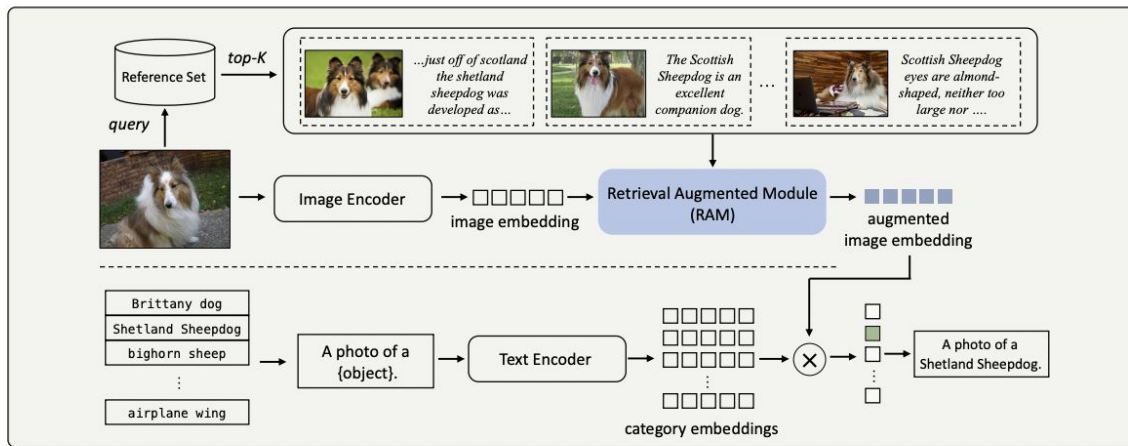
- RA-CLIP (CVPR 2023)



Figure 2. Overview of the proposed RA-CLIP. Given an input image, RA-CLIP retrieves K most similar image-text pairs from the reference set, which provide informative descriptions for the input image, thus we feed them into the Retrieval Augmented Module (RAM) to enrich the representation of input image.

# 2. CLIP 발전시키는 방향 (Image

(5) Unifying vision & language encoder
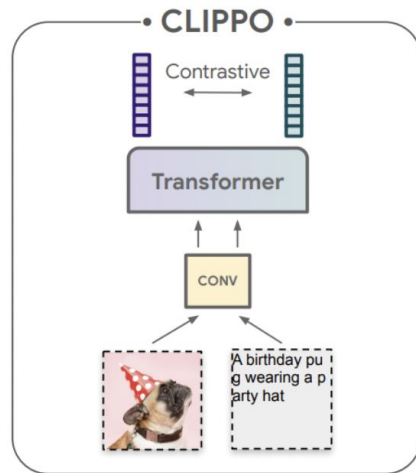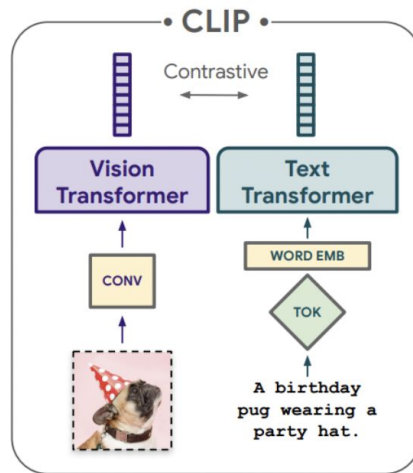


- CLIPPO (CVPR 2023)

  - CLIP-Pixels Only

    (text & image 모두 "픽셀" 취급)

  - CLIP의 한계점: Independent text & image towers

  - 목적: data modality간의 efficient communication을 위해!

  - 해결책: By rendering the alt-text as image!

# 2. CLIP 발전시키는 방향 (Image-Text CL)

(5) Unifying vision & language encoder

| Method | Formulation | |
|---|---|---|
| ITC | $\mathcal{F}(I)$ | $\mathcal{F}(T)$ |
| XMC | $(\mathcal{F}(I) + \mathcal{F}(T))/2$ | $(\mathcal{F}(I) + \mathcal{F}(T))/2$ |
| CIC | $(\mathcal{F}(I\|T) + \mathcal{F}(T\|I))/2$ | $(\mathcal{F}(I) + \mathcal{F}(T))/2$ |
| CMC | $\mathcal{F}(I,T\|I,T)$ | $(\mathcal{F}(I) + \mathcal{F}(T))/2$ |

Table 2: Summary of the contrastive objectives.

- OneR (AAAI 2023)

    - One-tower Model

    - 기본 가정: Image & Text는 같은 본질이 다른 식으로 발현된 것일 뿐

    - 네 가지 Task

        - Image-text Contrastive (ITC)

        - Cross-modal Mixup Contrastive (XMC)

        - Contextual Invariance Contrastive (CIC)

        - Contextual Mixup Contrast (CMC)

| Imagenet 0-shot | Top-1 Acc. | Top-5 Acc. |
|---|---|---|
| ITC | 1.65 | 5.25 |
| ITC (two heads) | 17.46 | 35.32 |
| ITC + XMC | 22.12 | 42.12 |
| ITC + XMC + CIC | 22.86 | 42.88 |
| ITC + CMC | **23.70** | **43.15** |

# 2. CLIP 발전시키는 방향 (Image-Text

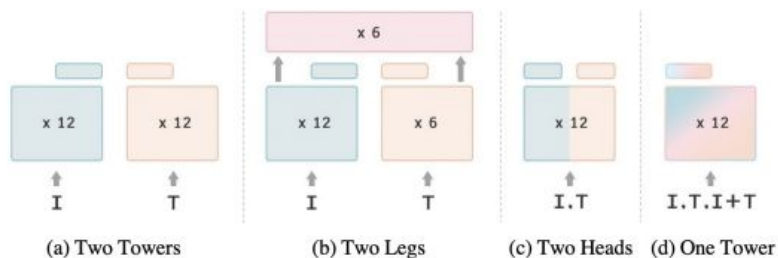## (5) Unifying vision & language encoder

- OneR (AAAI 2023)

Figure 2: Typical architectures of vision-language models. (a) is the basic form, with one transformer encoder and a projector for each modality. (b) adds fusion encoder blocks on top. (c) uses a single transformer encoder, but has separate projections. (d) unifies the two modalities with a generic one-tower model (OneR).
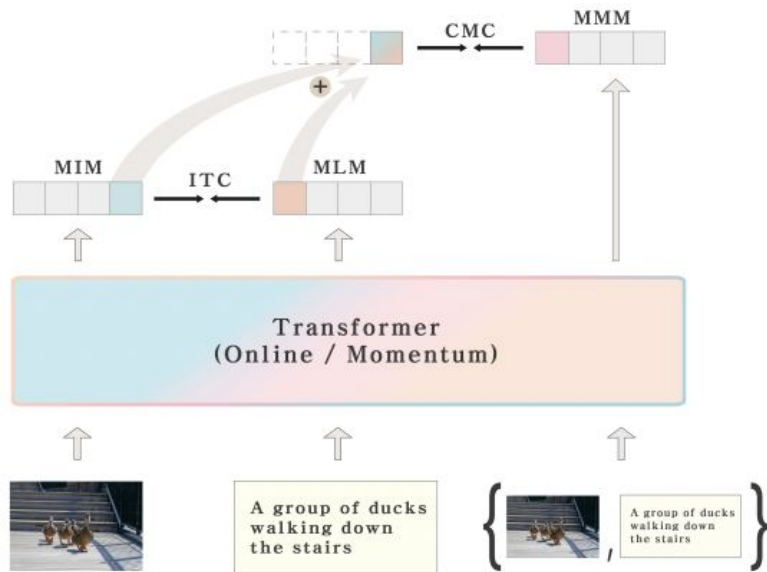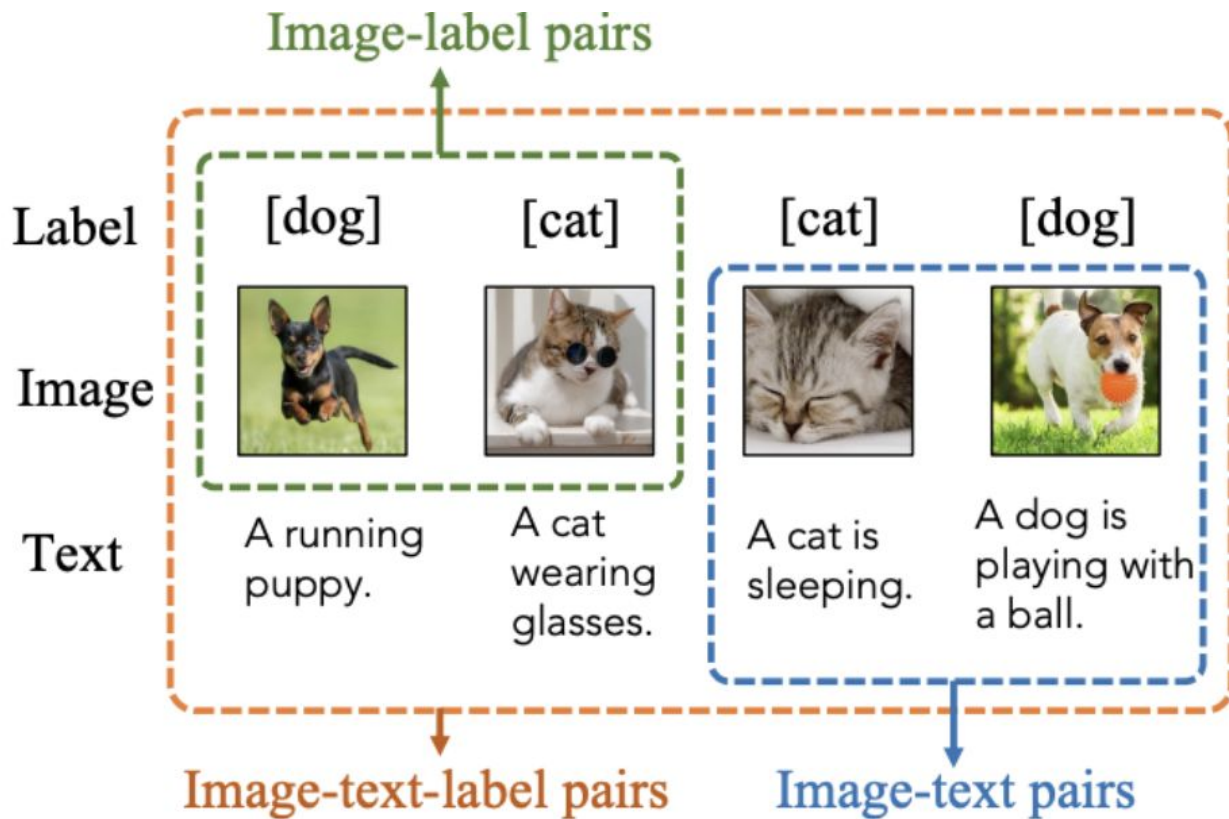
Figure 5: Overview of OneR. Image-text contrastive and contextual mixup contrastive objective provide guidance in parallel with masked modeling for three input types: image, language and multi-modal (image+text).

# 3. CLIP 발전시키는 방향 (Image-Text-Label CL)

# 3. CLIP 발전시키는 방향 (Image-Text-Label CL)
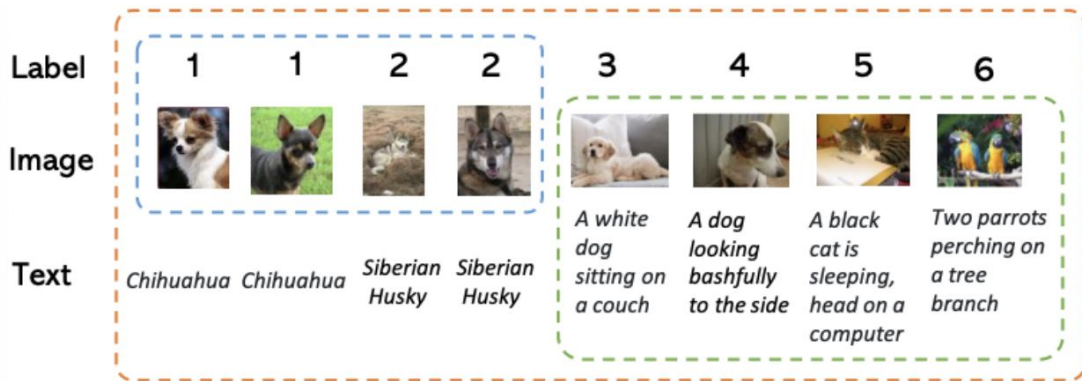
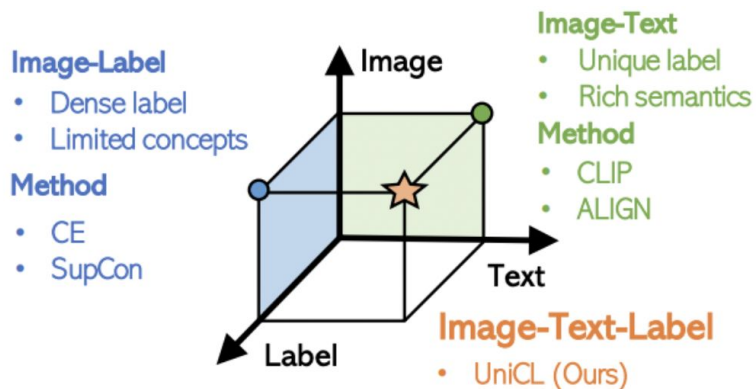## (1) UniCL (CVPR 2022)

- UniCL = Unified CL



Figure 1. Unified contrastive learning paradigm in the **image-text-label** space, which recovers the supervised learning (*e.g.*, Cross-Entropy (CE) [47] or Supervised Contrastive Learning (Sup-Con) [30]) on **image-label** data, and language-image contrastive learning (*e.g.*, CLIP [48] or ALIGN [29]) on **image-text** data.

# 3. CLIP 발전시키는 방향 (Image-Text-Label CL)

## (1) UniCL (CVPR 2022)

- UniCL = Unified CL