# Text and 3D Representation Learning

# Overview

1. Learning text and image representations

2. Extending to text and 3D

3. Efficient text and 3D

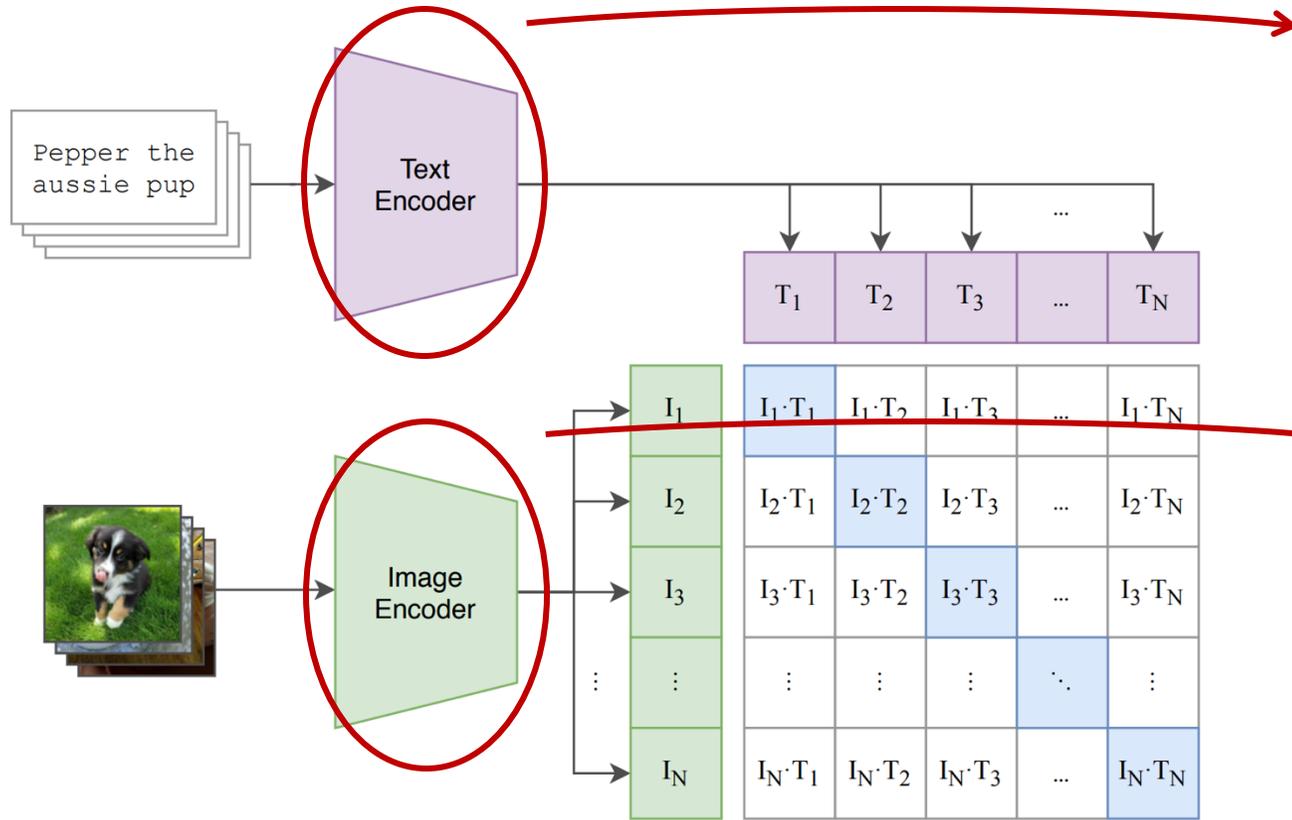4. Applications

# Preliminary (CLIP)

① Circular glass coffee table with two sets of wooden legs that clasp over the round glass edge. ①

① 

② A brown wooden moon shaped table with three decorative legs with a wooden vine shaped decoration base connecting the legs. 2a

Wooden half round table. 2b

③ Dark brown wooden chair with adjustable back rest and gold printed upholestry. 3a

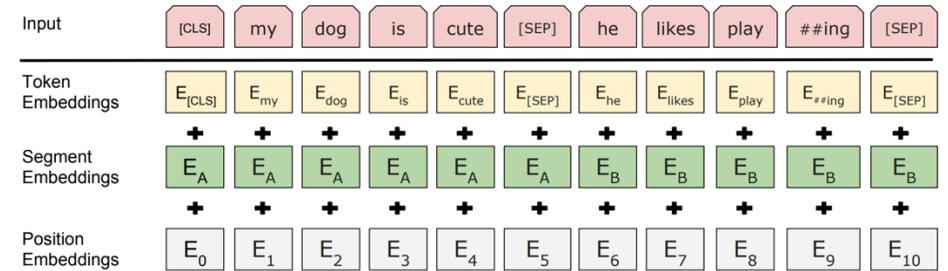Wooden recliner chair with patterned fabric. 3b

1. Have a bunch of images as well as captions from the Internet

2. How to learn a model to align embeddings for images and text?

# Preliminary (CLIP)
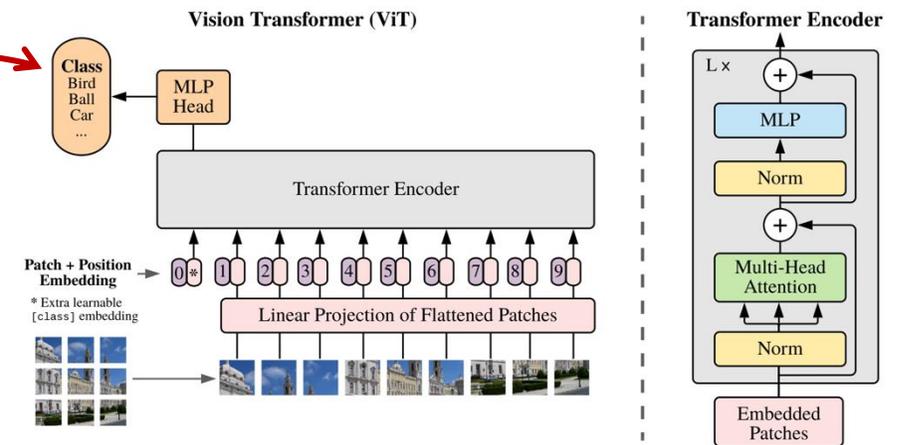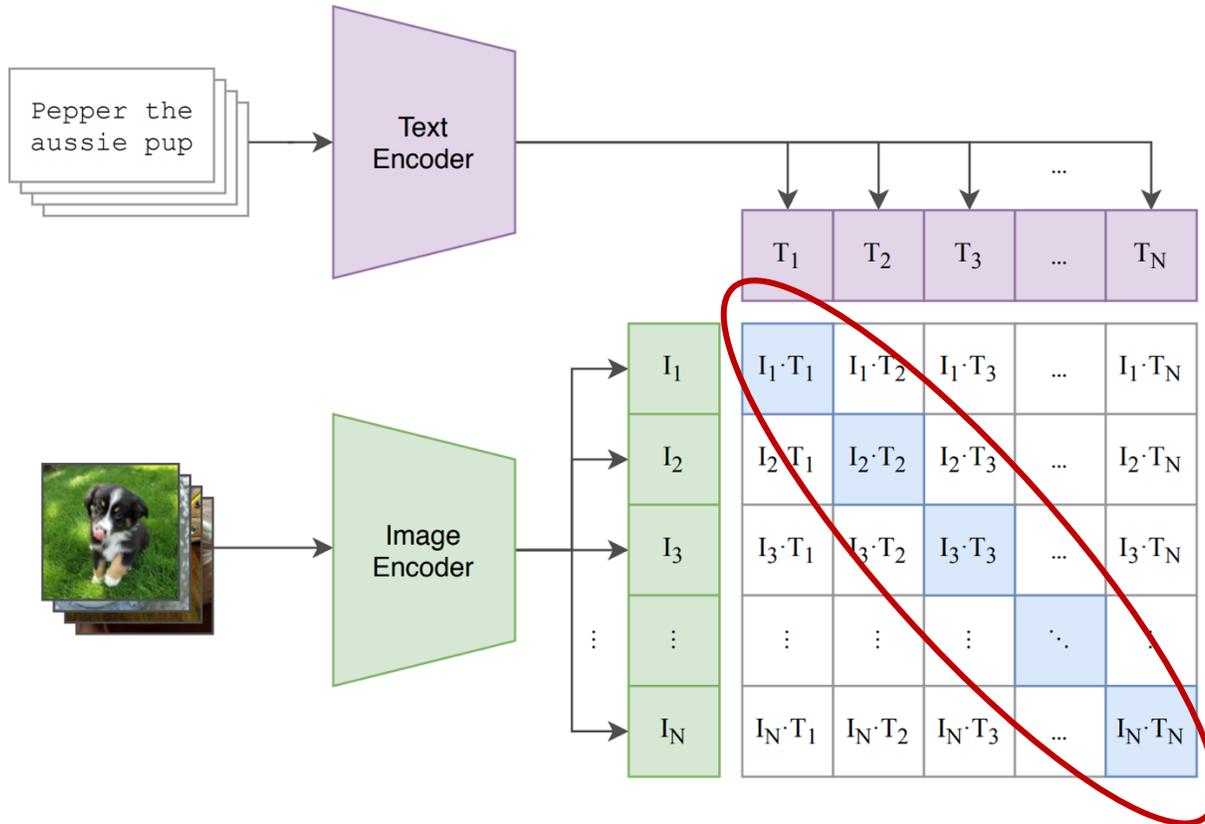


Decoder Only Transformer

CLIP (Radford et al. 2021)

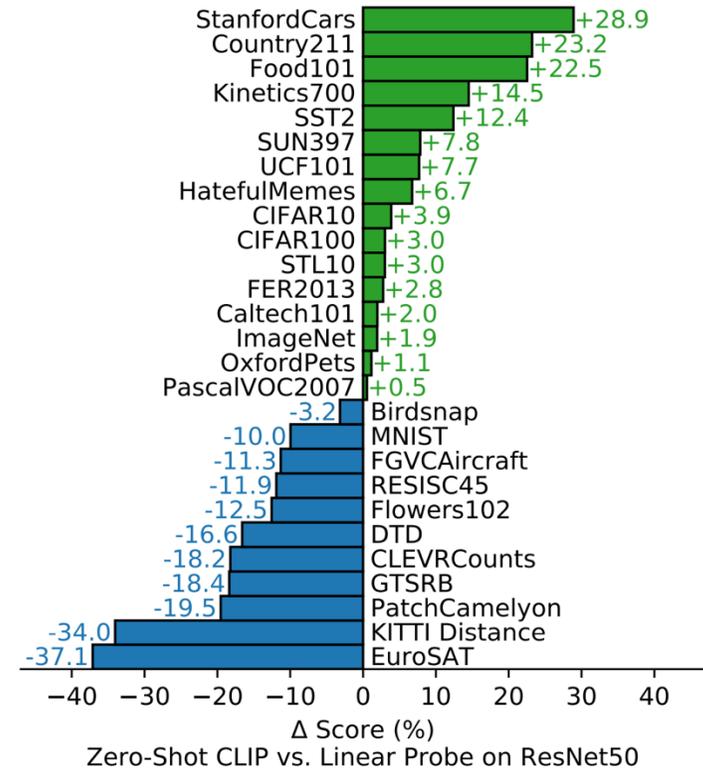ViT (Dosovitskiy et al. 2021)
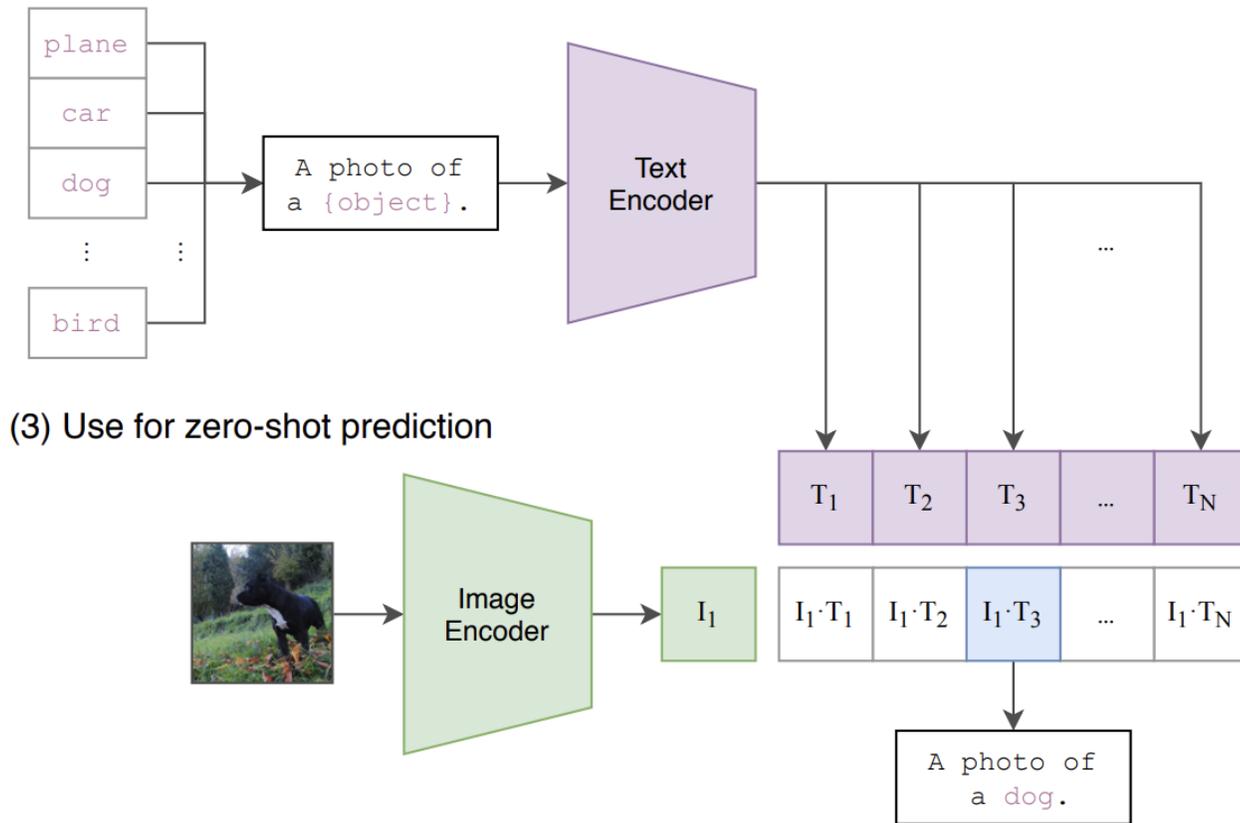
# Preliminary (CLIP)



CLIP (Radford et al. 2021)

$$\mathcal{L}_{ITC} = -\frac{1}{|\mathcal{B}|} \sum_{i \in \mathcal{B}} \log \frac{\exp\left(s\left(I_i, T_i\right)/\tau\right)}{\sum_{k \in \mathcal{B}} \exp\left(s\left(I_i, T_k\right)/\tau\right)}$$

$$-\frac{1}{|\mathcal{B}|} \sum_{j \in \mathcal{B}} \log \frac{\exp\left(s\left(I_j, T_j\right)/\tau\right)}{\sum_{k \in \mathcal{B}} \exp\left(s\left(I_k, T_j\right)/\tau\right)}$$

Contrastive Losses

# Preliminary (CLIP)

Can perform zero-shot inference using text.



CLIP (Radford et al. 2021)

# Extend to 3D and Text?



Mesh

Point Cloud

Voxel

Primitives

# Extend to 3D and text?

A mental cradle chair ...

Pretrained CLIP Models

Text Encoder ❄

Image Encoder ❄

1 View

PCD Encoder 🔥

Point Cloud

OpenShape / Uni3D

Contrastive Loss

Contrastive Loss

Contrastive Loss

Trainable Point Bert Model

$$l_i^{a \to b} = -\log \frac{\exp(\langle f_i^a, f_i^b \rangle)/\tau}{\Sigma_{k=1}^{N} \exp(\langle f_i^a, f_k^b \rangle)/\tau}$$
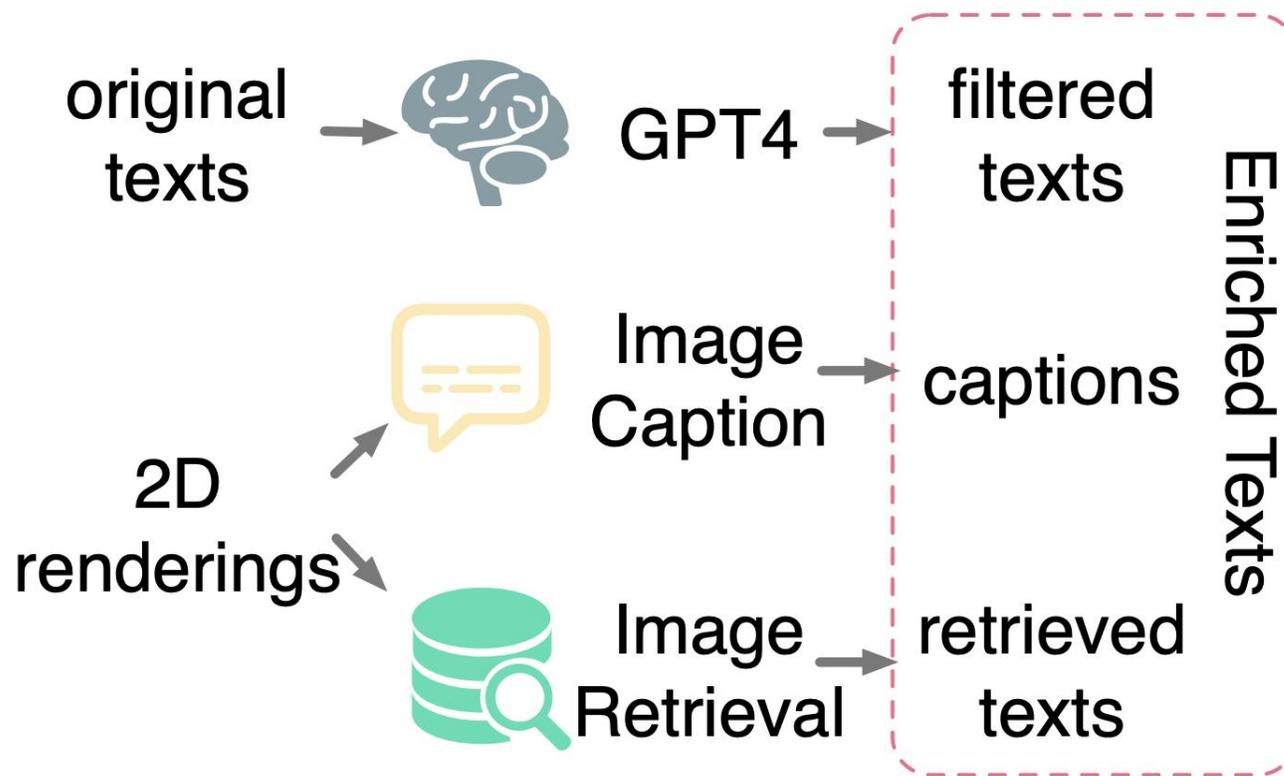
$$L_{CON} = \frac{1}{4N} \Sigma_{i=1}^{N} (l_i^{S \to T} + l_i^{T \to S} + l_i^{S \to I} + l_i^{I \to S})$$

# Data



(a) Ensemble Datasets

(b) Text Filtering & Enrichment

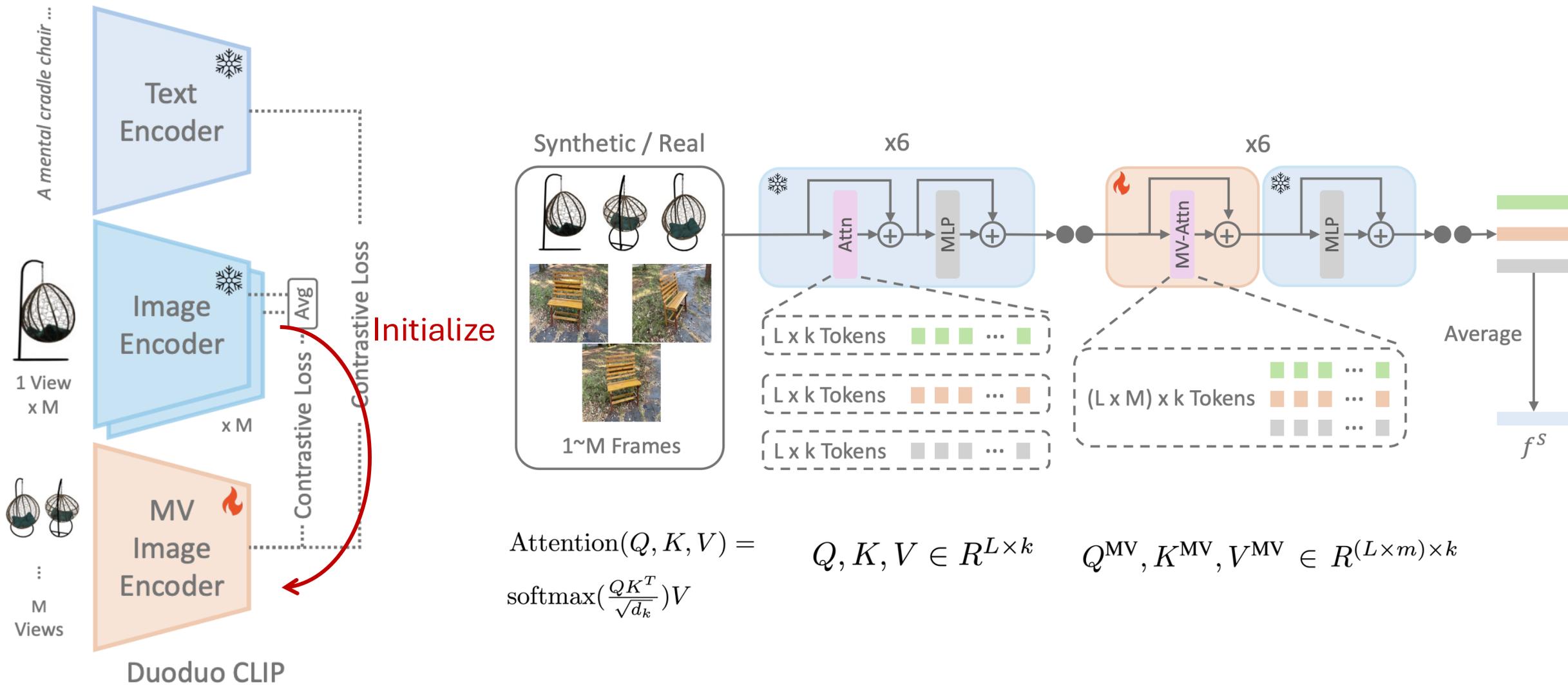OpenShape (Liu et al. 2023)

# More Efficient Training?

1. Point clouds are harder to acquire for real world objects
2. Domain gap between images and point clouds



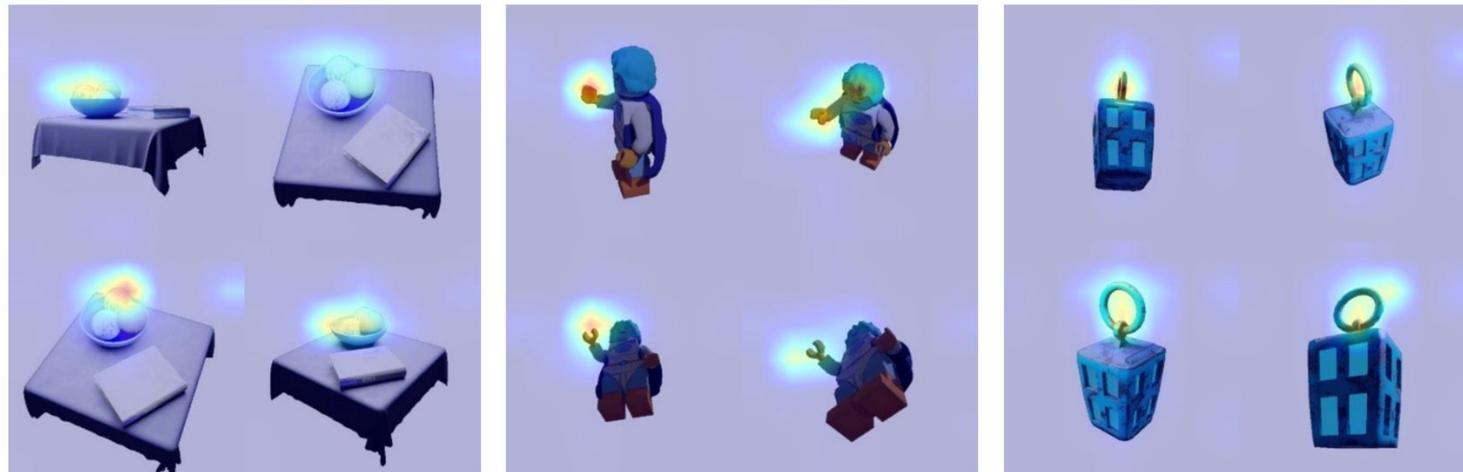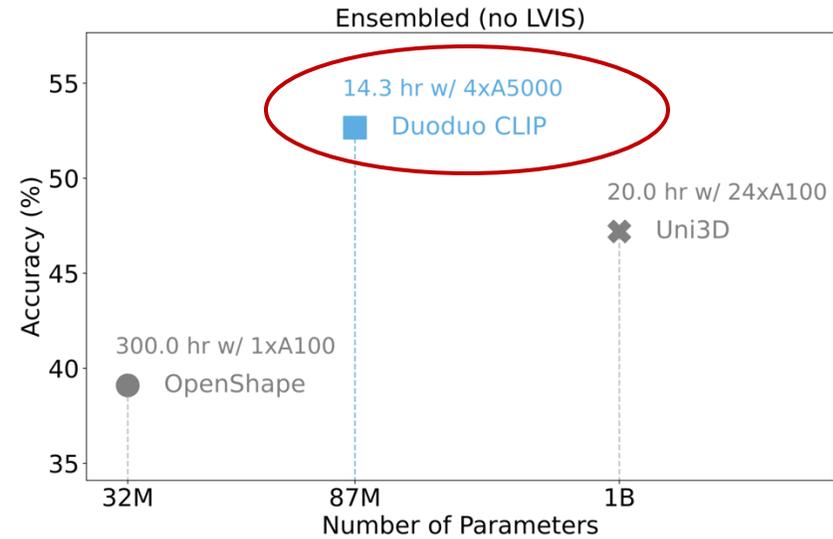1~M Frames

Use multi-view images instead!

# DuoduoCLIP



$$\text{Attention}(Q, K, V) = \text{softmax}(\frac{QK^T}{\sqrt{d_k}})V$$

$$Q, K, V \in R^{L \times k}$$

$$Q^{\text{MV}}, K^{\text{MV}}, V^{\text{MV}} \in R^{(L \times m) \times k}$$
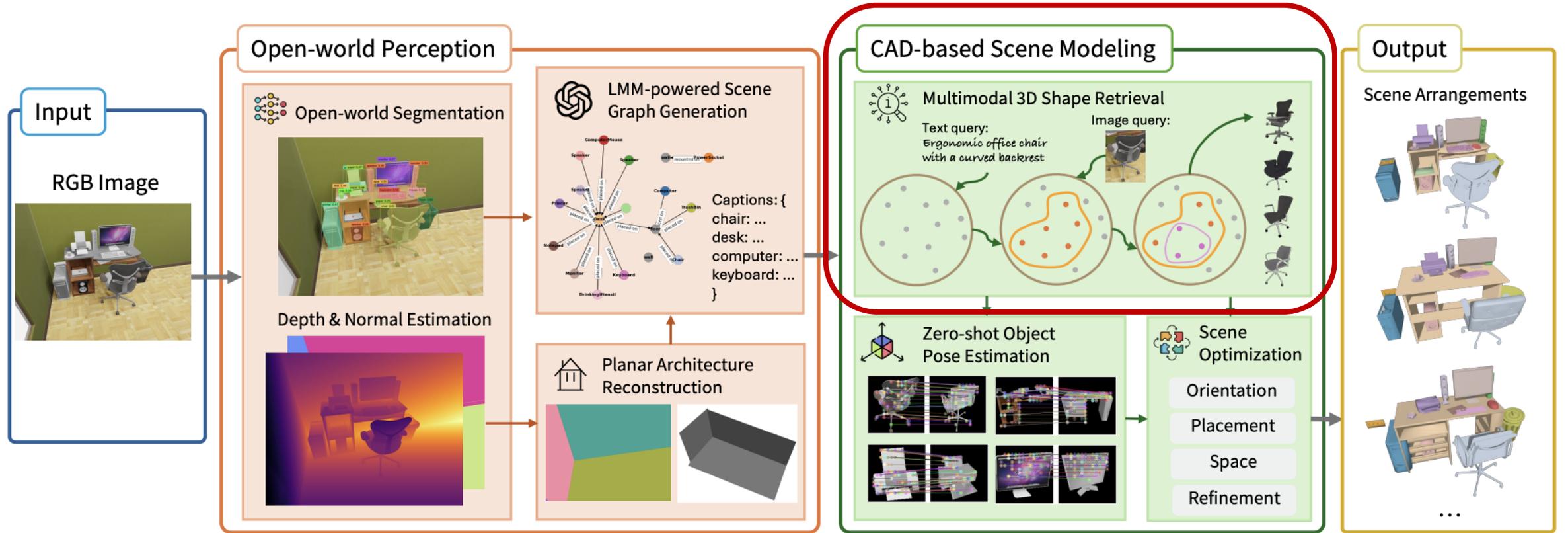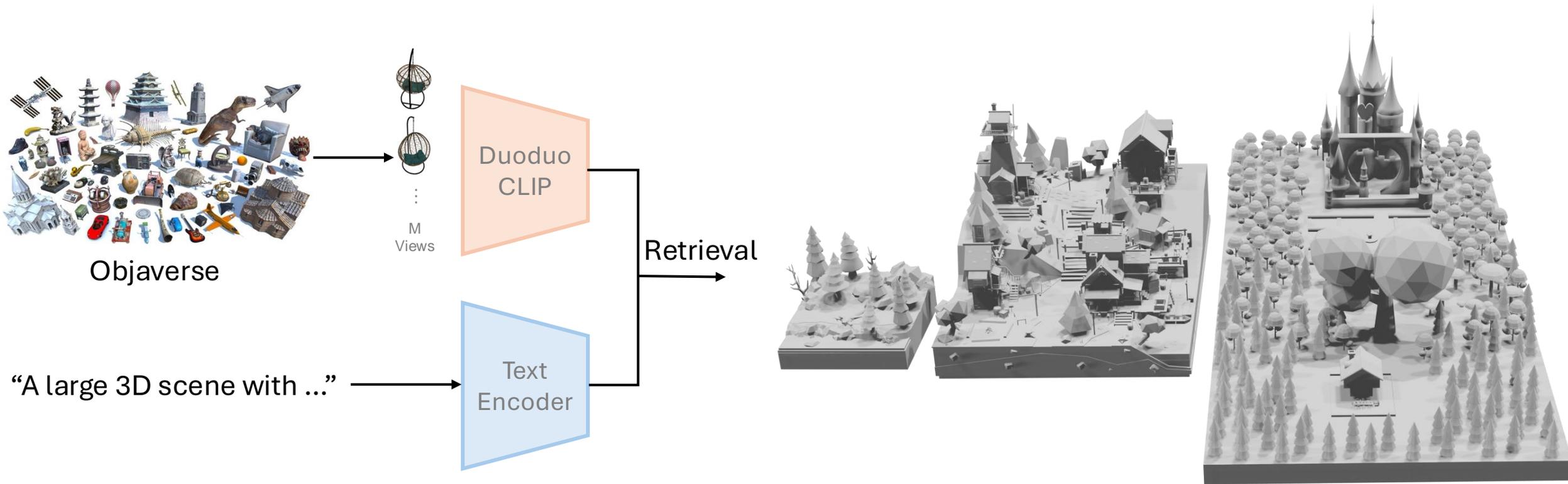
# DuoduoCLIP

# Applications (Digital Twin)



Diorama (Wu et al. 2025)

# Applications (Digital Twin)



Diorama (Wu et al. 2025)

# Applications (Dataset Filtering)



Objaverse

M Views

Duoduo CLIP

"A large 3D scene with ..."

Text Encoder

Retrieval

NuiScene (Lee et al. 2025)

# Applications (Dataset Filtering)

NuiScene (Lee et al. 2025)

# Check out the projects!



DuoduoCLIP



Diorama



NuiScene