# Improved Probabilistic Image-Text Representations
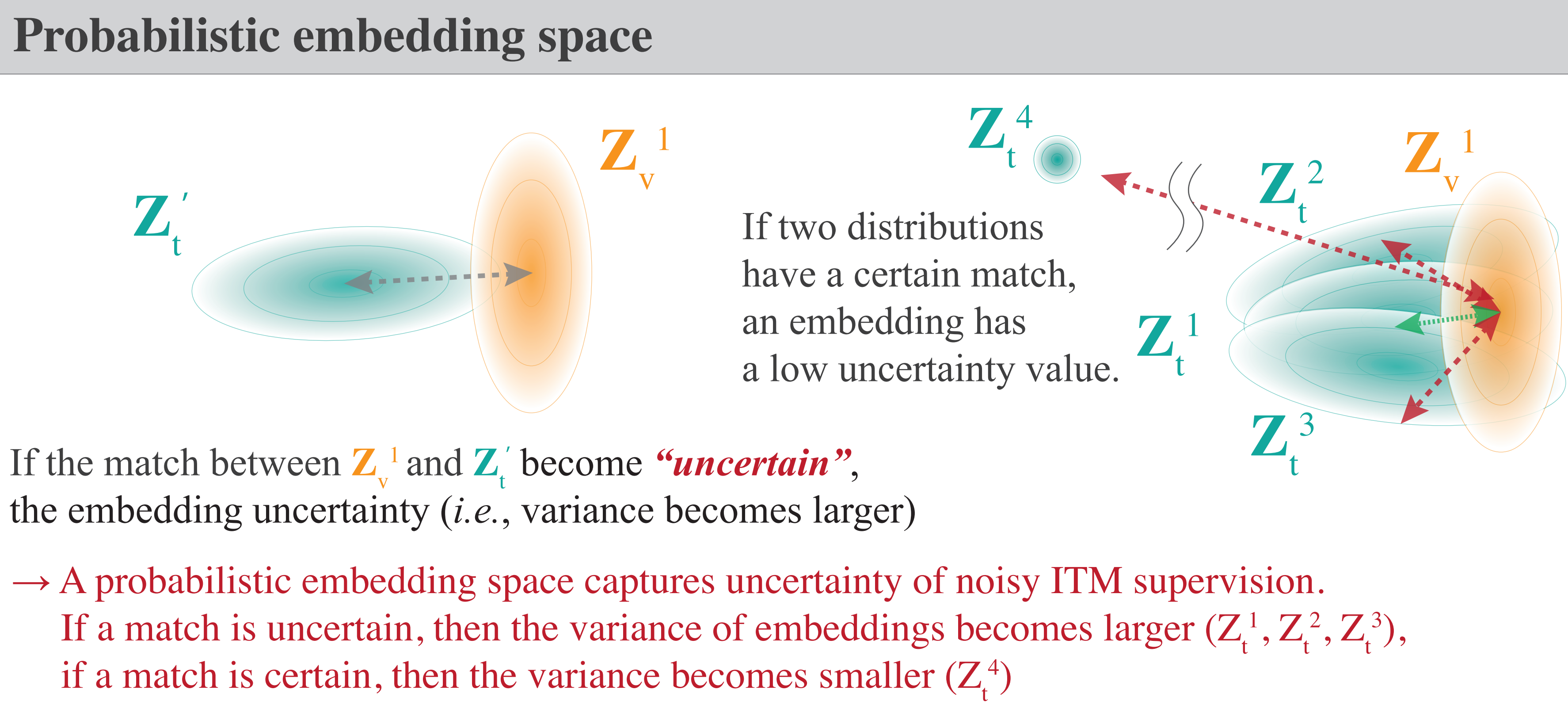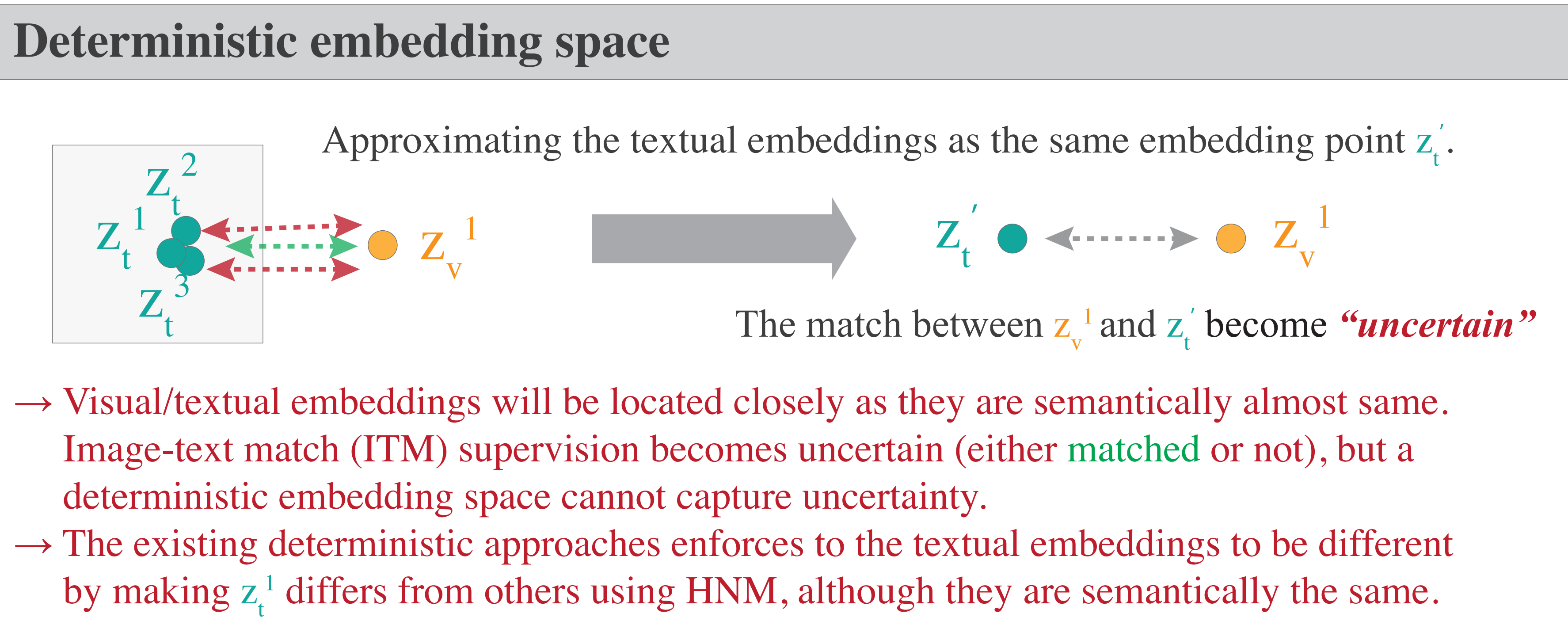
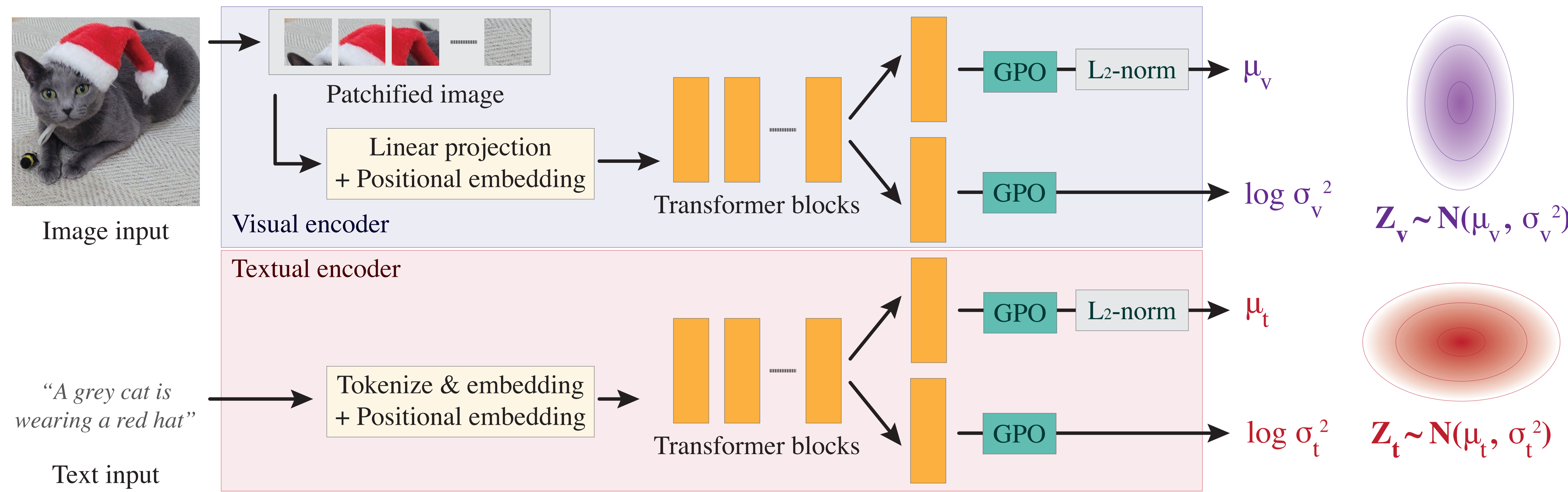Sanghyuk Chun — NAVER AI LAB

## Motivation: Inherent ambiguity in image-text matching problem

- - - → positive match
- - - → negative match

*A snowboarder is flying through the air doing stunts on his snowboard.*

*A person on a snow board flying in the air down the snow*

*A person on a snowboard jumping up in the air.*

| Image | Match | Caption |

**There are many false negatives (FNs) in dataset, and it naturally leads to ambiguous supervision to the model**

### Deterministic embedding space

Approximating the textual embeddings as the same embedding point $z_t'$.

The match between $z_v^1$ and $z_t^1$ become **"uncertain"**

→ Visual/textual embeddings will be located closely as they are semantically almost same. Image-text match (ITM) supervision becomes uncertain (either matched or not), but a deterministic embedding space cannot capture uncertainty.
→ The existing deterministic approaches enforces to the textual embeddings to be different by making $z_t^1$ differs from others using HNM, although they are semantically the same.

### Probabilistic embedding space

If two distributions have a certain match, an embedding has a low uncertainty value.

If the match between $Z_v^1$ and $Z_t^1$ become **"uncertain"**, the embedding uncertainty (i.e., variance becomes larger)

→ A probabilistic embedding space captures uncertainty of noisy ITM supervision. If a match is uncertain, then the variance of embeddings becomes larger ($Z_t^1, Z_t^2, Z_t^3$), if a match is certain, then the variance becomes smaller ($Z_t^4$)

## PCME++ An improved probabilistic VL model

Patchified image
Linear projection + Positional embedding
Image input
Visual encoder
Transformer blocks
GPO → L2-norm → $\mu_v$
GPO → $\log \sigma_v^2$
$Z_v \sim N(\mu_v, \sigma_v^2)$

"A grey cat is wearing a red hat"
Text input
Textual encoder
Tokenize & embedding + Positional embedding
Transformer blocks
GPO → L2-norm → $\mu_t$
GPO → $\log \sigma_t^2$
$Z_t \sim N(\mu_t, \sigma_t^2)$

How can we train a proper probabilistic embedding space? In other words, how to measure a distance between two distributions? Here, we expect three properties for the probabilistic embeddings:

1) There exists a proper probabilistic distance between two distributions
2) If the match between two distributions is certain, then the distributions should have small variances
3) If the match between two distributions is uncertain, then the variances should be large.

Popular distances such as Wasserstien dist, KL divergence cannot satisfy (2).

Monte-Carlo approximimation of matching probability (Chun et al. 2021)

$$p_\theta(m|x_\alpha, x_\beta) \approx \frac{1}{J^2} \sum_j^J \sum_{j'}^J s(-a\|z_\alpha - z_\beta\|_2 + b)$$
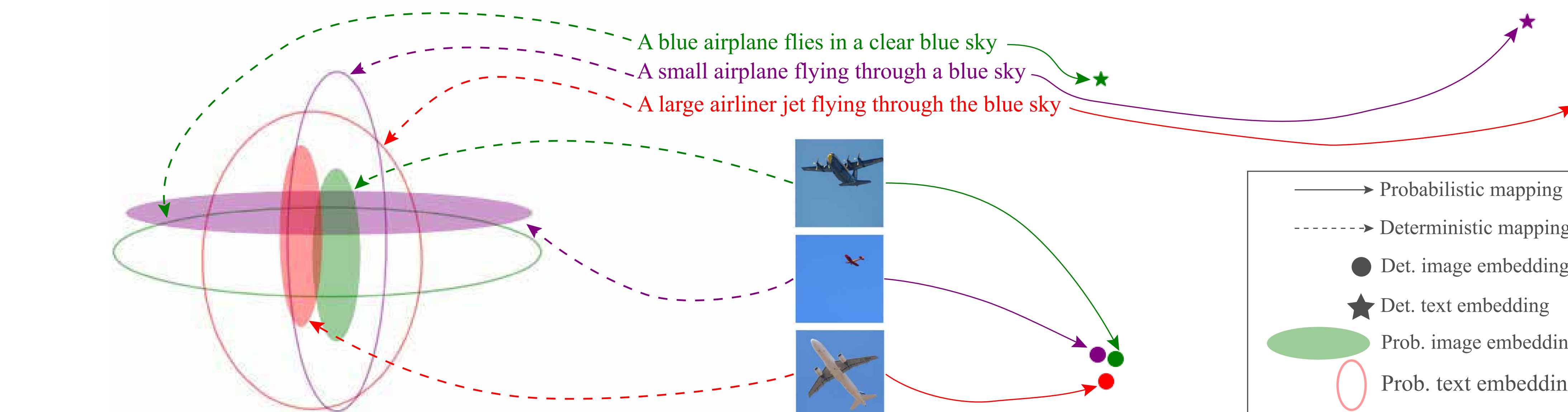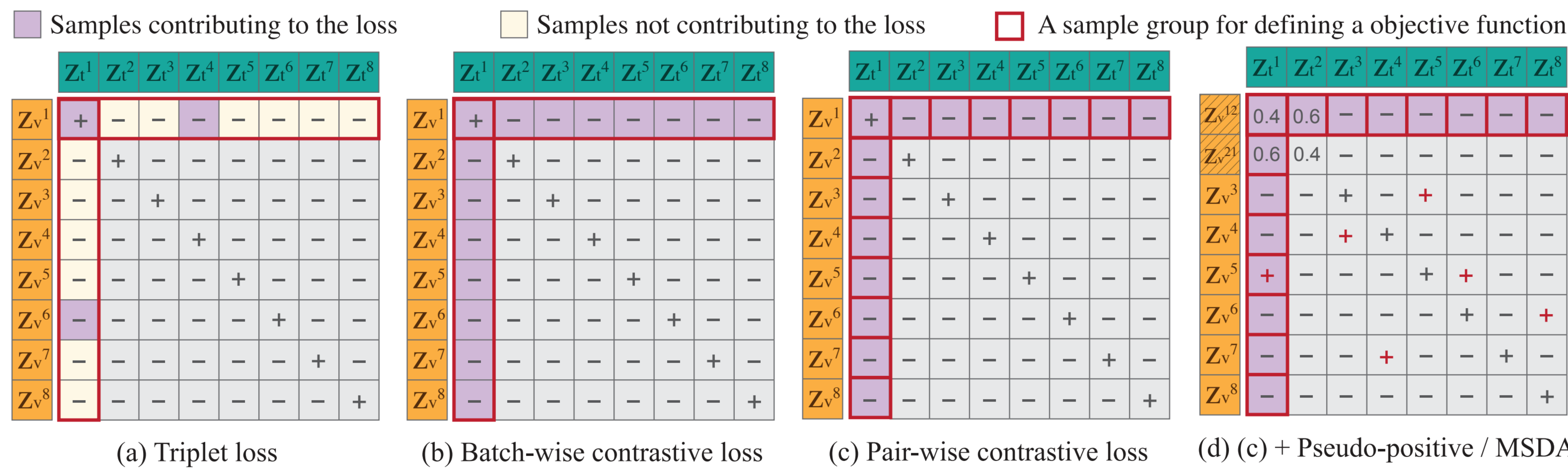
It satisfies all conditions, but the MC sampling itself has flaws in computation and accuracy

### Improvement 1: Closed-form Sampled Distance (CSD)

$$d(Z_v, Z_t) = \mathbb{E}_{Z_v, Z_t}\|Z_v - Z_t\|_2^2 = \|\mu_v - \mu_t\|_2^2 + \|\sigma_v^2 + \sigma_t^2\|_1$$

$$\mathcal{L}_{match} = -m_{vt}\log \text{sigmoid}(-a \cdot d(Z_v, Z_t) + b) - (1 - m_{vt})\log \text{sigmoid}(a \cdot d(Z_v, Z_t) - b)$$

### Improvement 2: Pseudo-Positive (PP) and Mixed data sample augmentation (MSDA)

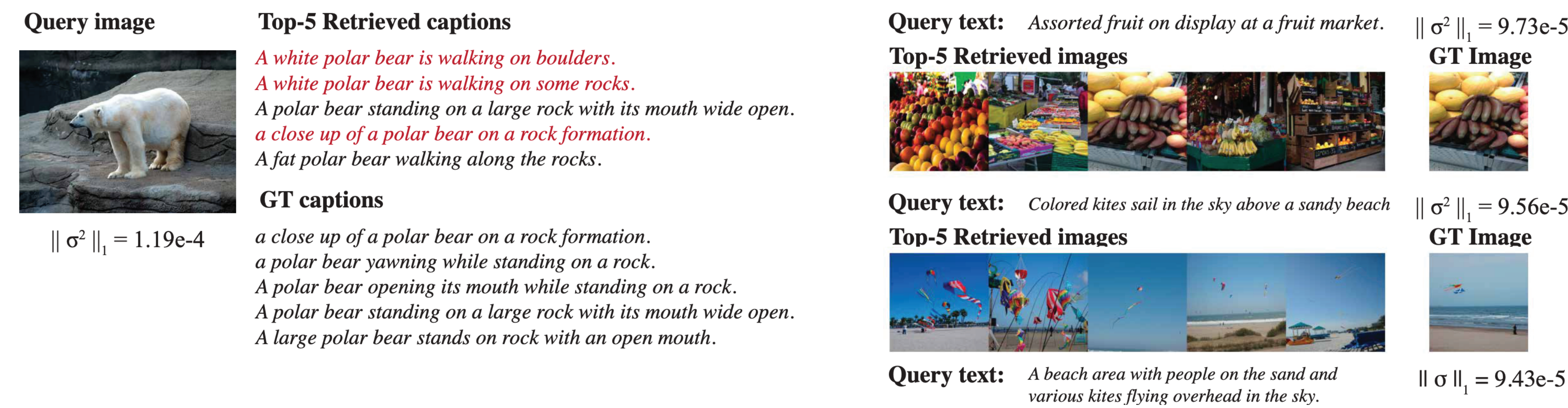☐ Samples contributing to the loss   ☐ Samples not contributing to the loss   ☐ A sample group for defining a objective function

(a) Triplet loss   (b) Batch-wise contrastive loss   (c) Pair-wise contrastive loss   (d) (c) + Pseudo-positive / MSDA

A blue airplane flies in a clear blue sky
A small airplane flying through a blue sky
A large airliner jet flying through the blue sky

— Probabilistic mapping
- - Deterministic mapping
● Det. image embedding
★ Det. text embedding
● Prob. image embedding
◯ Prob. text embedding

## Experimental results

### COCO Caption results

| Backbone | Method | Prob? | ECCV Caption mAP@R | ECCV Caption R-P | R@1 | COCO 5K R@1 | COCO RSUM |
|---|---|---|---|---|---|---|---|
| ViT-B/32 (151M) | CLIP ZS[†] | ✗ | 26.8 | 36.9 | 67.1 | 40.3 | 471.9 |
| | VSE∞ | ✗ | 40.0 | 49.5 | 83.1 | 55.2 | 536.5 |
| | P2RM | ✗ | 39.0 | 48.7 | 82.0 | 51.7 | 530.2 |
| | DAA | ✗ | 39.2 | 49.0 | 82.0 | 52.9 | 530.9 |
| | InfoNCE | ✗ | 39.0 | 48.7 | 81.7 | 53.0 | 532.6 |
| | PCME | ✓ | 39.1 | 48.9 | 81.4 | 53.0 | 532.0 |
| | PCME++ ($\mu$ only) | ✓ | 39.5 | 49.1 | 82.7 | 55.2 | 536.2 |
| | PCME++ | ✓ | 40.1 | 49.7 | 83.1 | 55.1 | 537.0 |
| | PCME++ (SWA) | ✓ | 40.2 | 49.8 | 82.9 | 55.2 | 537.3 |
| ViT-L/14 (428M) | CLIP ZS[†] | ✗ | 28.0 | 37.8 | 72.2 | 46.4 | 491.6 |
| | VSE∞ | ✗ | 20.2 | 31.5 | 46.2 | 22.7 | 424.3 |
| | InfoNCE | ✗ | 35.6 | 45.8 | 75.6 | 45.9 | 520.6 |
| | PCME | ✓ | 41.2 | 50.3 | 86.0 | 61.9 | 550.4 |
| | PCME++ | ✓ | 42.1 | 50.8 | 88.8 | 64.3 | 554.7 |

### COCO Noisy correspondence results

| Noise | Method | ECCV Caption mAP@R | ECCV Caption R-P | R@1 | COCO 5K R@1 | COCO RSUM |
|---|---|---|---|---|---|---|
| 20% | VSE∞ | 37.0 | 46.3 | 79.7 | 51.8 | 518.6 |
| | DAA | 6.7 | 12.5 | 18.5 | 6.0 | 212.8 |
| | InfoNCE | 35.9 | 46.3 | 76.1 | 45.8 | 514.6 |
| | PCME | 37.6 | 47.6 | 79.2 | 48.7 | 520.7 |
| | PCME++ (ours) | 37.7 | 47.6 | 80.0 | 50.4 | 524.6 |
| | NCR[†] | 35.9 | 46.0 | 78.0 | 48.8 | 518.6 |
| | DECL[†] | - | - | - | 49.4 | 518.2 |
| 50% | VSE∞ | 18.0 | 28.5 | 43.7 | 19.1 | 394.1 |
| | InfoNCE | 33.6 | 44.1 | 73.0 | 41.4 | 499.5 |
| | PCME | 35.2 | 45.5 | 75.7 | 44.4 | 508.0 |
| | PCME++ (ours) | 35.7 | 45.8 | 76.3 | 45.5 | 511.0 |
| | NCR[†] | 34.0 | 44.3 | 75.1 | 45.5 | 508.5 |

**Observation 1:** When scaling-up backbone, det. methods are overfitted, but prob. methods are not.
**Observation 2:** PCME++ successfully handles NC although it is not designed for tackling NC.

Query image
Top-5 Retrieved captions
*A white polar bear is walking on boulders.*
*A white polar bear is walking on some rocks.*
*A polar bear standing on a large rock with its mouth wide open.*
*a close up of a polar bear on a rock formation.*
*A fat polar bear walking along the rocks.*
GT captions
*a close up of a polar bear on a rock formation.*
*a polar bear yawning while standing on a rock.*
*A polar bear opening its mouth while standing on a rock.*
*A polar bear standing on a large rock with its mouth wide open.*
*A large polar bear stands on rock with an open mouth.*
$\|\sigma^2\|_1 = 1.19e\text{-}4$

Query text: *Assorted fruit on display at a fruit market.*
Top-5 Retrieved images
GT Image
$\|\sigma^2\|_1 = 9.73e\text{-}5$
Query text: *Colored kites sail in the sky above a sandy beach*
Top-5 Retrieved images
GT Image
$\|\sigma^2\|_1 = 9.56e\text{-}5$
Query text: *A beach area with people on the sand and various kites flying overhead in the sky.*
$\|\sigma\|_1 = 9.43e\text{-}5$

### Large-scale pre-training with PCME++

| Model | Prompts | Top-1 ImageNet Zero-shot acc |
|---|---|---|
| CLIP | "A photo of { · }" | 31.85 |
| | All 80 prompts | 35.50 |
| PCME++ | "A photo of { · }" | 30.43 |
| | All 80 prompts | 34.22 |
| | Top-K certain prompts | 34.22 |
| | Best top-K for each class | 41.82 |

### ViT-B/16 trained on CC3M,12M,RedCaps

**Bookcase**
a close-up photo of a {},
a close-up photo of the {}

**Jack-o'-lantern**
a photo of the clean {}. art of the {}. a cropped photo of the {}. a close-up photo of a {}. a photo of a clean {}. a cropped photo of a {}

**Rapeseed**
a sculpture of the {}. a {} in a video game. a sculpture of a {}. art of the {}. the {} in a video game. a tattoo of the {}. the plushie {}. a tattoo of a {}. a drawing of a {}. a drawing of the {}. a sketch of the {}. a close-up photo of the {}. art of a {}. a photo of a clean {}. a plushie {}. a close-up photo of a {}. a photo of the clean {}. a rendering of the {}. a photo of the large {}. a rendering of a {}. a sketch of a {}. a cropped photo of the {}. a rendition of a {}. graffiti of a {}. a rendition of the {}. a photo of the small {}. a photo of one {}. a photo of the dirty {}. ... 24 more prompts

### Check out more details in...

Paper   Code

### Check out my other related works!

PCME (CVPR'21)   ECCV Caption (ECCV'22)