

Deep Learning (COSE474)

2024 Fall

Final Project Topics

Hyunwoo J. Kim
(hyunwoojkim@korea.ac.kr)

Agenda

- Rules/Evaluation
- Possible project topics?

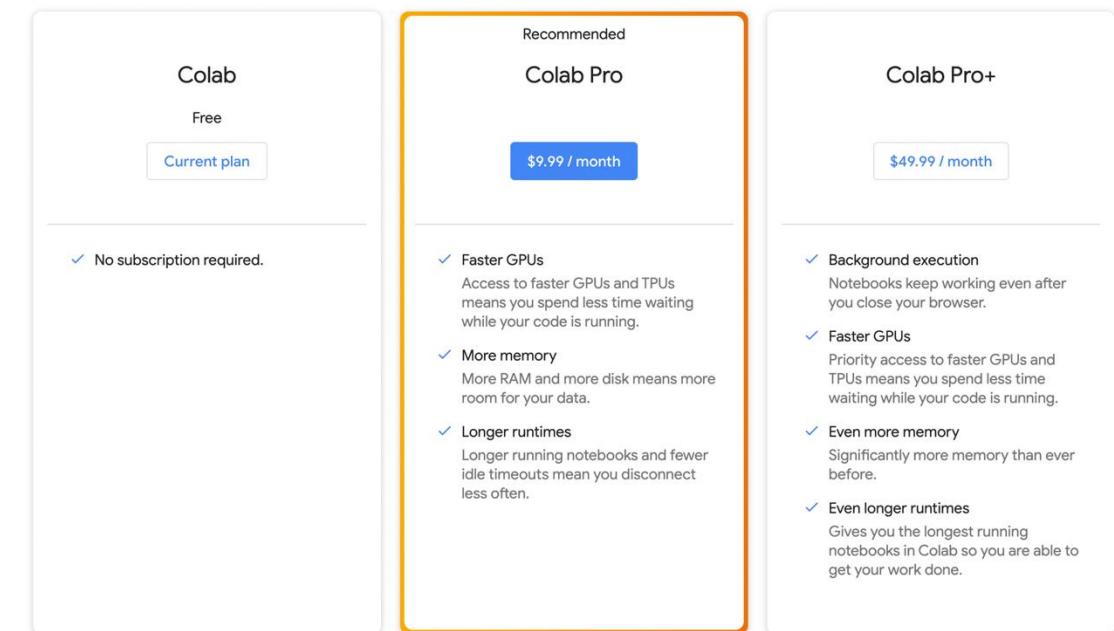
Rule and Evaluation

- Rule? Resources?
 - Computing resource: Colab/AWS/local GPU Servers (default: colab pro)
- Evaluation
 - **Absolute Grading**
 - We do **NOT** expect the state-of-the-art performance
 - We expect you to learn how to design an AI project and develop an AI algorithm.
 - It includes motivations, related works, methods (formulation, architecture), experiments (data preparation, hyperparameter tuning, quantitative/qualitative experimental results), discussion & future direction.

Computing resources

- Colab/Colab Pro/Colab Pro+/your own GPUs/AWS
- Describe your computing resources and reformulate your problem manageable

- Colab Pro \$10/month
- Colab Pro+ \$50/month



Evaluation:report (tentative)

- Evaluation (2024) – KOR or ENG
 - basic points (overall) (2): length (0.5), format (0.5), clarity of writing (1)
 - Introduction (5): motivation (2), problem definition (2), concise description of contribution (1)
 - Methods (5): significance/novelty (2), figure (1), reproducibility (2)-algorithm
 - Experiments (7): dataset (1), computer resource (CPU, GPU, OS, pytorch etc.) & experimental design (1), quantitative results (1), qualitative results (1), Figures (plots)/Tables and their analysis (2), discussion why the proposed method is successful or unsuccessful (1)
 - Future direction (1).
 - Github history (2)
 - Overleaf history (2)
 - ~~(Bonus+1)~~ pre-trained foundation models beyond ImageNet-pretrained CNNs (distillation, adaption, pseudo-labeling, baseline etc.), CLIP, BERT, RoBERTa

Prep Preliminary : 1/24/2024, 051- page 1 2024 X

Good Examples for your final report

- Please check papers from top-tier AI conferences

Graph Transformer Networks

Seongjin Yun, Minhyul Jeong, Raehyun Kim, Jaewoo Kang*, Hyunwoo J. Kim*
Department of Computer Science and Engineering
Korea University
(sys5419, minhyuljeong, raehyun, kangi, hyunwoojkim)@korea.ac.kr

Abstract

Graph neural networks (GNNs) have been widely used in representation learning on graphs and achieved state-of-the-art performance in tasks such as node classification and link prediction. However, most existing GNNs are designed to handle node representations on the *fixed* and *homogeneous* graphs. The limitations especially become problematic when learning representations on a misspecified graph or a *heterogeneous* graph that consists of various types of nodes and edges. In this paper, we propose Graph Transformer Networks (GTNs), which can generate general new graph structures, which involve identifying useful connections between unconnected nodes on the original graph, while learning effective node representation on the new graphs in an end-to-end fashion. Graph Transformer layer, a core layer of GTNs, learns a soft set of edge types to connect relations for generating homogeneous meta-paths on the new graphs. Our experiments show that GTNs learn new graph structures, based on data and tasks without domain knowledge, and learn powerful node representation via convolution on the new graphs. Without domain-specific graph redefining, GTNs achieved the best performance in all three selected node classification tasks, the same as the other methods that require pre-defined meta-paths from domain knowledge.

1 Introduction

In recent years, Graph Neural Networks (GNNs) have been widely adopted in various tasks over graphs, such as graph classification [11, 21, 40], link prediction [18, 30, 42] and node classification [3, 14, 33]. The representation learned by GNNs has been proven to be effective in achieving state-of-the-art performance in a variety of graph tasks, such as social network analysis [14, 19], citation network [19], protein-protein interaction [20], recommendation systems [1, 27, 31]. The underlying graph structure is utilized by GNNs to operate convolution directly on graphs by passing node features [12, 14] to neighbors, or perform convolutions in the spectral domain using the Fourier basis of a given graph, i.e., eigenfunctions of the Laplacian operator [9, 15, 19].

However, one limitation of GNNs is that they are unable to operate GNNs on a noisy graph or a graph with wrong neighbors on the graph. In addition, in some applications, constructing a graph to operate GNNs is challenging due to their small and sparse labeled graph datasets. One prevalent remedy to address this problem is data augmentation. Data augmentation increases the diversity of data and improves the generalization power of machine learning models trained on randomly augmented samples. It is widely used to enhance the generalization ability of models in many domains. For instance, in image recognition, advanced methods like [13–15] as well as simple transformations such as random cropping, cutout, Gaussian noise, or blurring have been used to achieve competitive performance.

However, unlike image recognition, designing effective and label-preserving data augmentation for individual graph datasets is challenging due to their non-localized nature and the dependencies between data samples. In image recognition, it is straightforward to identify operations that preserve labels. For instance, human can verify that rotation, translation, and small color jittering do not change the labels in image classification. In contrast, graphs are less interpretable and it is non-trivial

Metropolis-Hastings Data Augmentation for Graph Neural Networks

Hyeonjin Park^{1*}, Seunghun Lee¹, Sihyeon Kim¹, Jinyoung Park¹
Jisu Jeong^{2,3}, Kyung-Min Kim^{2,3}, Jung-Woo Ha^{2,3}, Hyunwoo J. Kim^{1,†}
Korea University¹, NAVER CLOVA², NAVER AI LAB³
(hyeonjin961030, llsashh319, sh_be15, lpm678, hyunwoojkim)@korea.ac.kr
(jisuj.eoeng, kyungmin.kim.ml, jungwoo.ha)@navercorp.com

Abstract

Graph Neural Networks (GNNs) often suffer from weak-generalization due to sparsely labeled data despite their promising results on various graph-based tasks. Data augmentation is a prevalent remedy to improve the generalization ability of models in many domains. However, due to the non-Euclidean nature of data space and the dependency between nodes, designing effective data augmentation graphs is challenging. In this paper, we propose Metropolis-Hastings Data Augmentation (MH-Aug) that augments graphs from an explicit target distribution for semi-supervised learning. MH-Aug produces a sequence of augmented graphs from the target distribution enables flexible control of the strength and diversity of augmentation. Since the direct sampling from the complex target distribution is challenging, we adopt the Metropolis-Hastings algorithm to obtain the augmented graphs. To further improve the generalization ability of GNNs via self-supervised learning strategy with generated samples from MH-Aug. Our extensive experiments demonstrate that MH-Aug can generate a sequence of samples according to the target distribution to significantly improve the performance of GNNs.

1 Introduction

Graph Neural Networks (GNNs) [1] have been widely used for representation learning on graph-structured data due to their superior performance in various applications such as node classification [2–4], link prediction [5–7] and graph classification [8, 9]. They have been proven effective by achieving impressive performance for diverse datasets such as social networks [10], citation [4], physics [11], molecular [12] and image [13] datasets. However, GNNs often suffer from weak generalization due to their small and sparse labeled graph datasets. One prevalent remedy to address this problem is data augmentation. Data augmentation increases the diversity of data and improves the generalization power of machine learning models trained on randomly augmented samples. It is widely used to enhance the generalization ability of models in many domains. For instance, in image recognition, advanced methods like [13–15] as well as simple transformations such as random cropping, cutout, Gaussian noise, or blurring have been used to achieve competitive performance. However, unlike image recognition, designing effective and label-preserving data augmentation for individual graph datasets is challenging due to their non-localized nature and the dependencies between data samples. In image recognition, it is straightforward to identify operations that preserve labels. For instance, human can verify that rotation, translation, and small color jittering do not change the labels in image classification. In contrast, graphs are less interpretable and it is non-trivial

*First two authors have equal contribution.
†is the corresponding author.

Copyright © 2022, Association for the Advancement of Artificial Intelligence (www.aaai.org). All rights reserved.

The Thirty-Sixth AAAI Conference on Artificial Intelligence (AAAI-22)

Deformable Graph Convolutional Networks

Jinyoung Park, Sungho Yeo, Jibwan Park, Hyunwoo J. Kim*
Department of Computer Science and Engineering, Korea University
(jpmn678, ysd424, jyeven7071, hyunwoojkim)@korea.ac.kr

Abstract

Graph neural networks (GNNs) have significantly improved the representation power for graph-structured data. Despite the success of GNNs, there are still many challenges. Most GNNs have two limitations. Since the graph convolution is performed in a small local neighborhood on the input graph, they are unable to capture long-range dependencies between distant nodes. In addition, a node has neighbors that belong to different classes, i.e., *heterogeneity*. To address these two problems, we propose Deformable Graph Convolutional Networks (Deformable GCNs) that adaptively perform convolution in multiple latent spaces and capture short-long range dependencies between nodes with different features. Our framework simultaneously learns the *node positional embeddings* (coordinates) to determine the relation between nodes. By learning the node position and its position, the convolution kernels are deformed by deformation vector and apply different transformations to its neighborhood. Our proposed Deformable GCNs flexibly handles the heterogeneity and achieves the best performance in node classification tasks on six heterogeneous graph datasets. Our code is publicly available at <https://github.com/kakaobrain/dgcnn>.

Introduction

Graphs are flexible representations for modeling relationships in data analysis problems and are widely used in various domains such as social network analysis (Wang, Cui, and Zhang 2016), recommender system (Berg, Kipf, and Welling 2017), chessboard (Zhou et al. 2018), and image processing (Ezakian and Ramer 2004), etc. Graph convolutional networks (GCNs) have achieved great success in many graph-based learning tasks such as node classification [1], link prediction [Zhang and Chen 2018; Schlichtkrull et al. 2018], and graph classification [Erica et al. 2019; Ying et al. 2018]. In particular, GCNs have shown remarkable representations via message passing schemes, which iteratively learn

*is the corresponding author.
Copyright © 2022, Association for the Advancement of Artificial Intelligence (www.aaai.org). All rights reserved.

74



This CVPR 2021 paper is the Open Access version, provided by the Computer Vision Foundation.
Except for this watermark, it is identical to the accepted version.
The final published version of the proceedings is available on IEEE Xplore.

HO/TR: End-to-End Human-Object Interaction Detection with Transformers

Bumsoo Kim^{1,2} Junhyun Lee² Jaewoo Kang² Eun-Sol Kim^{1,†} Hyunwoo J. Kim^{2,†}
¹Kakao Brain ²Korea University
(bumsoso.brain, eunsol.kim)@kakaobrain.com
(meliketox, ijhyun33, kangj, hyunwoojkim)@korea.ac.kr

Abstract

Human-Object Interaction (HOI) detection is a task of identifying “a set of interactions” in an image, which involves the (i.e., subject (i.e., humans) and target (i.e., objects) of interaction) and the classification of the interactions. Most existing HOI detection methods have indirectly addressed this task by detecting human and object instances and individually inferring every pair of the detected pairs via the association module. However, on heterogeneous graphs where unconnected nodes have both features and different labels, the conventional graph convolutional neural networks often underperform simple methods such as a multi-layer perceptron (MLP) or completely ignore the semantic relationships between nodes. In addition, graph convolution receives messages from local neighbors, which it has the limitation of capturing long-range dependencies, which is critical for the HOI tasks.

To address these limitations, we propose a Deformable Graph Convolutional Network (Deformable GCN) that softly changes the receptive field of each node by adaptively changing the kernel size and applying different convolutions in multiple latent spaces. Started from a general definition of the discrete convolution with finite support, we extend the discrete convolution to 2D convolution (Dai et al. 2017) in a latent space for graph-structured data. Due to the limitation defined on a grid space for images, our convolution kernel generates different convolutions for various relations. This allows us to model useful relational node nodes represented by the difference of learned *node positional embeddings*. Our contributions are as follows:

• We propose a Deformable Graph Convolutional Network (Deformable GCN) that simultaneously performs convolution in a latent space and adaptively deforms the convolution kernels to handle heterogeneity and capture long-range dependencies between nodes.

• We evaluate our method on the Parrot dataset. Deformable Graph Convolution Networks (Deformable GCN) that simultaneously learn node representations (features) and node positional embeddings (coordinates) and efficiently perform Deformable GCN in multiple latent spaces using



Figure 1. Time vs. Performance analysis for HOI detectors on V-COCO dataset. HOI recognition time is measured by subtracting the time of detection from the end-to-end inference time. Blue circle represents sequential HOI detectors, orange circle represents parallel HOI detectors and red star represents HOI detection benchmarks with an inference time under 1 ms after object detection.

1. Introduction

Human-Object Interaction (HOI) detection has been formally defined in [1] as the task to predict a set of possible object-object interactions within an image. Previous methods have addressed this task in an indirect manner by performing object detection first and associating (human-object) pairs afterward with separate post-processing steps. Especially for sequential HOI detectors [2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12], they have performed this association with a subsequent neural network, thus being time-consuming and computationally expensive.

To overcome the redundant inference structure of sequential HOI detectors, researchers [10, 11, 12] proposed parallel HOI detectors. These works explicitly localize interactions with either interaction boxes (i.e., the tightest box that covers both the center point of an object

*corresponding author
†corresponding author

74

NeurIPS 19

NeurIPS 21

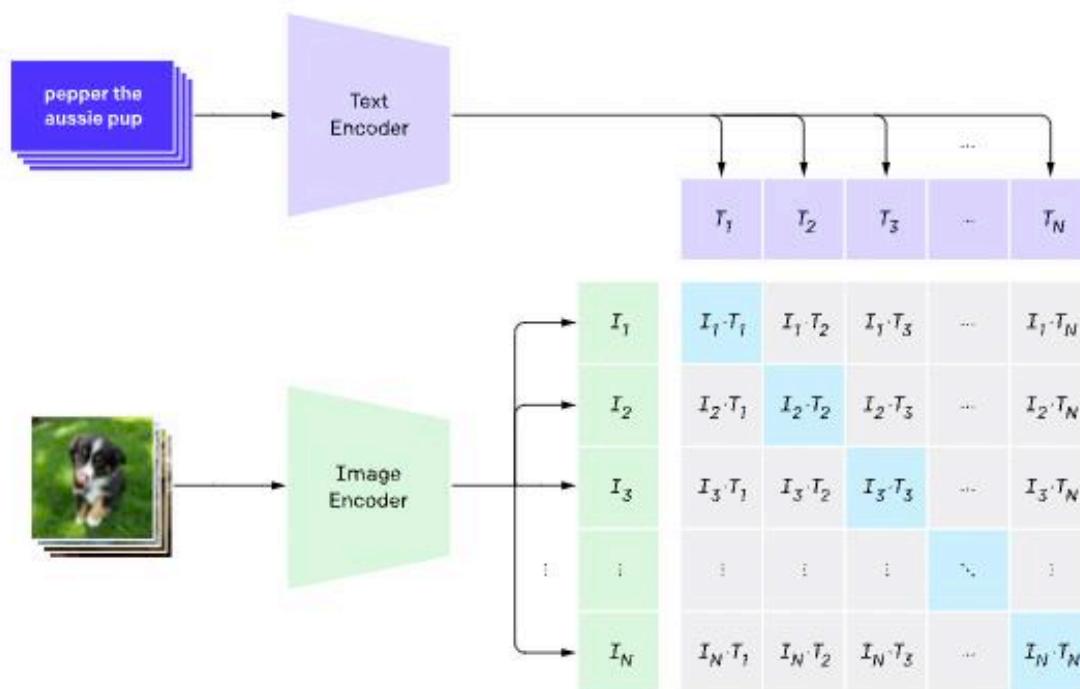
AAAI 22

CVPR 21 (Oral)

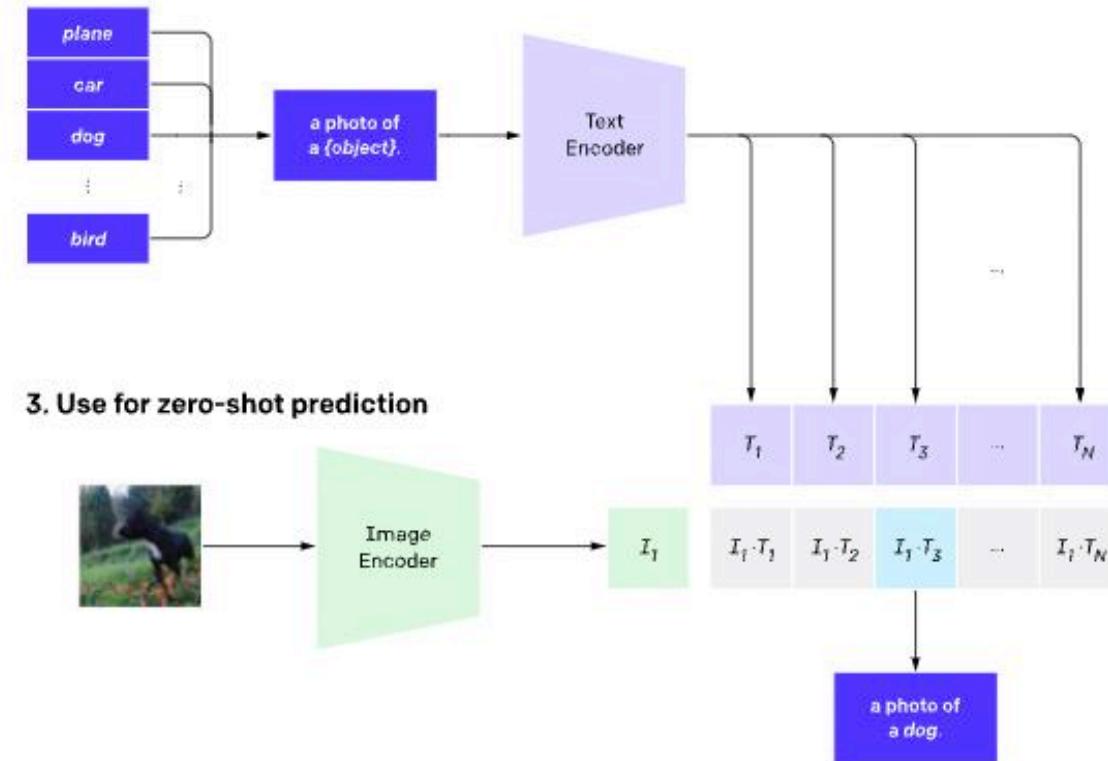
Foundation Models

CLIP: Connecting Text and Images

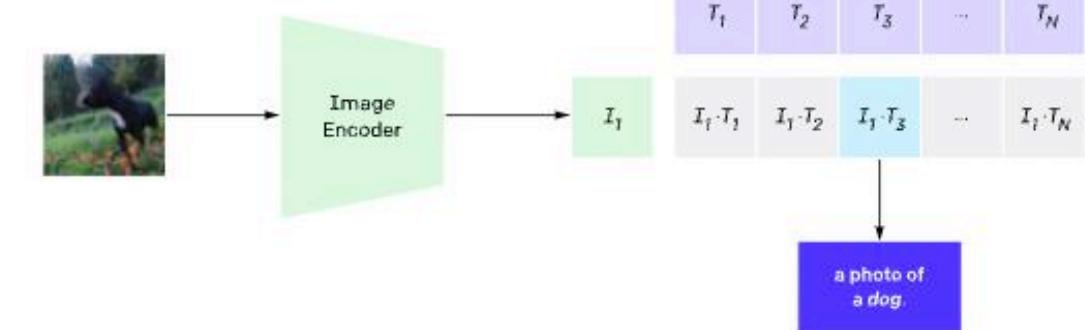
1. Contrastive pre-training



2. Create dataset classifier from label text

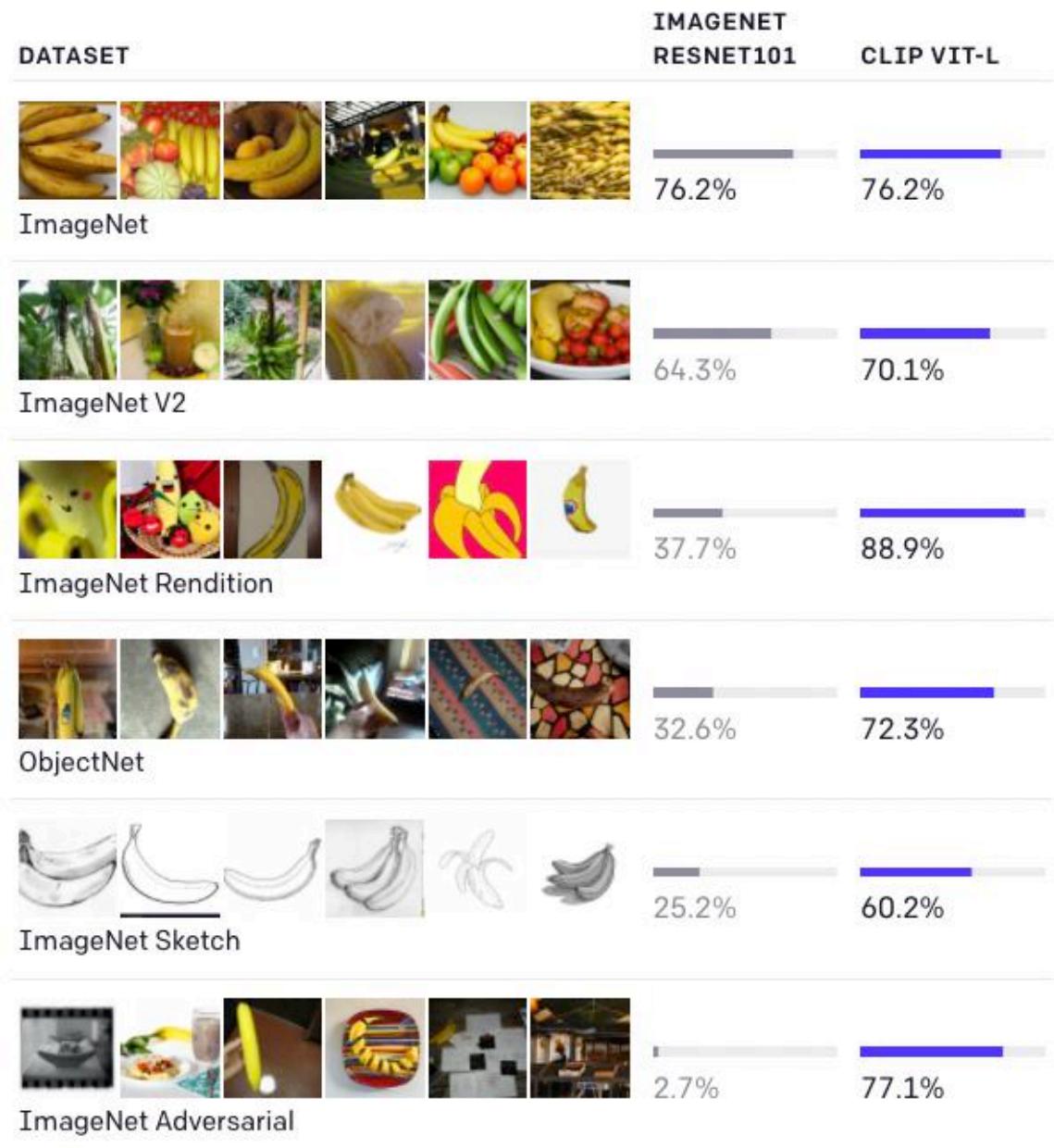


3. Use for zero-shot prediction



CLIP: Connecting Text and Images

- CLIP which efficiently learns visual concepts from natural language supervision.
- CLIP can be applied to any visual classification benchmark by simply providing the names of the visual categories to be recognized, similar to the “zero-shot” capabilities of GPT-2 and GPT-3.



kangaroo (99.8%) Ranked 1 out of 102



- ✓ a photo of a **kangaroo**.
- ✗ a photo of a **gerenuk**.
- ✗ a photo of a **emu**.
- ✗ a photo of a **wild cat**.
- ✗ a photo of a **scorpion**.

Siberian Husky (76.0%) Ranked 1 out of 200



- ✓ a photo of a **siberian husky**.
- ✗ a photo of a **german shepherd dog**.
- ✗ a photo of a **collie**.
- ✗ a photo of a **border collie**.
- ✗ a photo of a **rottweiler**.

airplane, person (89.0%) Ranked 1 out of 23



- ✓ a photo of a **airplane**.
- ✗ a photo of a **bird**.
- ✗ a photo of a **bear**.
- ✗ a photo of a **giraffe**.
- ✗ a photo of a **car**.

annual crop land (12.9%) Ranked 4 out of 10



- ✗ a centered satellite photo of **permanent crop land**.
- ✗ a centered satellite photo of **pasture land**.
- ✗ a centered satellite photo of **highway or road**.
- ✓ a centered satellite photo of **annual crop land**.
- ✗ a centered satellite photo of **brushland or shrubland**.

Average linear probe score across 27 datasets

85%

80%

75%

70%

65%

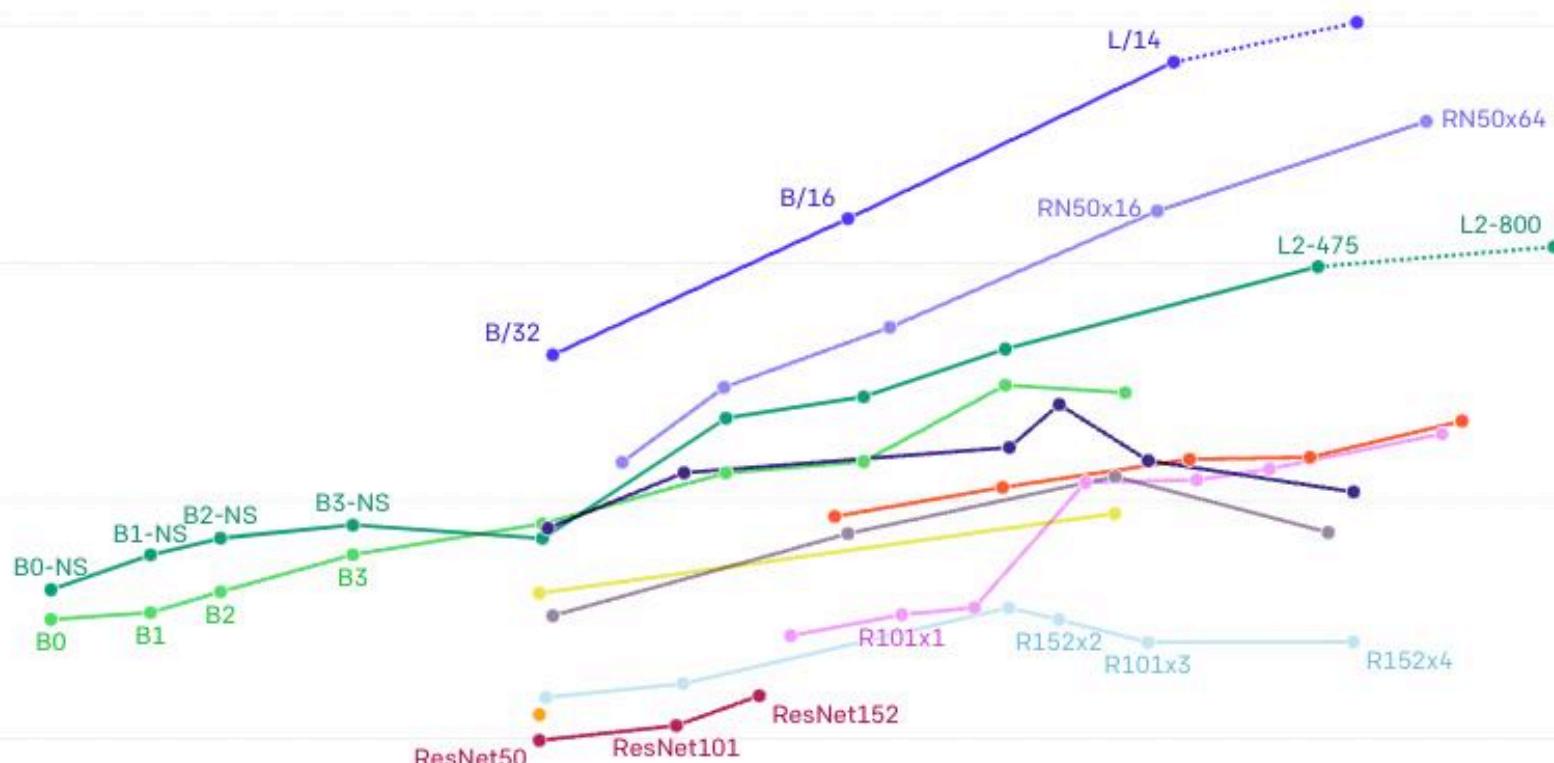
1 10 100

Forward-pass GFLOPs/image

- CLIP-ViT
- CLIP-ResNet
- EfficientNet-NoisyStudent
- EfficientNet

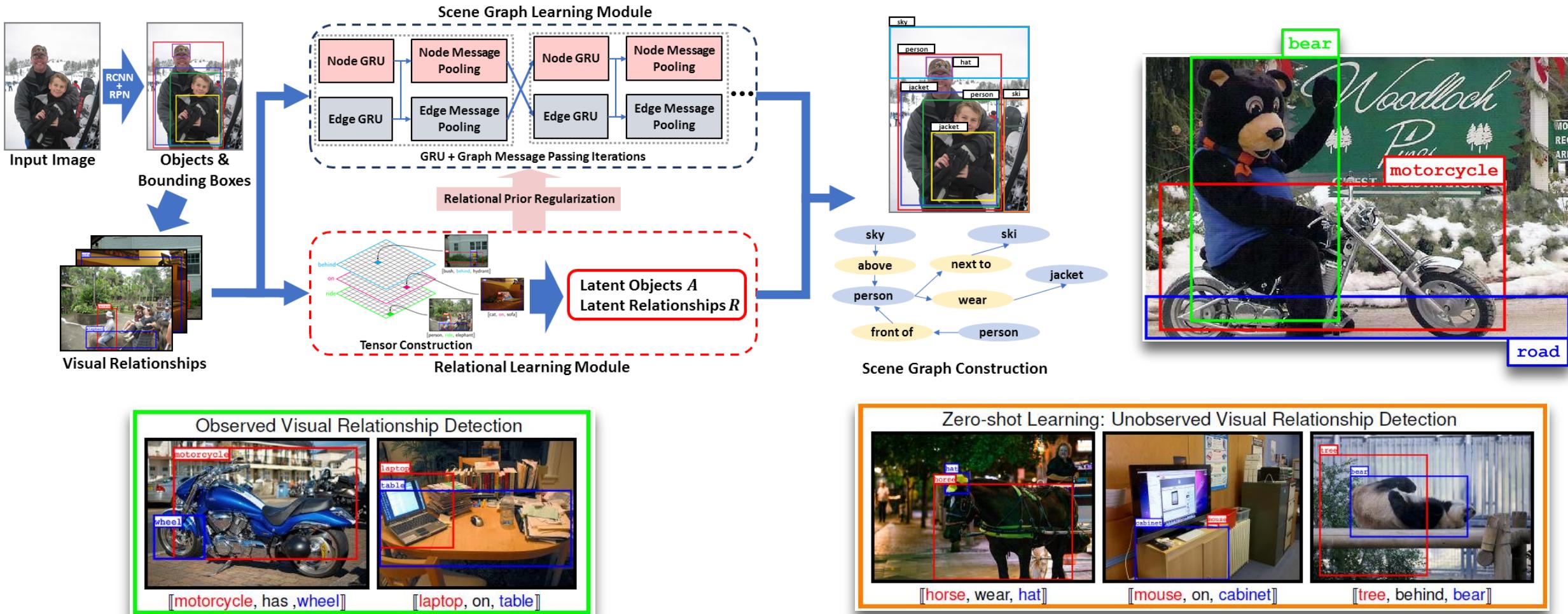
- Instagram
- SimCLRV2
- BYOL
- MoCo

- ViT (ImageNet-21k)
- BiT-M
- BiT-S
- ResNet

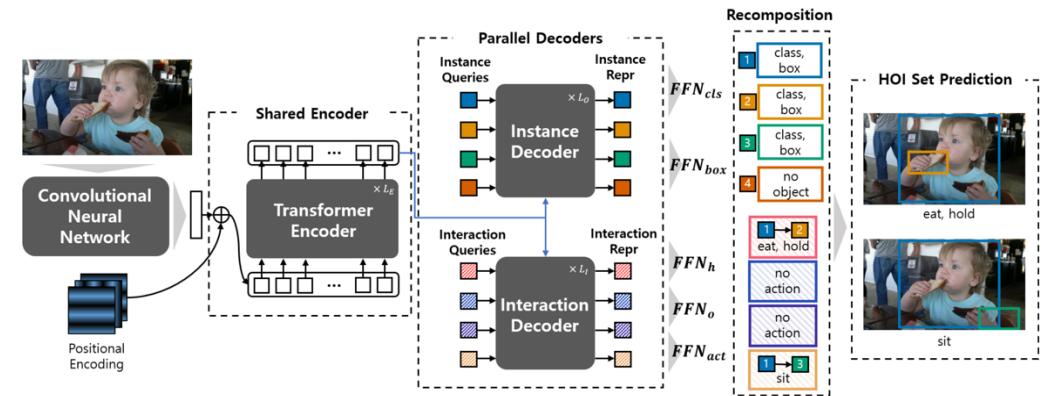
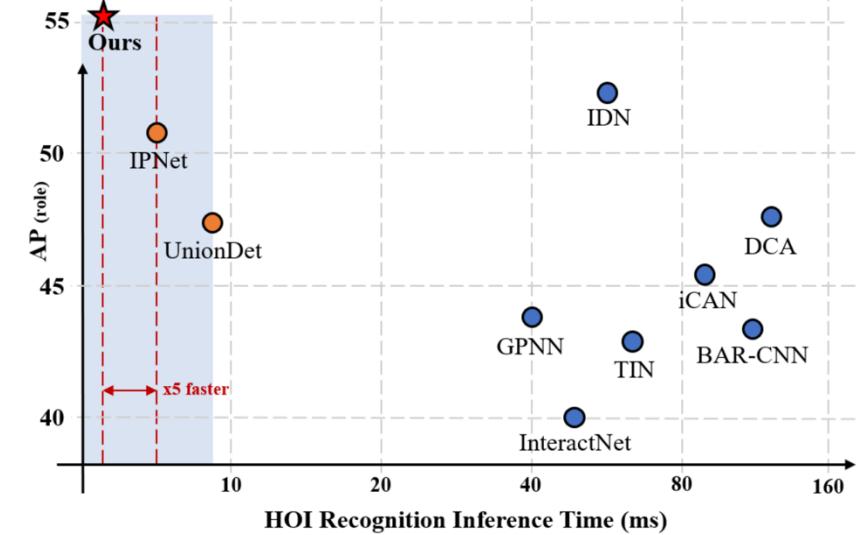
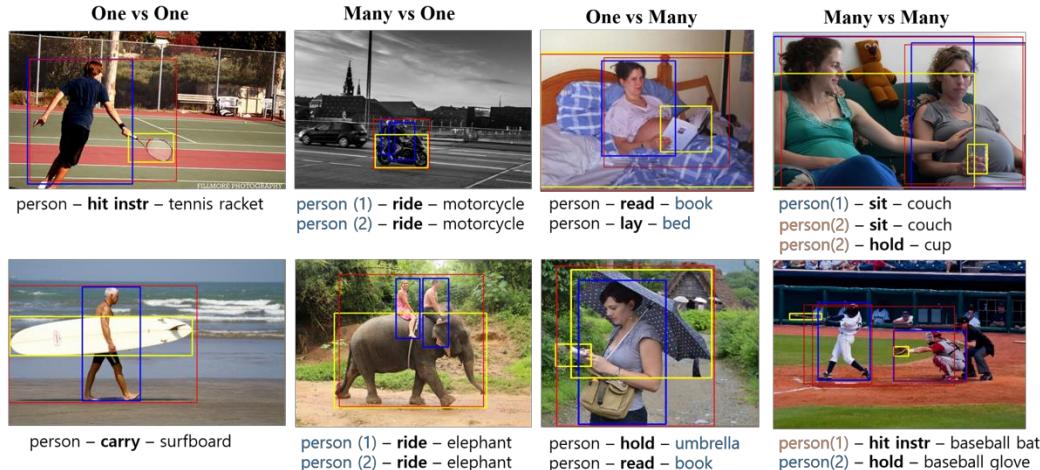
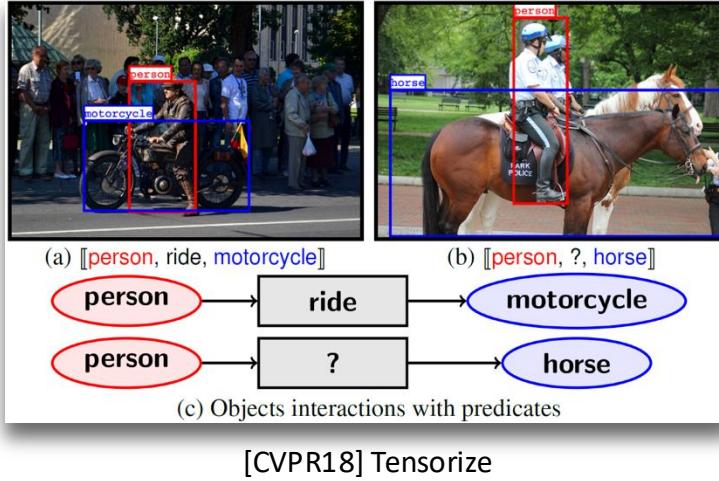


Scene understanding

Scene graph generation



Scene understanding and object detection



[CVPR '21 (ORAL)] HOTR

3초

WBSC

이스라엘 0

대한민국 3

25

1-1

TOKYO 2020



hold
hit_obj



hold
hit_inst

야구 본선 2라운드 연장 중계
김나진·허구연·김선우

esos

물체 탐지 강의 (데이터 과학원)

KUIDS 인공지능은 물체를 어떻게 이해할까? : 물체 탐지 원리

Watch later Share

MORE VIDEOS

Wu, Zhenyu, et al. ICCV 2019

데이터 과학원 온라인 강의: 물체 탐지 원리 6/25

YouTube

1. 물체 탐지 원리
[\[Link\]](#)



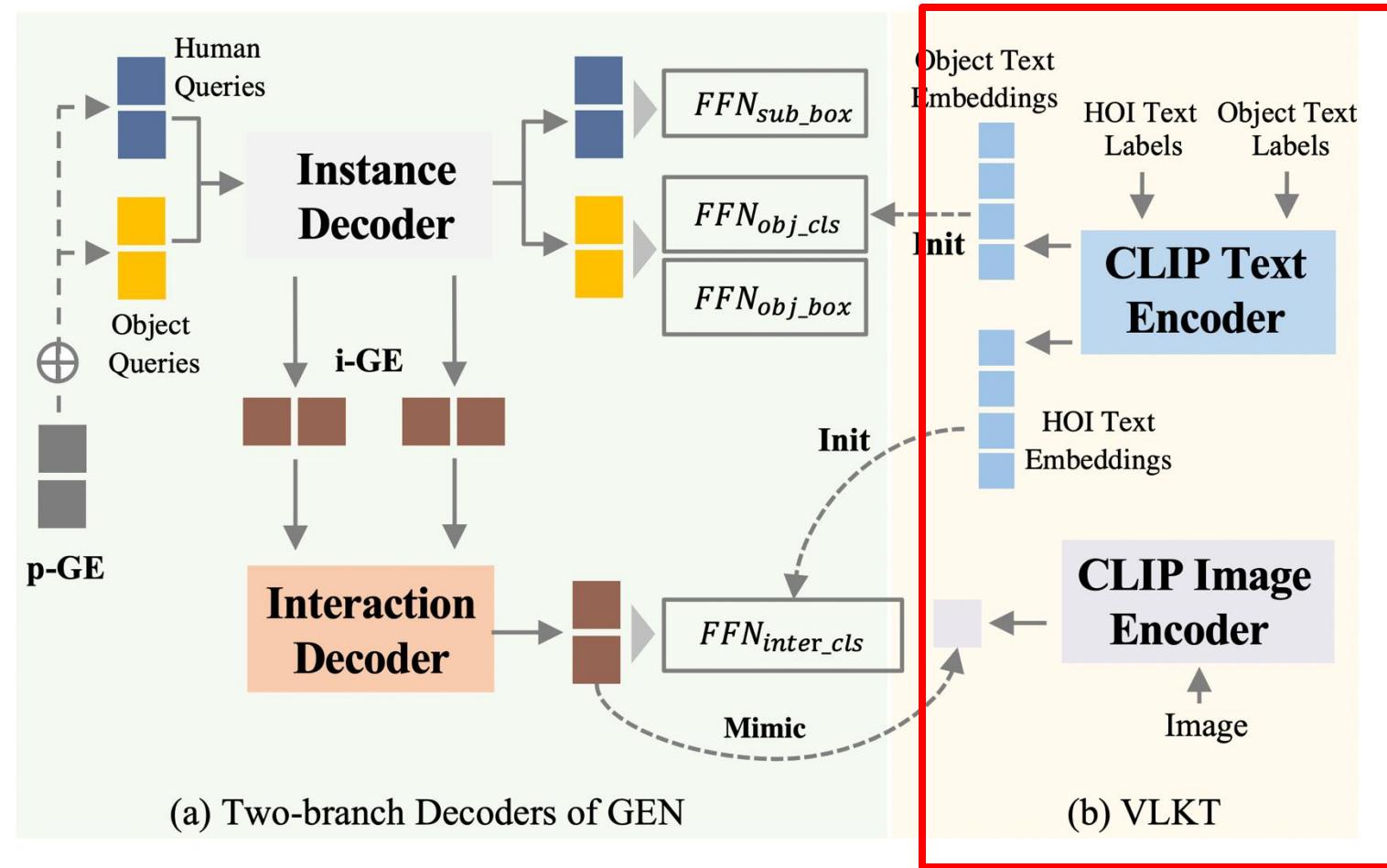
2. 물체 탐지 기법
[\[Link\]](#)



3. 물체간 관계 추론
[\[Link\]](#)

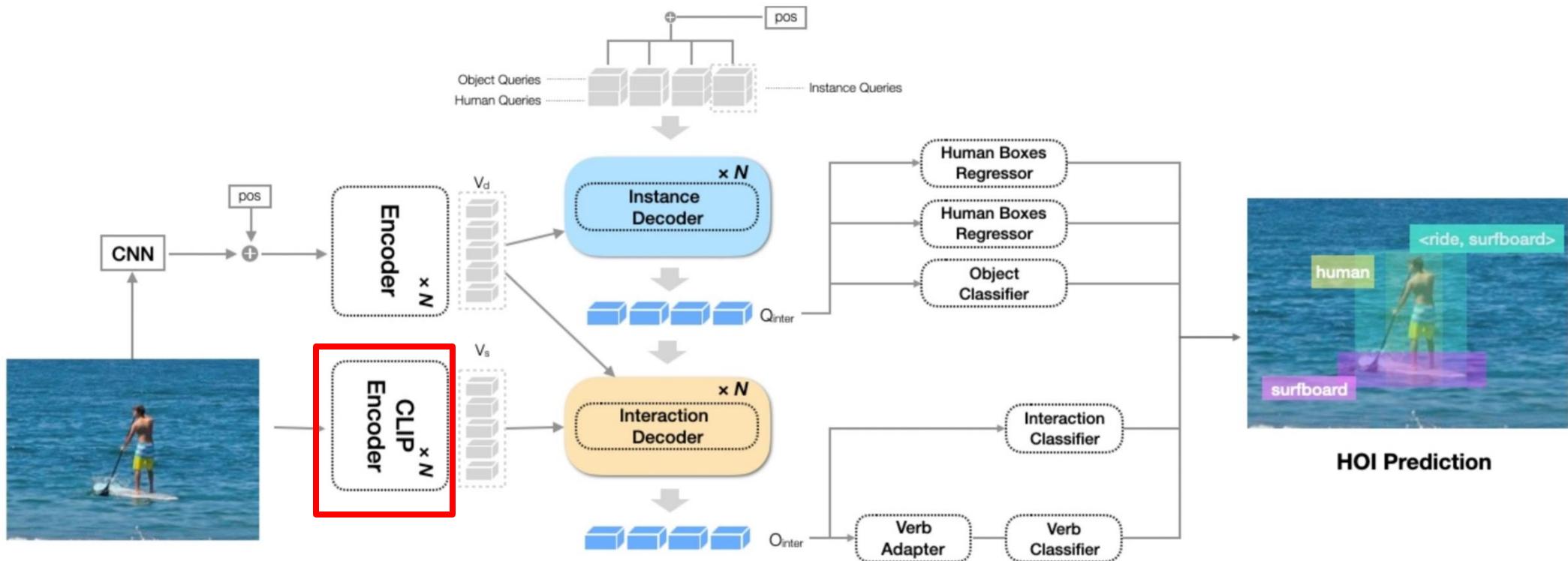


초거대 AI + 장면이해? Latent space alignment



[CVPR22] GEN-VLKT: Simplify Association and Enhance Interaction Understanding for HOI Detection

초거대 AI + 장면이해?



HOICLIP: Efficient Knowledge Transfer for HOI Detection with Vision-Language Models

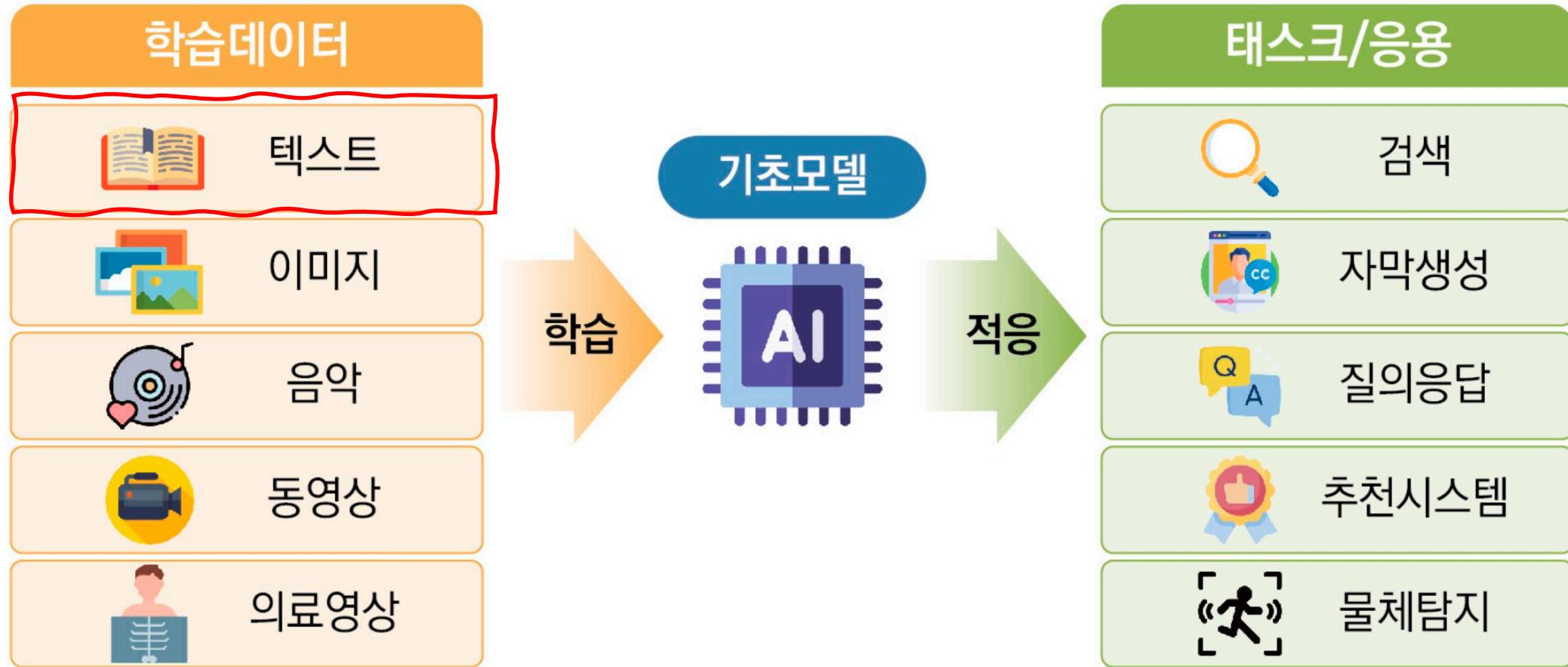
Examples: MLV top-tier conference papers (2023)

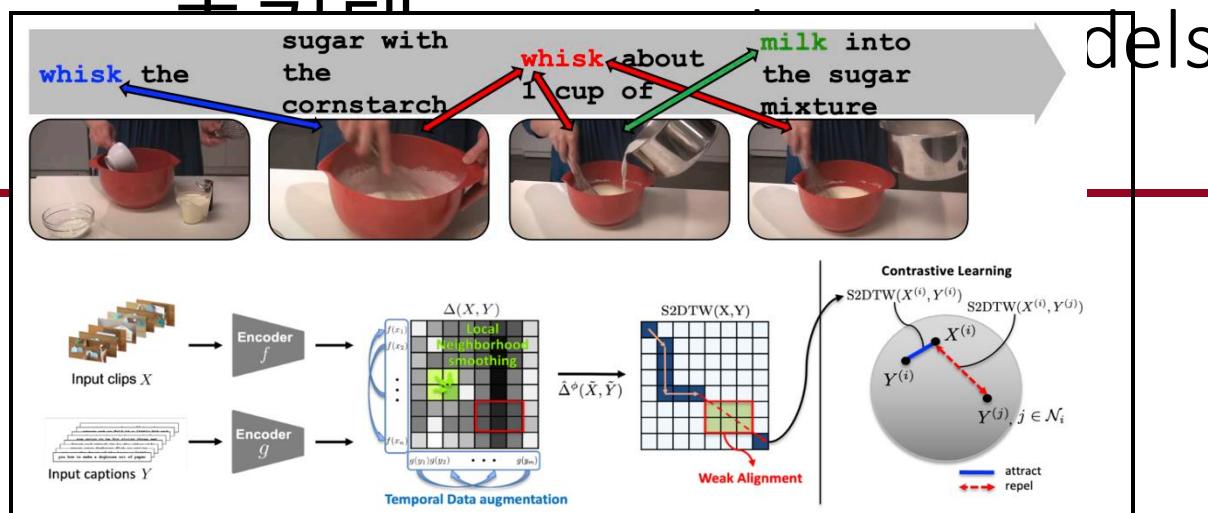
1. [EMNLP23] Large Language Models are Temporal and Causal Reasoners for Video Question Answering
2. [NeurIPS23] NuTrea: Neural Tree Search for Context-guided Multi-hop KGQA
3. [NeurIPS23] Advancing Bayesian Optimization via Learning Smooth Latent Spaces
4. [NeurIPS23] Unconstrained Pose Prior-Free Neural Radiance Field
5. [ICCV23] Open-Vocabulary Video Question Answering: A New Benchmark for Evaluating the Generalizability of Video Question Answering Models
6. [ICCV23] Distribution-Aware Prompt Tuning for Vision-Language Models
7. [ICCV23] Semantic-Aware Template Learning via Part Deformation Consistency
8. [ICCV23] Read-only Prompt Optimization for Vision-Language Few-shot Learning
9. [ICML23] Robust Camera Pose Refinement for Multi-Resolution Hash Encoding
10. [CVPR23] MELTR: Meta Loss Transformer for Learning to Fine-tune Video Foundation Models
11. [CVPR23] Self-positioning Point-based Transformer for Point Cloud Understanding
12. [AAAI23] Relation-aware Language-Graph Transformer for Question Answering

Examples: MLV top-tier conference papers (2024)

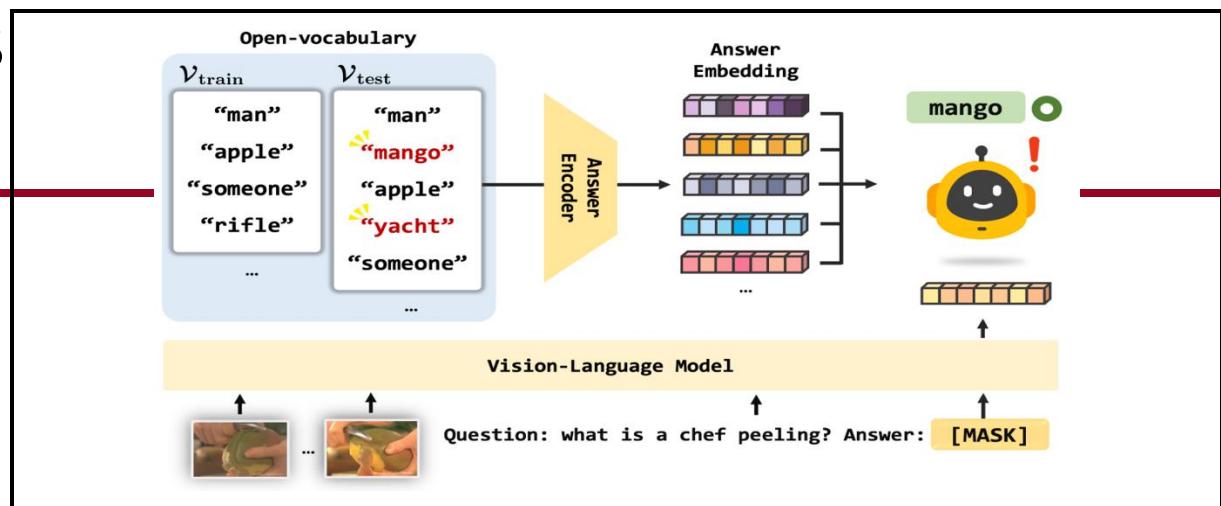
1. [NeurIPS] LLaMo: Large Language Model-based Molecular Graph Assistant
2. [NeurIPS] Alleviating Attention Bias for Visual-Informed Text Generation
3. [NeurIPS] Constant Acceleration Flow
4. [NeurIPS] Inversion-based Latent Bayesian Optimization
5. [EMNLP] Generative Subgraph Retrieval for Knowledge Graph–Grounded Dialog Generation
6. [ECCV] Diffusion Prior-Based Amortized Variational Inference for Noisy Inverse Problems
7. [ECCV] Understanding Multi-compositional learning in Vision and Language models via Category Theory
8. [ICML] Stochastic Conditional Diffusion Models for Robust Semantic Image Synthesis
9. [CVPR] vid-TLDR: Training Free Token merging for Light-weight Video Transformer
10. [CVPR] Multi-criteria Token Fusion with One-step-ahead Attention for Efficient Vision
11. [CVPR] Groupwise Query Specialization and Quality-Aware Multi-Assignment for Transformer-based Visual Relationship Detection
12. [CVPR] Prompt Learning via Meta-Regularization
13. [CVPR] Retrieval-Augmented Open-Vocabulary Object Detection
14. [ICLR] Domain-agnostic Latent Diffusion Models for Synthesizing High-Quality Implicit Neural Representations

Foundation Models (기초모델)

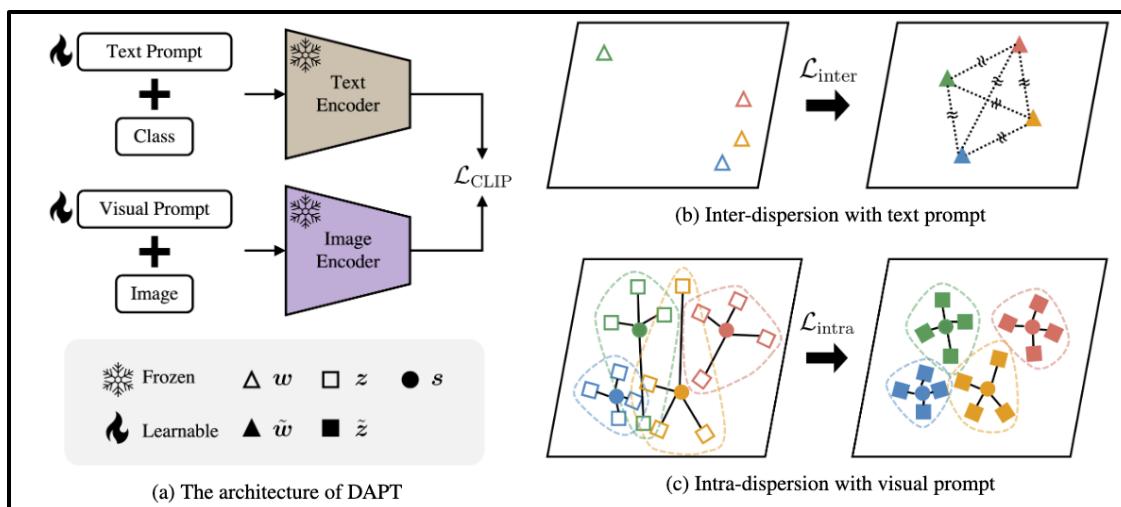




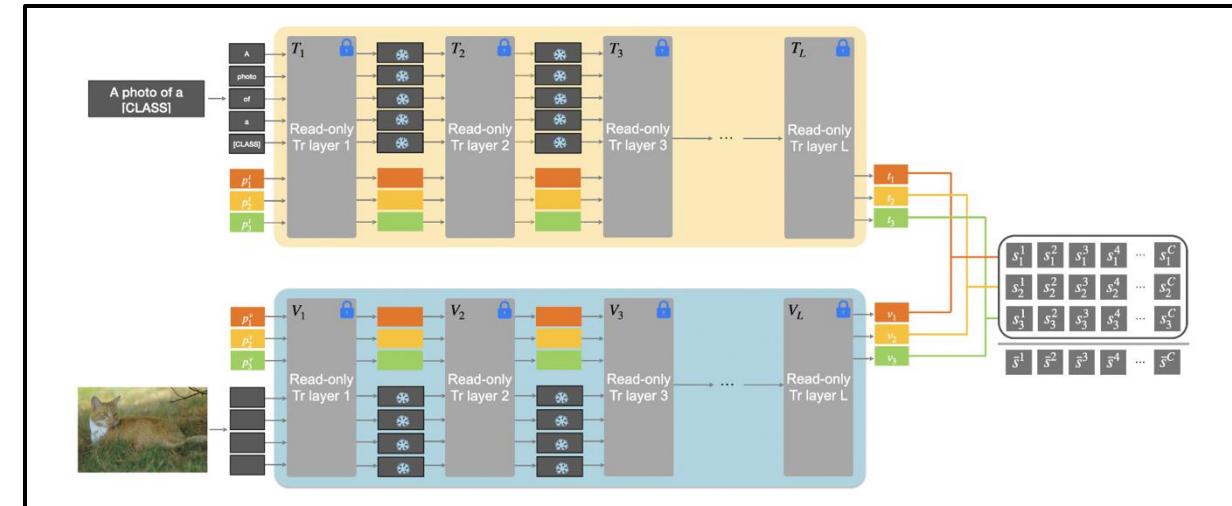
[CVPR22] VT-TWIN. (1억개 유튜브 클립학습. 비디오 기초모델)



[ICCV23] 거대 동영상 기초모델 기반 오픈-사전 비디오 QA 시스템



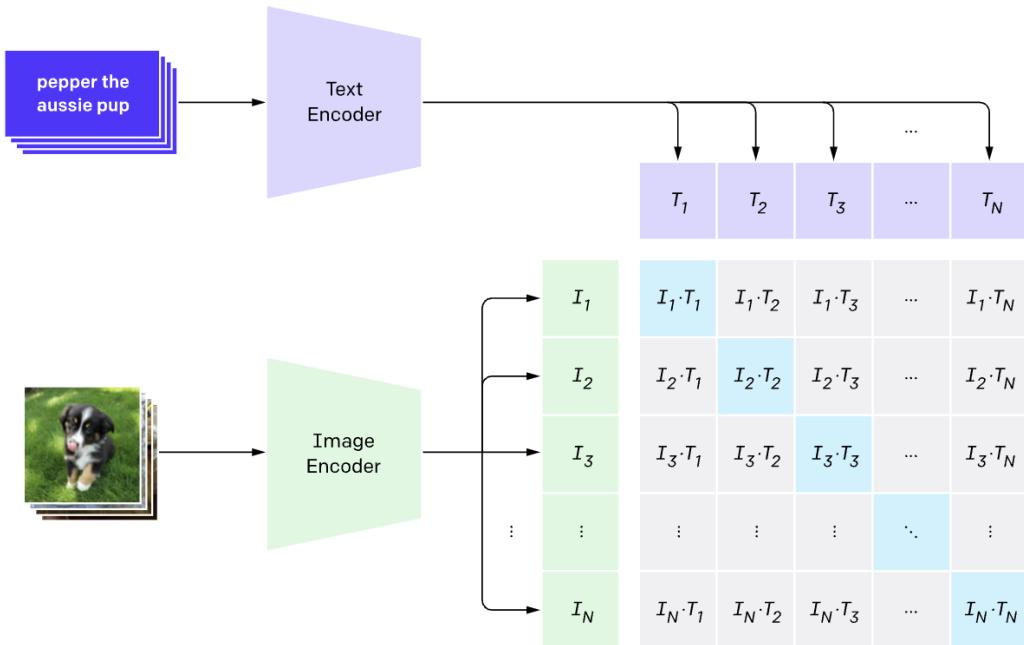
[ICCV23] DAPT. CLIP 기반 few-shot learning



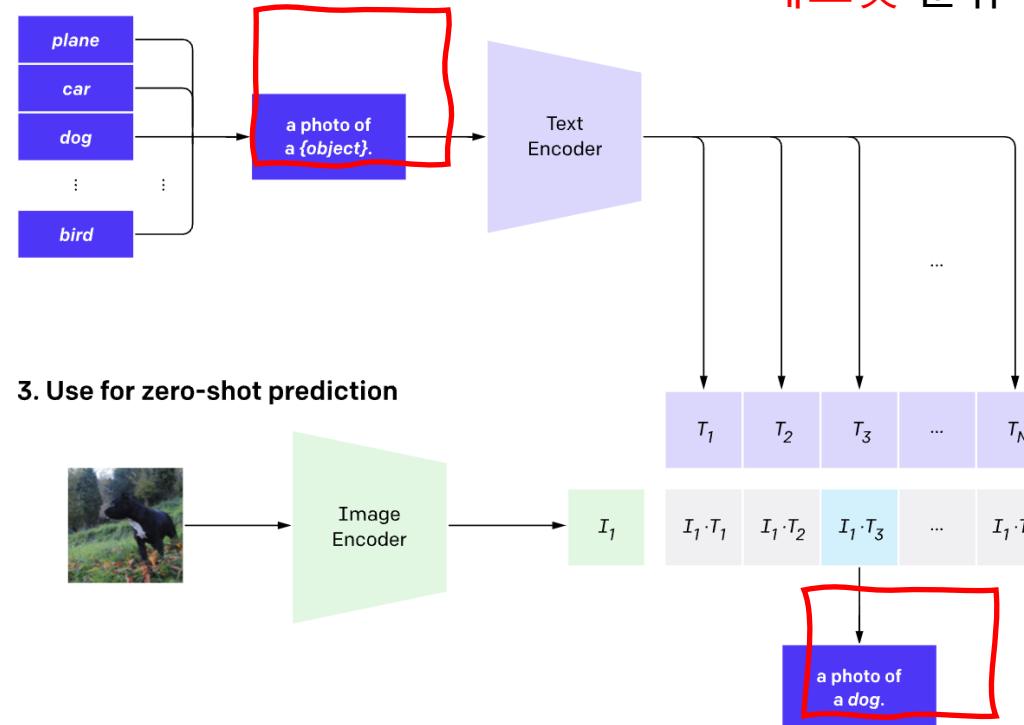
[ICCV23] RPT. 기초모델 위한 읽기 중심 프롬프트 튜닝

Contrastive Learning (CLIP)

1. Contrastive pre-training



2. Create dataset classifier from label text

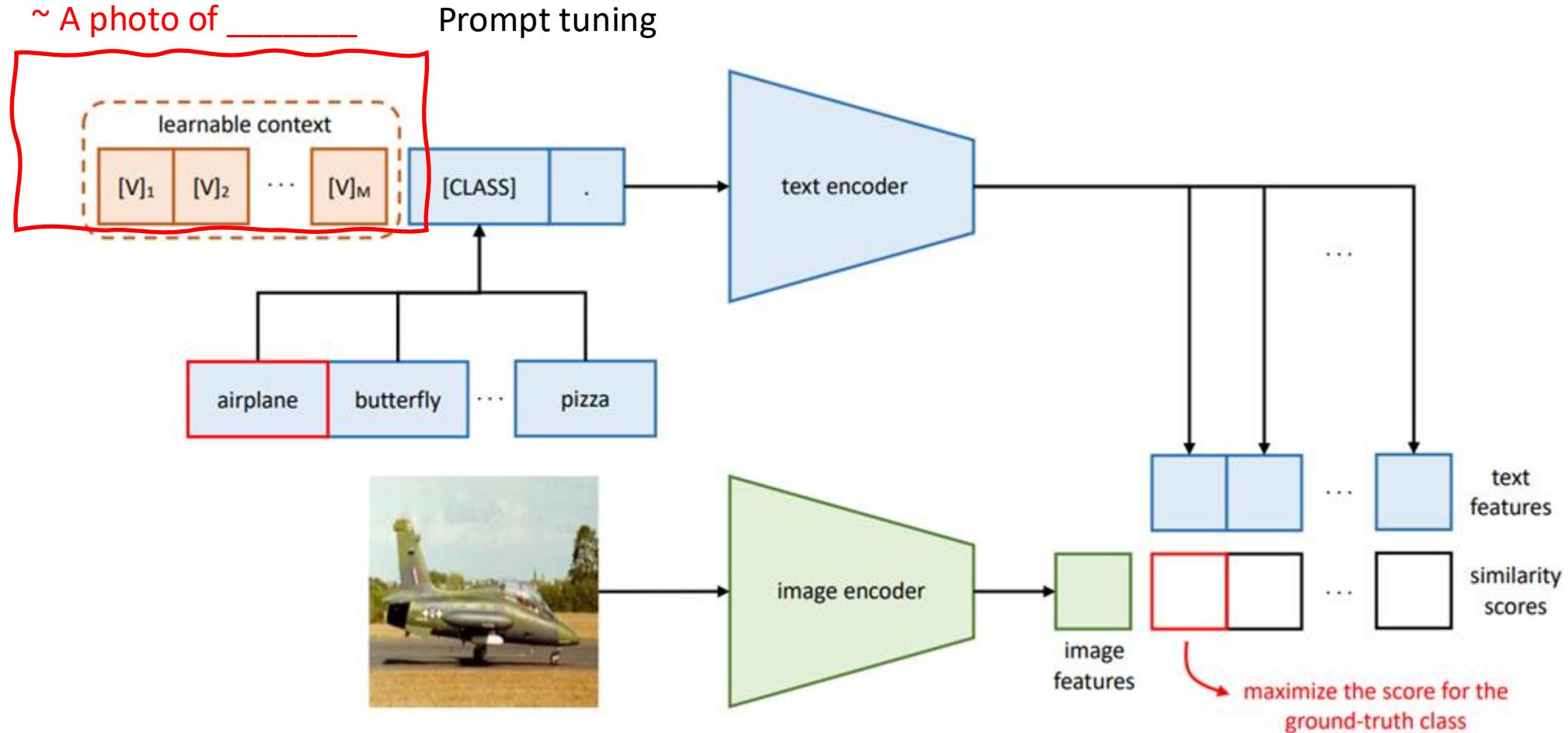


- A photo of _____
- **프롬프트 엔지니어링**
 - **제로샷 분류기**

Self-Supervised Learning (자기지도학습)

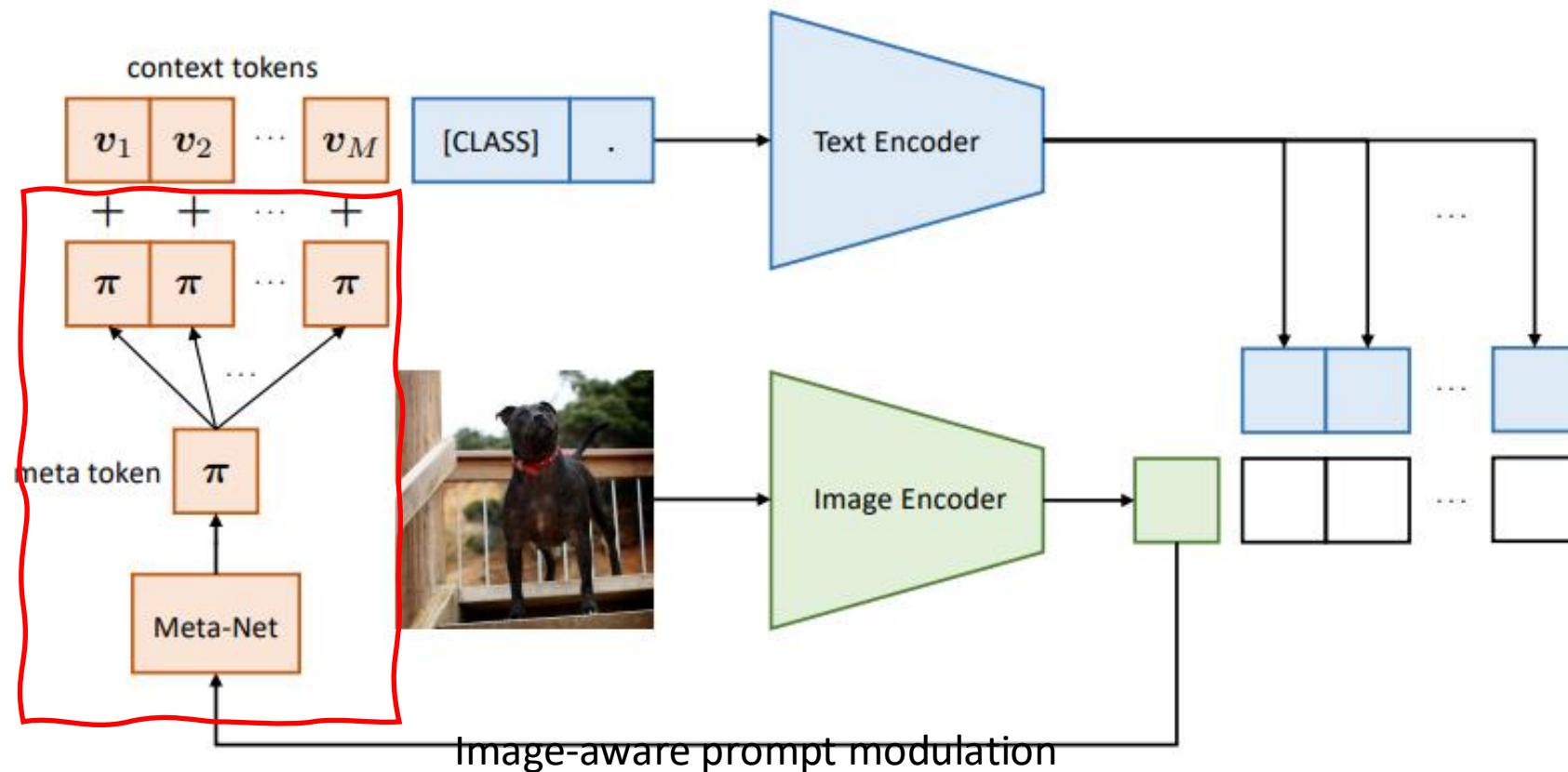
Radford, Alec, et al. "Learning transferable visual models from natural language supervision." ICML 2021.

Prompt Learning (CoOp)

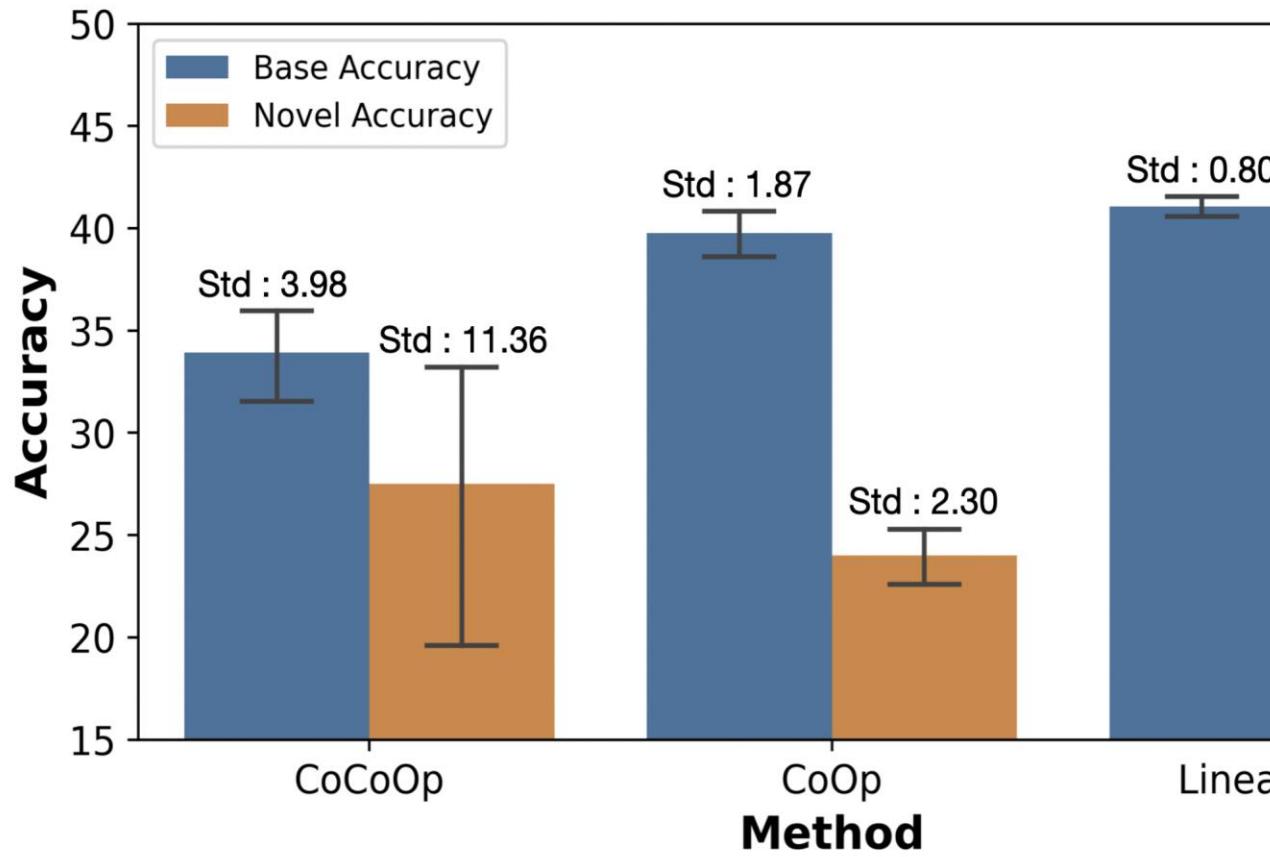


Learning to prompt for vision-language models (IJCV 2022)

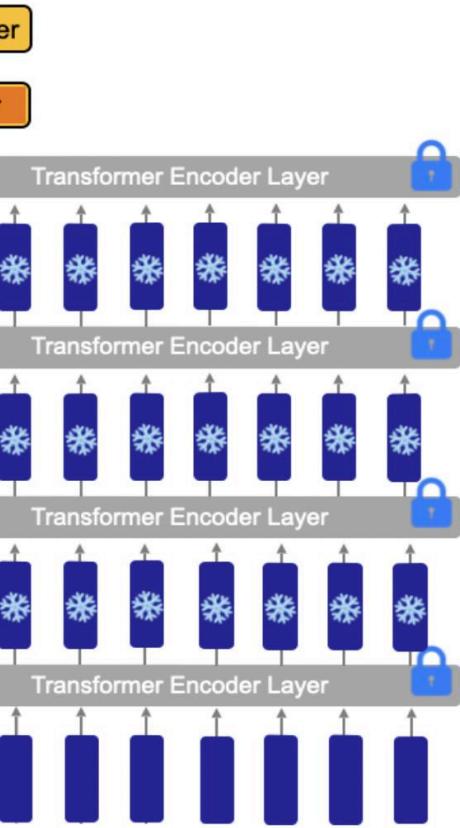
CoCoOp



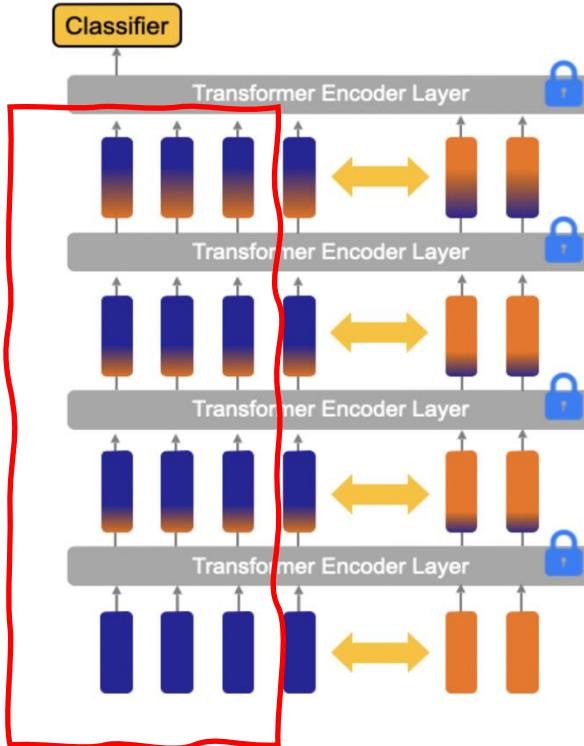
Our Observation



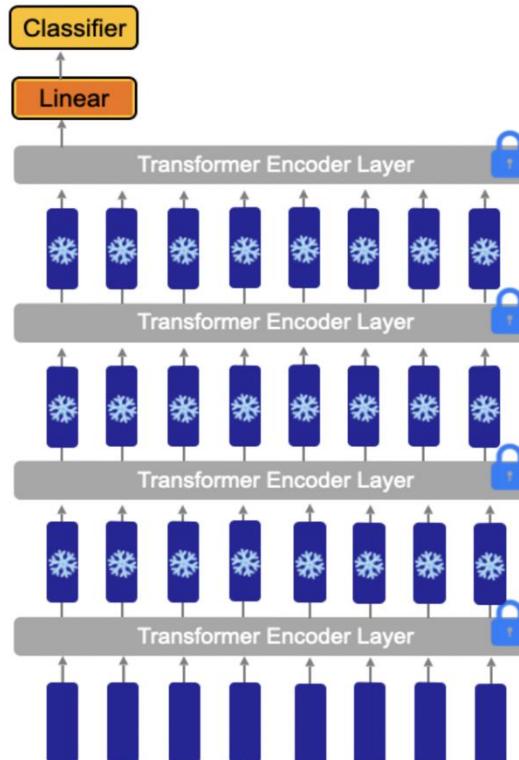
Linear Probe achieved better performance with small VAR (std).



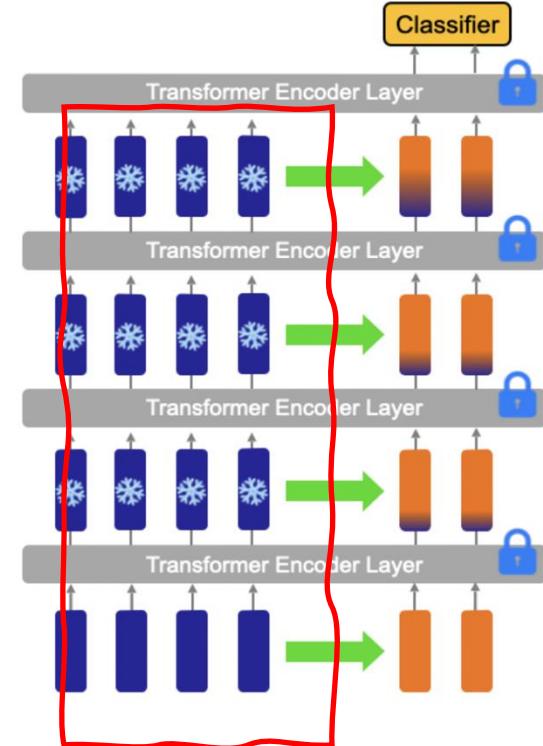
Read-only Prompt Optimization



(a) Conventional Prompt Tuning

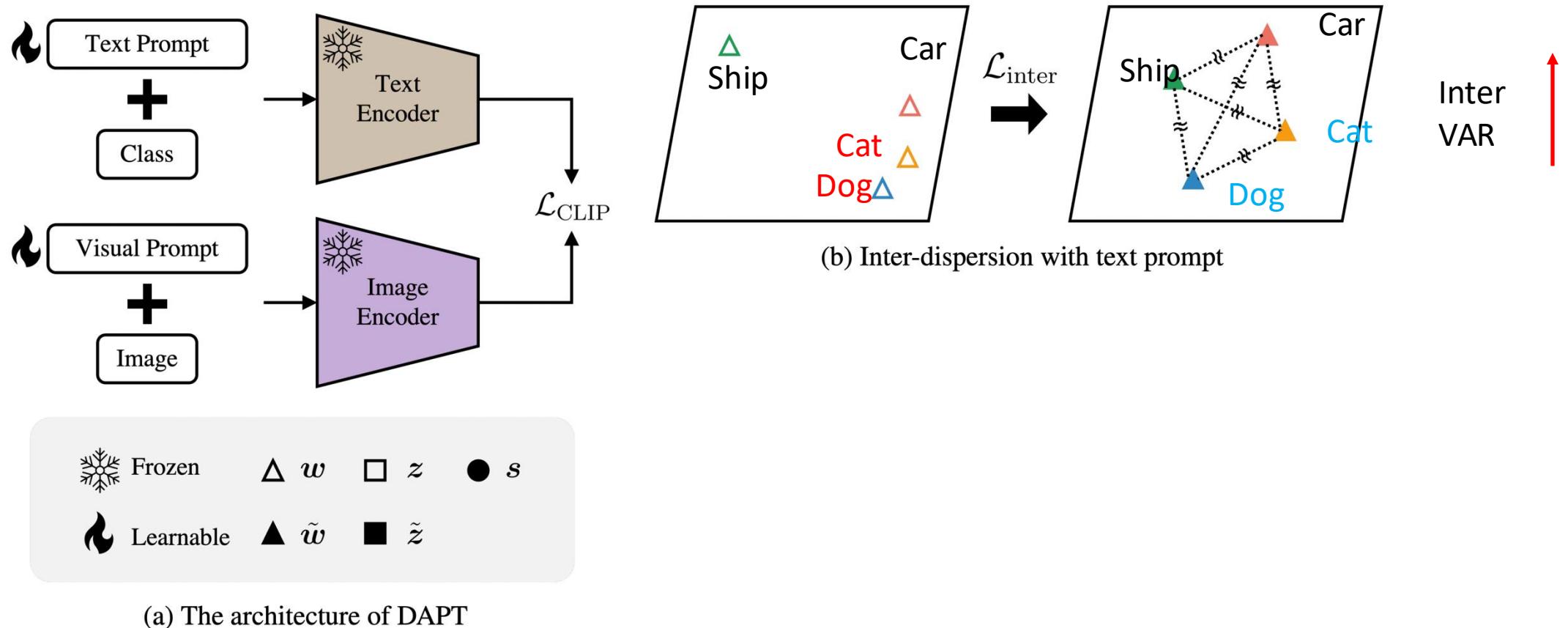


(b) Linear Probing

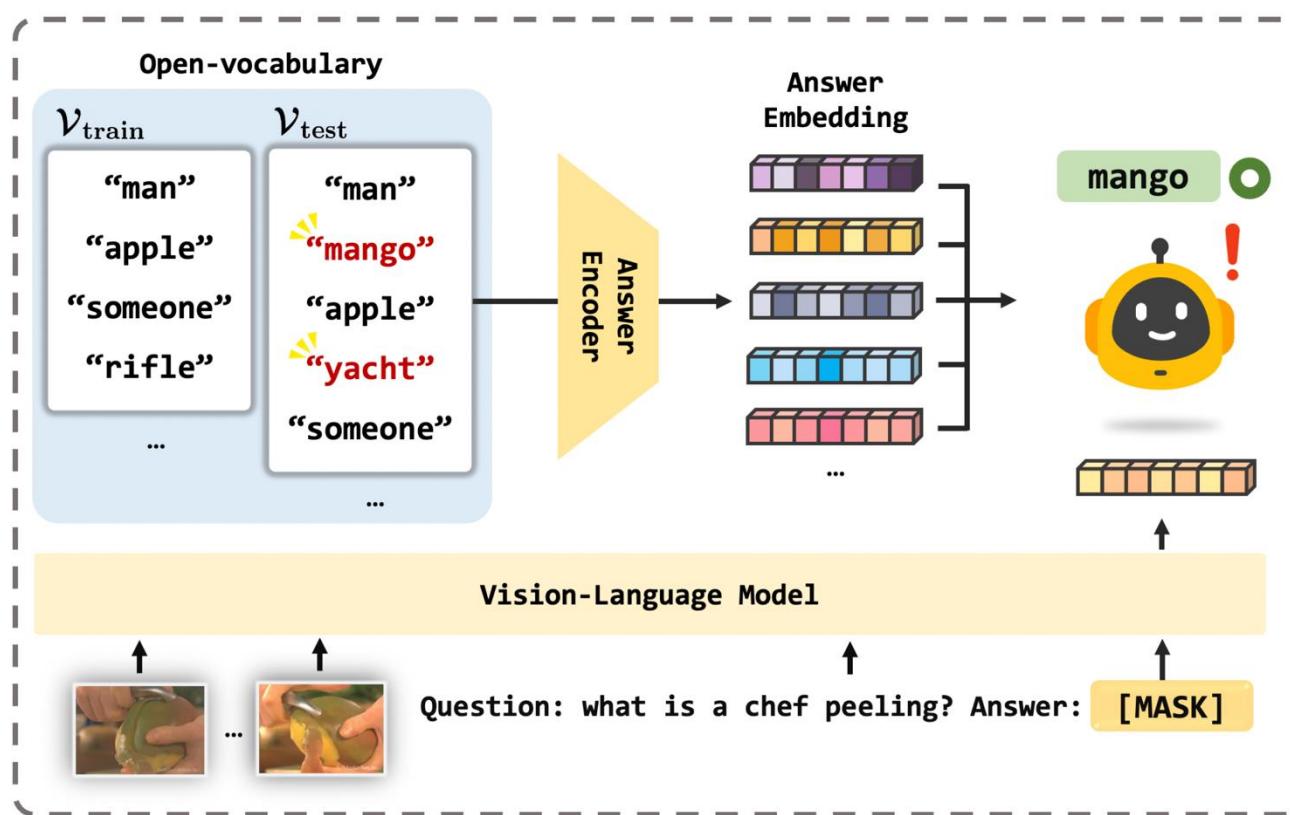


(c) Read-only Prompt Optimization

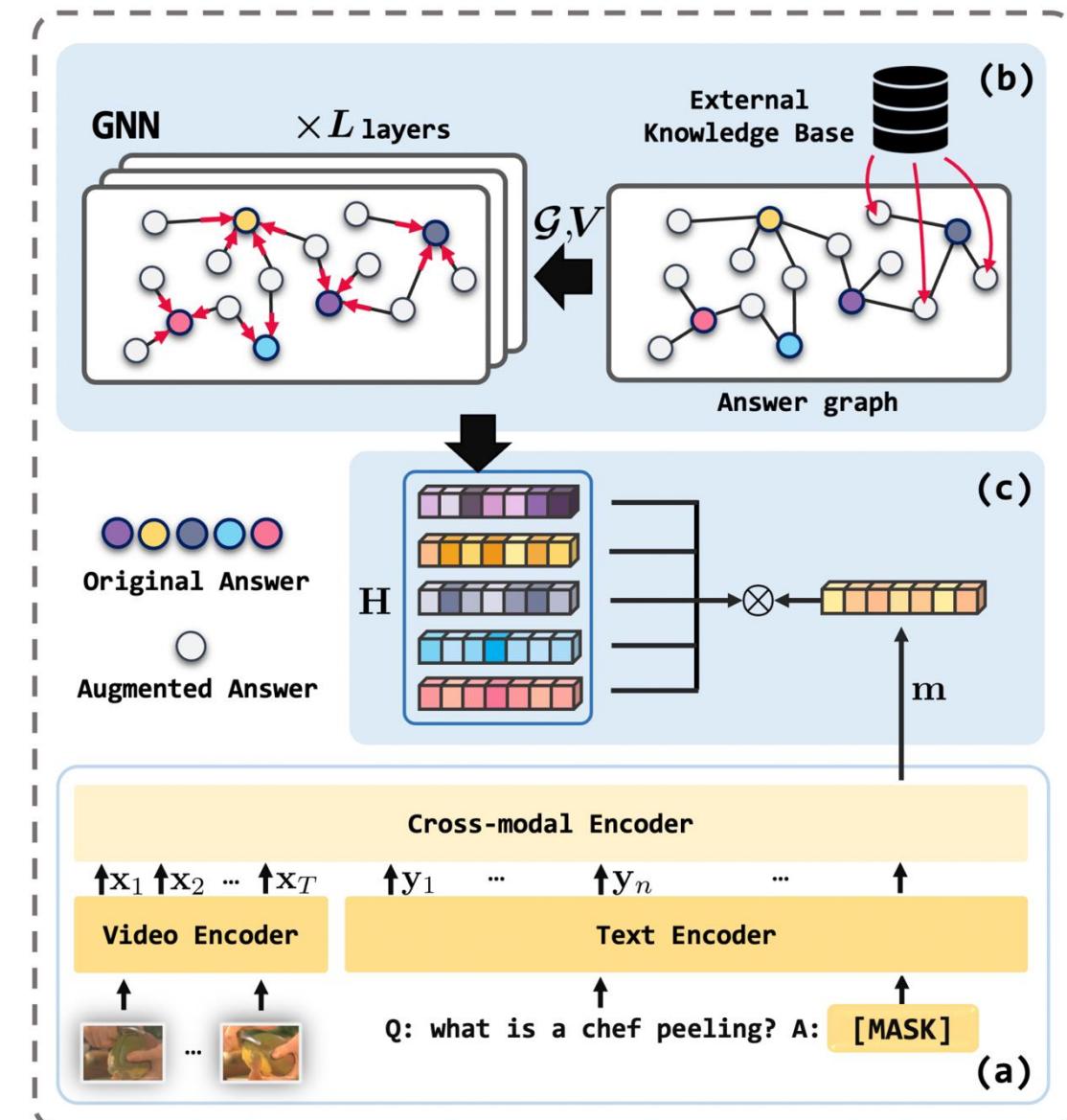
Distribution-Aware Prompt Tuning



Video Question Answering



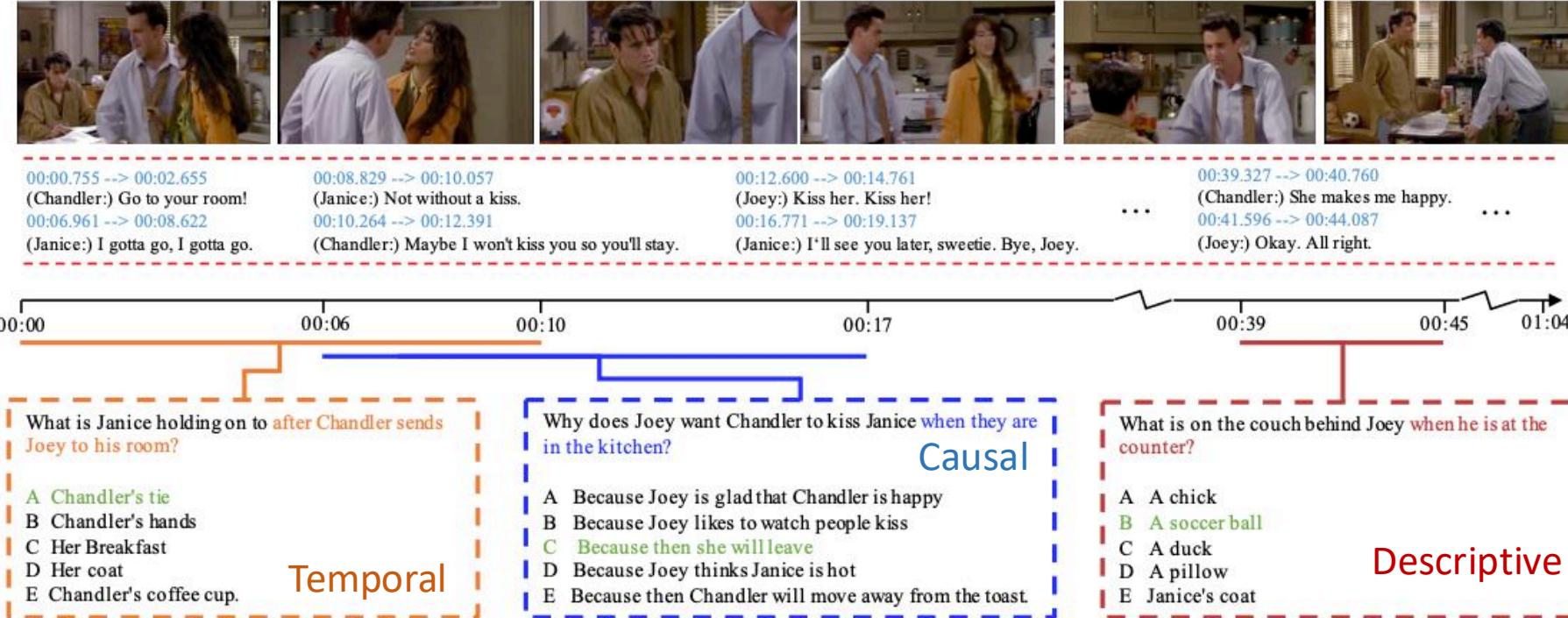
- New Benchmark
- Label smoothing by GNNs
- Vision-Language Model
- Cross-Modal Encoder thanks to Transformers



Video QA

Downstream tasks

- Video question answering



- Dataset:
 - TGIF
 - MSVD
- Evaluation metric:
 - Accuracy

Questions?
