

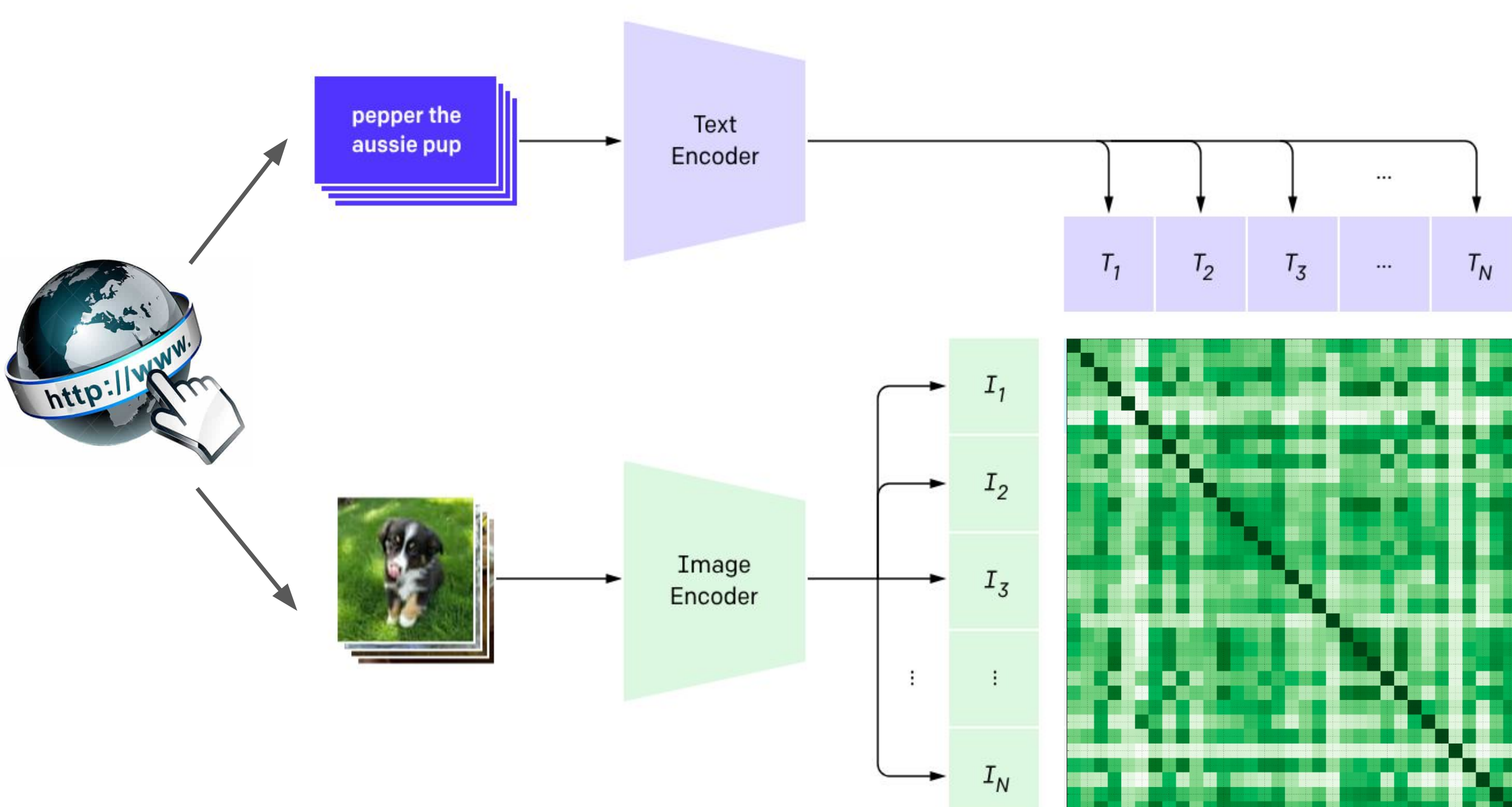


DIME-FM : Distilling Multimodal and Efficient Foundation Models

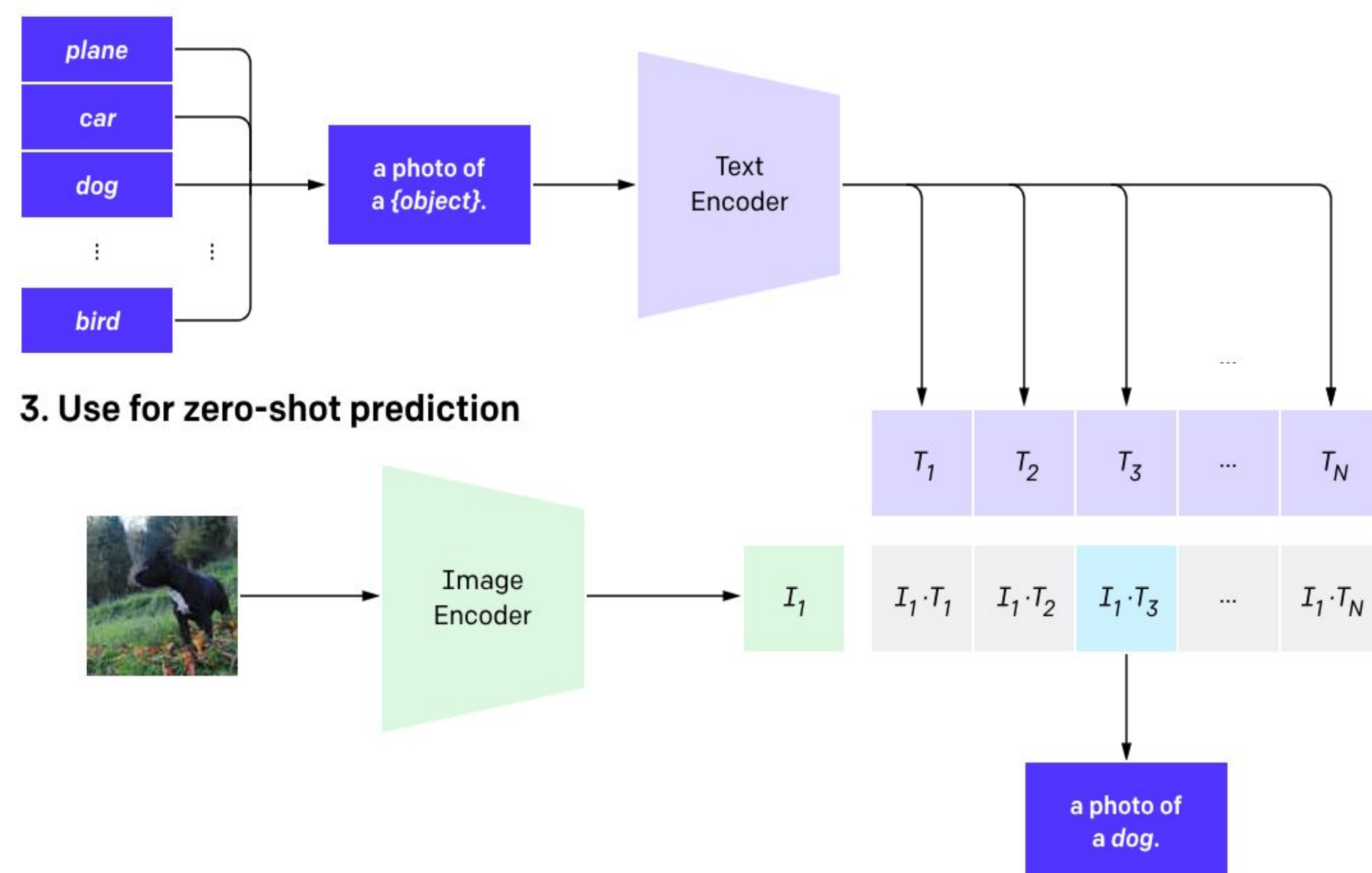
Ximeng Sun¹, Pengchuan Zhang², Peizhao Zhang², Hardik Shah², Kate Saenko^{1,2}, Xide Xia²
¹ Boston University, ² Meta AI

Background and Motivation

Contrastive Pretraining for Vision-Language Foundation models (such as CLIP)
 1. Contrastive pre-training



2. Create dataset classifier from label text



3. Use for zero-shot prediction

Superior Zero-Shot performance

Food101 Transfer to other dataset

guacamole (90.1%) Ranked 1 out of 101 labels

ImageNet-R (Rendition) Robustness to domain shifts

Siberian Husky (76.0%) Ranked 1 out of 200 labels

Real

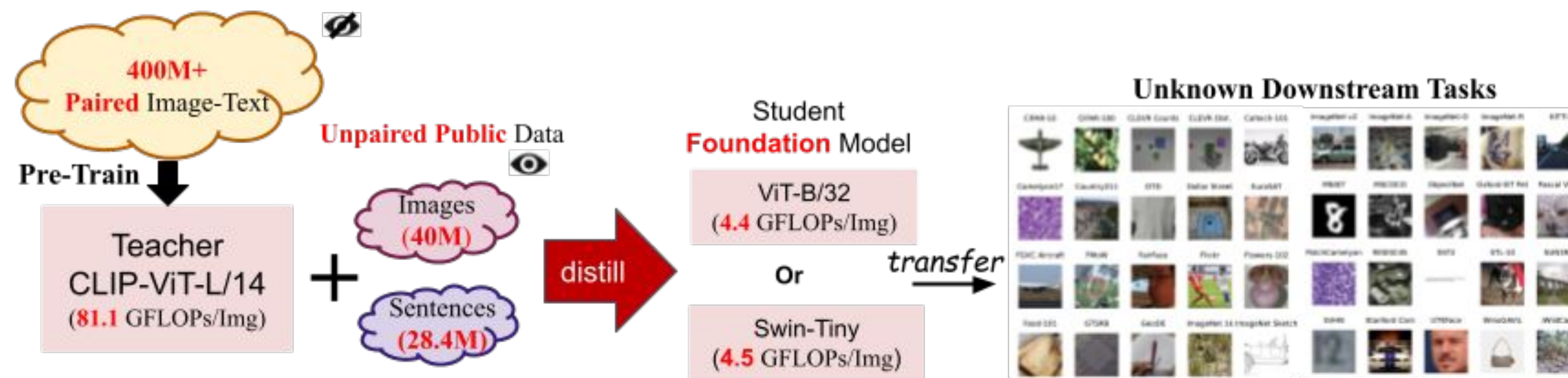
Domain shift

Problem: very expensive to train a CLIP-like model!!

- Heavy Data Consumption and Data Privacy**
 - 400M pretraining image-text pairs
 - private to OpenAI
- Expensive Training:**
 - batch-size = 32,768
 - largest ResNet model, RN50x64, took 18 days to train on 592 V100 GPUs
 - largest Vision Transformer took 12 days on 256 V100 GPUs



How to efficiently train a customizable CLIP-like model?



Our Remarkable Performance

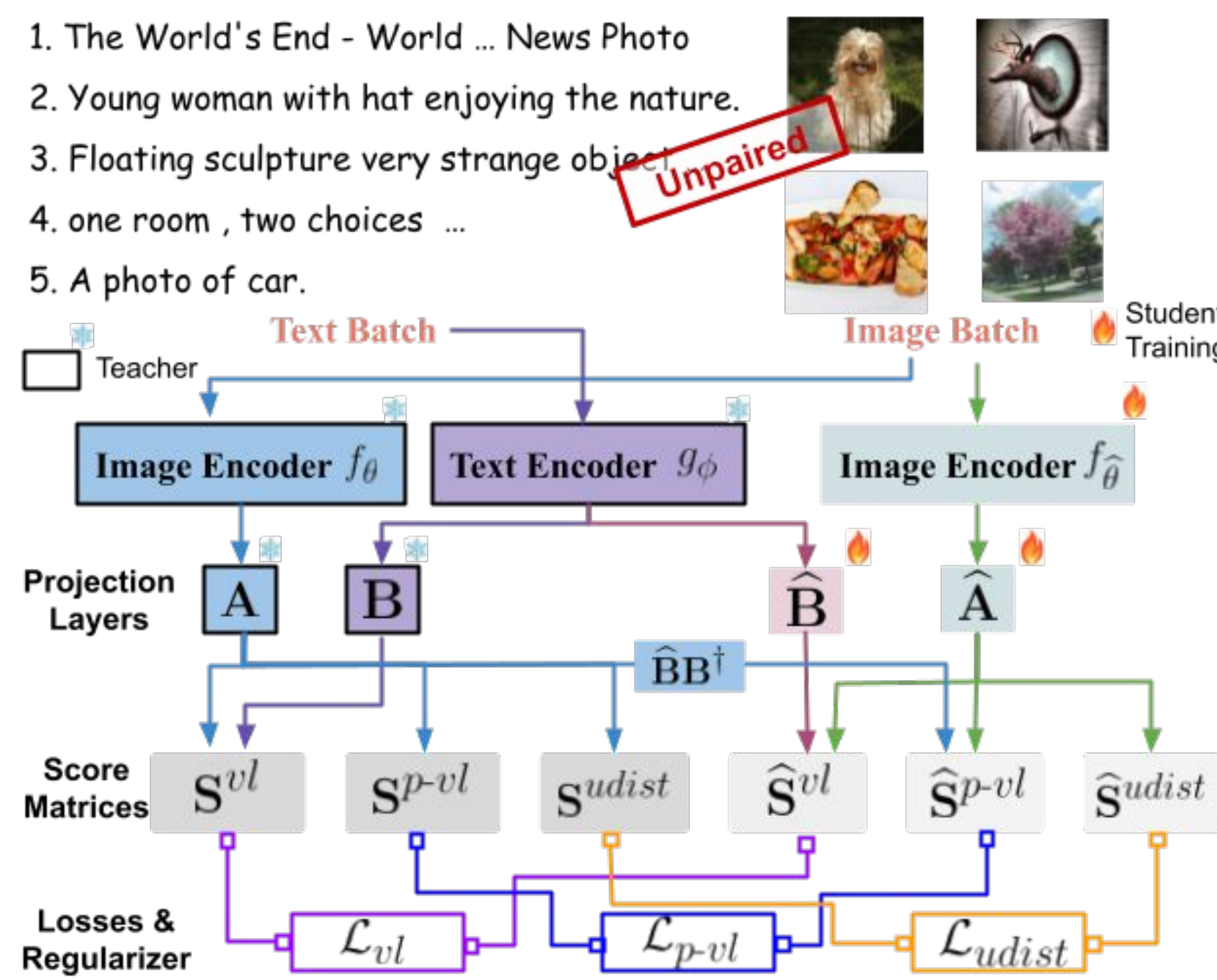
	ZS on IN-1K	ZS on ELEVATER	LP on ELEVATER	Robustness
OpenAI CLIP-ViT-B/32	63.4%	57.2%	78.2%	48.6%
Distill-ViT-B/32 (Constructed NLP Corpus)	66.5%	56.4%	79.2%	50.2%
Distill-ViT-B/32* (Caption Corpus)	64.8%	55.0%	78.6%	49.4%
Distill-ViT-B/32* (Task-Aware + Captions)	66.1%	57.7%	79.4%	50.4%

Dataset	Method	Zero-Shot		Linear Probing
		ELEVATER	IN-1K	ELEVATER
IN-21K	UniCL	27.2%	28.5%	74.8%
	UniCL*	40.9%	51.4%	75.3%
	Distill-UniCL*	45.6%	59.5%	76.2%
IN-21K + YFCC-14M	UniCL	37.1%	40.5%	77.1%
	UniCL*	44.6%	58.7%	75.4%
	Distill-UniCL*	47.6%	60.0%	76.6%

Question: What logits should we distill in the open-vocabulary image recognition task?

*** We find out the choice of text corpus is critical.

Key Contribution 1: Three Distillation Losses



We propose the Pseudo Text Corpus and Pseudo VL knowledge distillation loss

Key Contribution 2: Construct the visually-grounded Text Corpus

Algorithm 1: Constructing text corpus T

```

Input : image embeddings  $\mathcal{U}$  as defined in Eq.15. A large text corpus  $\mathcal{T}_{large}$ .
Output : Selected text corpus  $\mathcal{T}$ , and  $|\mathcal{T}| \approx |\mathcal{U}|$ 
1  $\mathcal{U}_{left} \leftarrow \mathcal{U}$ ,  $\mathcal{T}_{avail} \leftarrow \mathcal{T}_{large}$ ,  $\mathcal{T} \leftarrow \emptyset$ ,  $U_p = \infty$ 
2 while  $\mathcal{U}_{left} \neq \emptyset$  and  $|\mathcal{U}_{left}|/U_p < 0.95$  do
3    $U_p = |\mathcal{U}_{left}|$ ,  $Matched = dict()$ 
4   for  $u \in \mathcal{U}_{left}$ :
5     /* find the best text that matches the image */
6      $t(u) = \arg \max_{t \in \mathcal{T}_{avail}} s(u, B \cdot g_\phi(t))$  //
7    $Matched[u] = t(u)$ 
8 for  $u, t \in Matched.items()$ :
9   /* For all images matching to the same text, pick the first match */
10  if  $t \in \mathcal{T}_{avail}$ :
11     $\mathcal{U}_{left} \leftarrow \mathcal{U} \setminus \{u\}$ 
12     $\mathcal{T}_{avail} \leftarrow \mathcal{T}_{avail} \setminus \{t\}$ ,  $\mathcal{T}.add(t)$ 
    
```

