Vision-Language Models Part I: Representation learning CLIP

Introduction

From zero-shot Transfer to representation learning

Remember: Transfer Learning

		Source Data (not directly related to the task)			
		labelled	unlabeled		
Target Data	labelled	Fine-tuning Multitask Learning	Not considered here		
	unlabeled	Domain adaptation-adversarial training Zero-shot learning	Not considered here		

- Source data: $(x^s, y^s) \longrightarrow$ Training data
- Target data: (Ø) usually same domain

Different tasks

Training time:

cat



dog ···

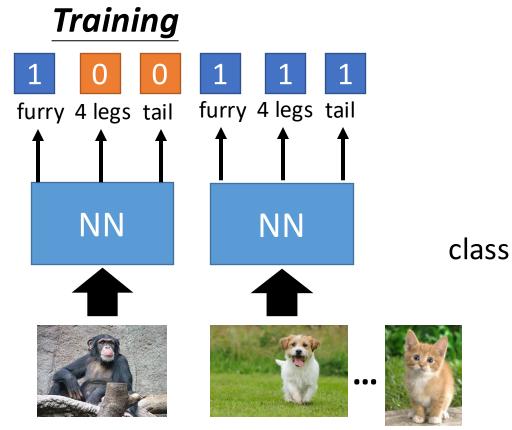
+ Class Information

Test time x^t :



=> Fish class!

Representing each class by its attributes

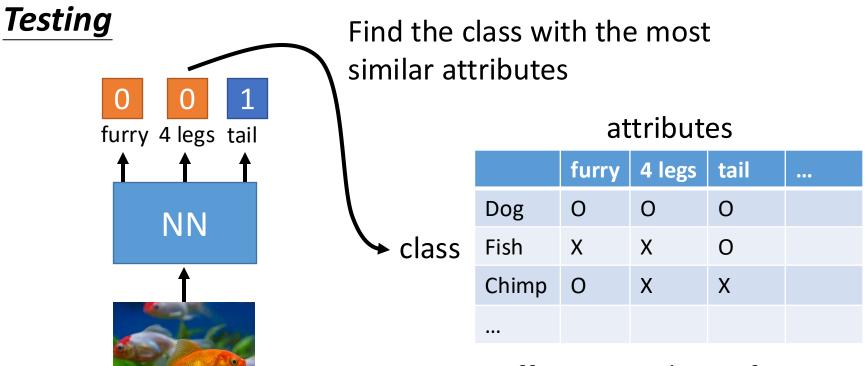


Database attributes

	furry	4 legs	tail	•••
Dog	0	0	0	
Fish	X	X	0	
Chimp	0	X	X	

sufficient attributes for one to one mapping

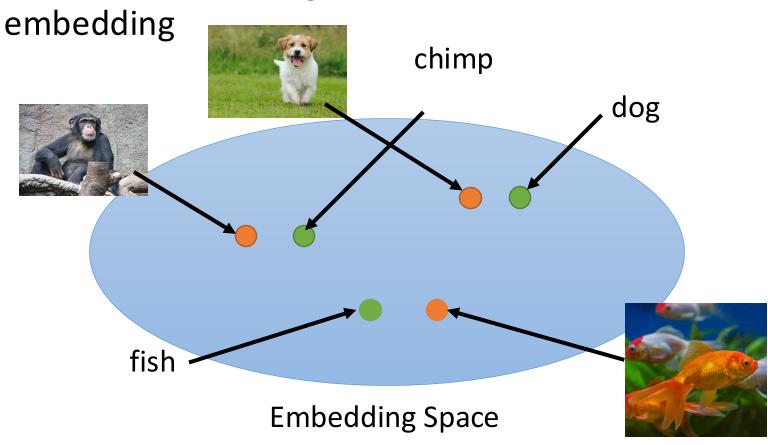
Representing each class by its attributes



sufficient attributes for one to one mapping

What if we don't have attribute database

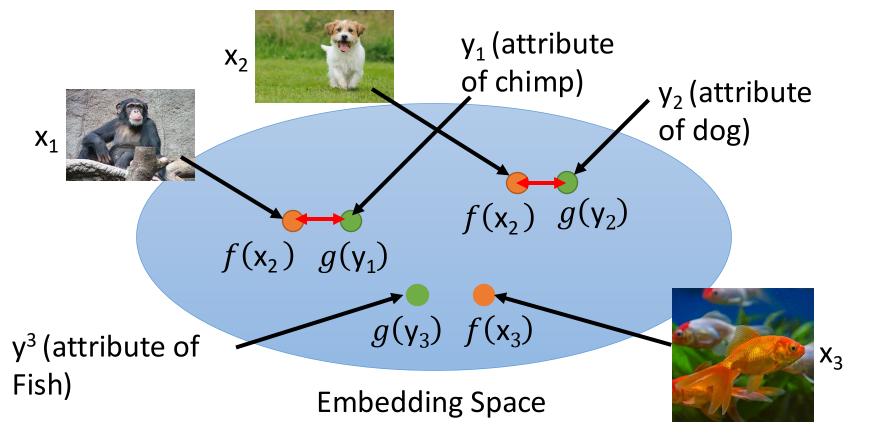
Attribute embedding + class (word name)



Attribute embedding

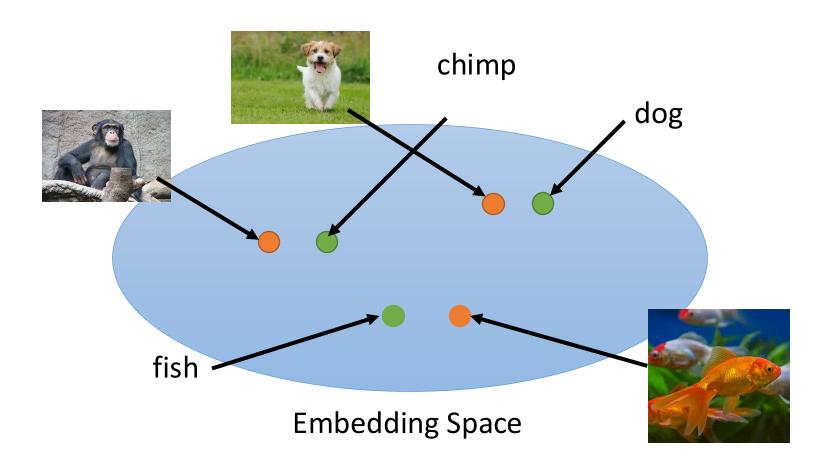
f(*) and g(*) can be NN. Training target:

 $f(x_n)$ and $g(y_n)$ as close as possible

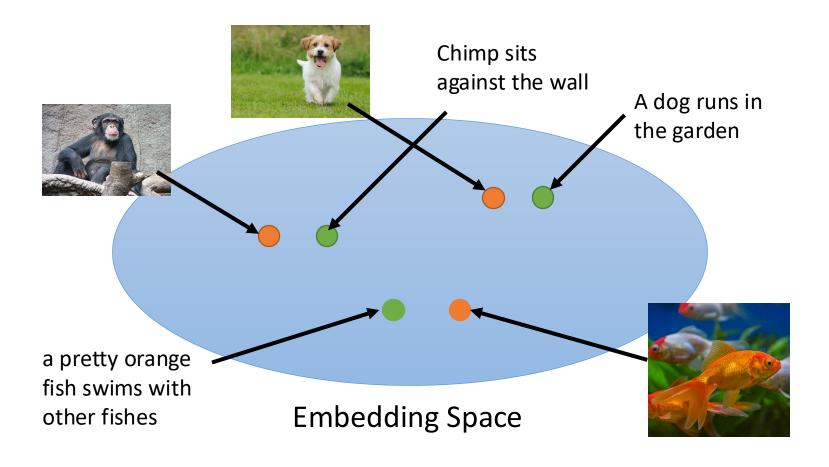


y_i are linked together by a class relationship (e.g. class name embedding as W2v)

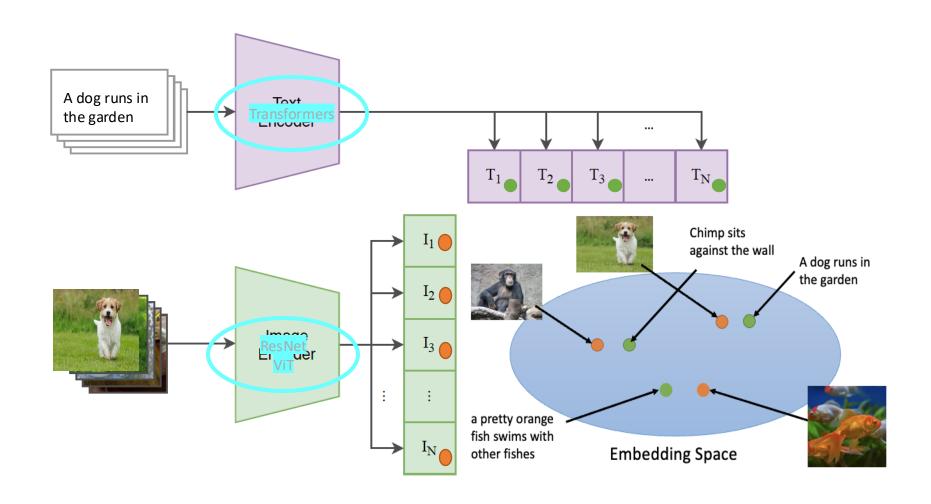
More on Vision-Language: Representation Learning



[Learning transferable visual models from natural language supervision. Radford/Sutskever ICML, 2021]



Dual architecture: Text encoder + Image encoder

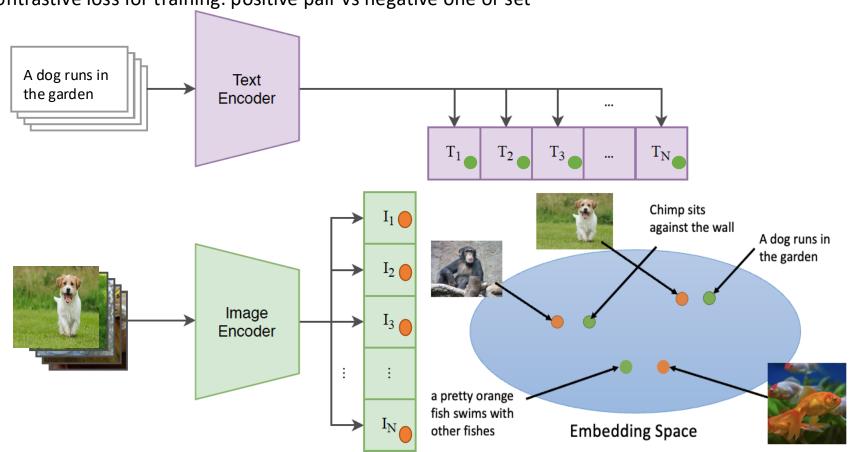


Learning strategy

Training set: $A = \{(\mathbf{I}_n, \mathbf{T}_n)\}_n$ of image/caption pairs (coherent!)

Massive Text+Image =400M pairs to train the model (from the Internet)

Contrastive loss for training: positive pair vs negative one or set

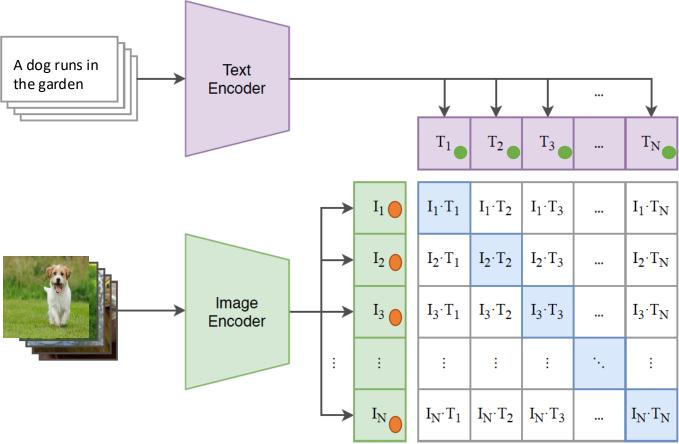


Learning strategy

Training set: $A = \{(\mathbf{I}_n, \mathbf{T}_n)\}_n$ of image/caption pairs (coherent!)

Massive Text+Image =400M pairs to train the model (from the Internet)

Contrastive loss for training: positive pair vs negative one or set



Learning strategy

Training set: $A = \{(\mathbf{I}_n, \mathbf{T}_n)\}_n$ of image/caption pairs (coherent!)

Massive Text+Image =400M pairs to train the model (from the Internet)

Contrastive loss for training: positive pair vs negative pair or more

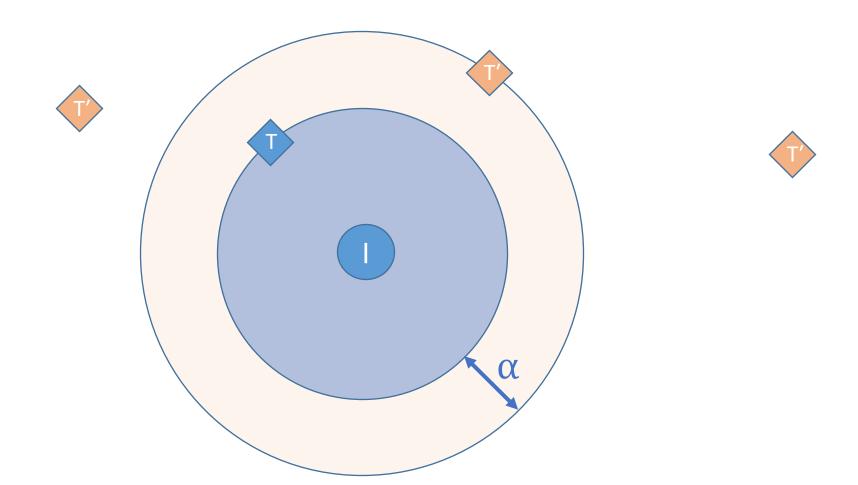
(contrastive) Triplet loss: A variant of the standard margin based loss (SVM)

- Triplet (I, T, T') (Batch = 3)
- Anchor: I (E.g image representation)
- Positive: T (E.g associated caption representation)
- Negative: T' (E.g contrastive caption representation)
- Margin parameter α

TripletLoss(I, T, T') =
$$\max\{0, \alpha + d(I, T) - d(I, T')\}$$

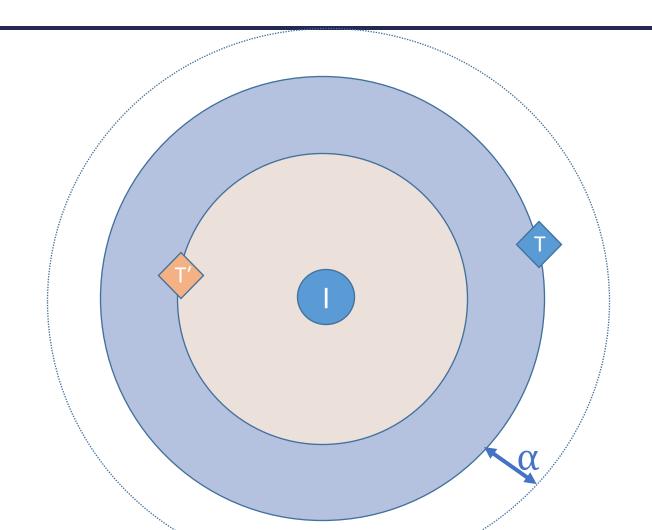
Learning strategy: triplet loss

TripletLoss(I, T, T') =
$$\max\{0, \alpha + d(I, T) - d(I, T')\}$$



Learning strategy: triplet loss

TripletLoss(I, T, T') =
$$\max\{0, \alpha + d(I, T) - d(I, T')\}$$



Learning strategy: triplet loss

Hard negative margin-based loss:

Loss for a **batch** $\mathcal{B} = \{(\mathbf{I}_n, \mathbf{T}_n)\}_{n \in B}$ of image/sentence pairs:

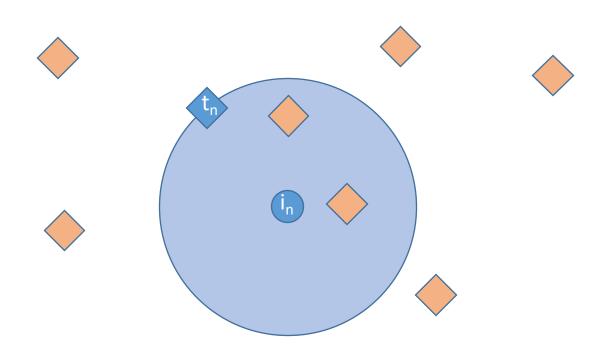
$$\mathcal{L}(\mathbf{\Theta}; \mathcal{B}) = \frac{1}{|B|} \sum_{n \in B} \begin{pmatrix} \max_{m \in C_n \cap B} \operatorname{loss}(\mathbf{I}_n, \mathbf{T}_n, \mathbf{T}_m) \\ + \max_{m \in D_n \cap B} \operatorname{loss}(\mathbf{T}_n, \mathbf{I}_n, \mathbf{I}_m) \end{pmatrix}$$

With C_n (resp. D_n) set of indices of caption (resp. image) unrelated to n-th element

Learning strategy: hard negative triplet loss

Mining hard negative contrastive example:

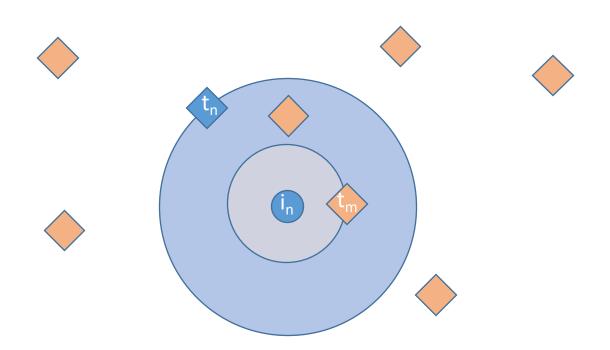
$$\mathcal{L}(\mathbf{\Theta}; \mathcal{B}) = \frac{1}{|B|} \sum_{n \in B} \begin{pmatrix} \max_{m \in C_n \cap B} \operatorname{loss}(\mathbf{I}_n, \mathbf{T}_n, \mathbf{T}_m) \\ + \max_{m \in D_n \cap B} \operatorname{loss}(\mathbf{T}_n, \mathbf{I}_n, \mathbf{I}_m) \end{pmatrix}$$



Learning strategy: hard negative triplet loss

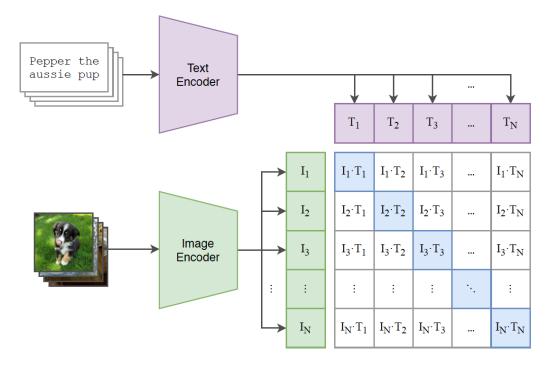
Mining hard negative contrastive example:

$$\mathcal{L}(\mathbf{\Theta}; \mathcal{B}) = \frac{1}{|B|} \sum_{n \in B} \begin{pmatrix} \max_{m \in C_n \cap B} \log (I_n, T_n, T_m) \\ + \max_{m \in D_n \cap B} \log (T_n, I_n, I_m) \end{pmatrix}$$



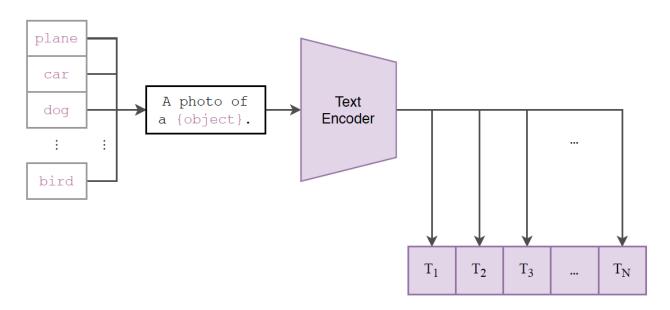
Massive Text+Image =400M pairs to train the model (from the Internet)

Contrastive loss for training

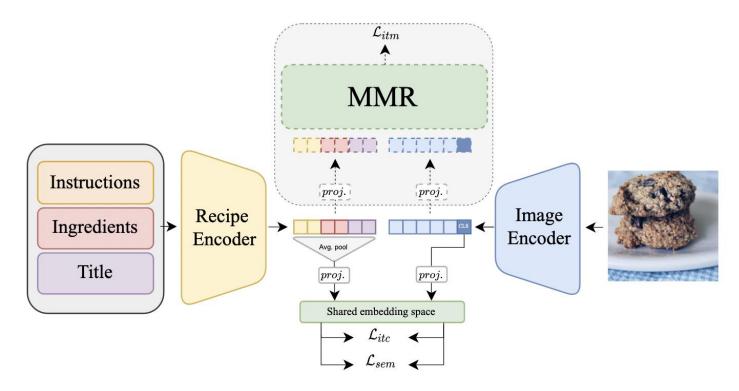


$$\mathcal{L}_{InfoNCE_{CLIP}} = -\sum_{i} \log \left(\frac{exp(\frac{sim(I_{i}, T_{i})}{\tau})}{\sum_{k=1}^{N} exp(\frac{sim(I_{i}, T_{k})}{\tau})} \right)$$

Pre-trained encoders = **dual encoders** (Text/Image) used for Zero-shot classifier, and other downstream tasks



A lot of variants



Title query	Ingredient query	Instruction query	
	1 cup Unsalted Butter,	Add sugar, cream, peppermint, and food coloring	
Mint Chocolate Chip Frosting.	2 Tablespoons Heavy Cream,	scoop the frosting and place on top of your cupcakes	
	2 drops Green Food Coloring, Chocolate	Source: Chocolate Cupcakes with Mint Chocolate Chip	
	1 broiler-fryer chicken, halved,	Place the halved chicken in a large, shallow container	
Honey-Grilled Chicken.	34 cup butter, melted,	Combine the remaining ingredients, stirring sauce well	
	14 cup honey	Grill chicken, skin side up	
	1/2 cup Kale,	Wash and cut kale off the stems	
The Best Kale Ever.	1 teaspoon Olive Oil,	Heat olive oil on medium heat and add garlic	
	1/4 teaspoons Red Pepper Flakes	Add in kale and red pepper flakes	

