

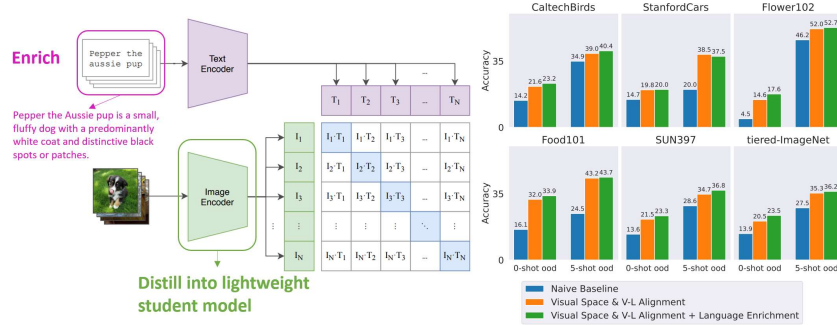
# Distilling Large Vision-Language Model with Out-of-Distribution Generalizability

Xuanlin Li\*, Yunhao Fang\*, Minghua Liu, Zhan Ling, Zhuowen Tu, Hao Su  
University of California San Diego



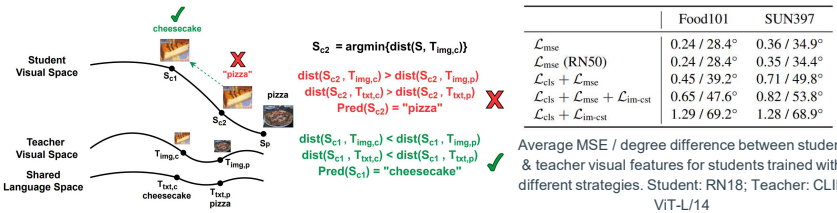
## Overview

- Goal: Distill visual representations in large vision-language teachers into lightweight student models, while allowing students to possess **strong open-vocabulary generalization ability towards out-of-distribution (OOD) concepts**.
  - Potential for deployment on mobile / IoT devices and robotics scenarios.
  - Experiments are conducted on small-to-medium datasets, allowing for fast research & development (R&D) cycles and is practical for lower-data regimes like 3D and robotics.
- Main Findings: **2 principles** to enhance student's OOD generalizability:
  - Better imitate teacher's visual representation space**, and carefully promote **better vision-language alignment coherence** with the teacher.
  - Enrich teacher's language representations** with more finegrained & meaningful attributes to effectively distinguish between different labels, both during distillation and inference. **This can be accomplished by prompting LLMs like ChatGPT.**

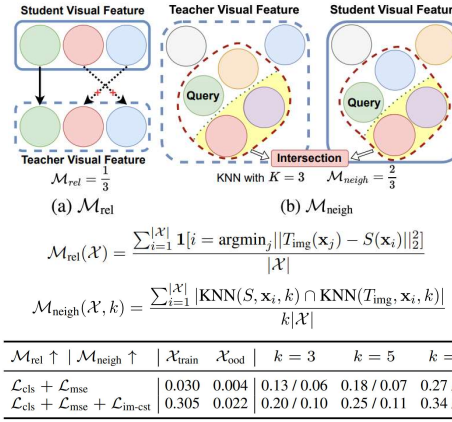


## Teacher-Student Visual Space Alignments

- Precisely matching teacher and student's high-dimensional visual spaces is inherently challenging (with high MSE feature matching loss).
- In this case, **minimizing student-teacher visual feature distance does NOT yield the best OOD generalization ability for students.**



- However, teacher-student coherence in **local visual space structures** and **relative visual feature relationships** is crucial for student OOD generalization, as they *implicitly* enhance vision-language alignments.
- We find that **Contrastive losses** like  $L_{im-cst}$  are especially helpful for student OOD generalization, as it **enables much better student-teacher proximity in both local and relative visual feature structures**.



## Teacher-Student Vision-Language Alignments

- Explicitly** enhance teacher-student vision-language alignment coherence.
- Carefully** preserve teacher's vision-language alignment structure by **accounting for teacher misalignments** (through e.g.,  $L_{vprox}$ ).

$$\mathcal{L}_{vprox}(\mathbf{x}, k) = I(\mathbf{x}) \cdot \mathcal{D}_{KL}(P_{T, \text{topk}}(\cdot|\mathbf{x}) || P_{S, \text{topk}}(\cdot|\mathbf{x}))$$

$$I(\mathbf{x}) = \mathbf{1}[\arg\max_y P_T(y|\mathbf{x}) = \text{label}(\mathbf{x})]$$

$$P_{\cdot, \text{topk}}(y|\mathbf{x}) = \frac{\mathbf{1}_{y \in Y_{\text{topk}}}}{\sum_{y \in Y_{\text{topk}}} P_{\cdot}(y|\mathbf{x})}; Y_{\text{topk}} = \arg\text{topk}_y P_T(y|\mathbf{x})$$

| $M_{valign} \downarrow$  | $k=2$       | $k=3$       | $k=5$       |
|--|-------------|-------------|-------------|
| $\mathcal{L}_{cls} + \mathcal{L}_{mse}$  | 0.20 / 0.50 | 0.68 / 1.45 | 2.67 / 4.73 |
| $\mathcal{L}_{cls} + \mathcal{L}_{mse} + \mathcal{L}_{im-cst}$                       | 0.18 / 0.43 | 0.62 / 1.3  | 2.52 / 4.24 |
| $\mathcal{L}_{cls} + \mathcal{L}_{mse} + \mathcal{L}_{im-cst} + \mathcal{L}_{vprox}$ | 0.17 / 0.39 | 0.59 / 1.17 | 2.17 / 4.20 |

$$M_{valign}(\mathcal{X}, k) = \frac{\sum_{i=1}^{|\mathcal{X}|} \#reverse\_pairs(arrS(i, k))}{|\mathcal{X}|}$$

$$arrS(i, k) = [|S(x_i) - T_{img}(\ell(y_j))|_2]_{j \in \mathcal{I}(i, k)}$$

$$\mathcal{I}(i, k) = \arg\text{topk}(|-T_{img}(x_i) - T_{img}(\ell(y_j))|_2)_{j=1}^{|\mathcal{Y}|}$$

## Language Representation Enrichment

- Prompt ChatGPT to enrich label descriptions**: "Use a single sentence to describe the appearance and shape of {cls}. Only describe the shape and appearance."
- ChatGPT-enriched descriptions are significantly more informative than e.g., auxiliary captions generated by OFA**, supplying comprehensive details to distinguish finegrained categories.



Original Description:

"A photo of an Acura Integra Type R 2001"

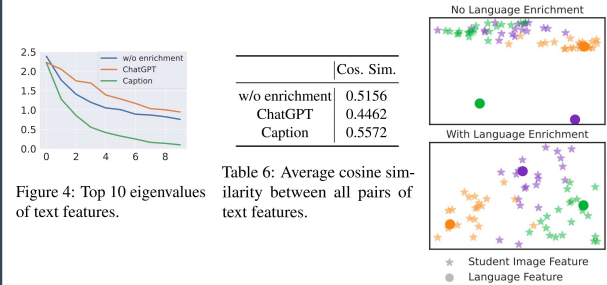
ChatGPT-Enriched:

"A photo of an Acura Integra Type R 2001. The 2001 Acura Integra Type R features a compact and sporty design with a sleek, aerodynamic body, sharp angles, and a distinctive rear spoiler."

Auxiliary Caption from OFA:

"A white car is parked in a field."

- ChatGPT-enriched language representation space confers more **independent & meaningful attributes to distinguish labels**, allowing OOD text features to be more precisely aligned with image features.



## Applications (e.g., Robotics)



|                          | $\mathcal{X}_{id}$ | $\mathcal{Y}_{id}$ on $\mathcal{X}_{ood}$ | $\mathcal{Y}_{ood}$ on $\mathcal{X}_{ood}$ |
|--------------------------|--------------------|---|--|
| Closed-Set               | 96.4 / 96.5        | NA / 86.2                                 | NA / 87.7                                  |
| $\mathcal{L}_{cls}$      | 96.9 / 97.2        | 79.3 / 85.3                               | 71.7 / 87.3                                |
| + $\mathcal{L}_{im-cst}$ | <b>99.2 / 99.2</b> | 84.0 / 91.9                               | 76.3 / 88.3                                |
| + Semantic Enrich        | 98.2 / 99.0        | <b>84.3 / 92.0</b>                        | <b>83.0 / 89.6</b>                         |

(a) Overall accuracy over all YCB objects

|                          | $\mathcal{X}_{id}$ | $\mathcal{Y}_{id}$ on $\mathcal{X}_{ood}$ | $\mathcal{Y}_{ood}$ on $\mathcal{X}_{ood}$ |
|--------------------------|--------------------|---|--|
| Closed-Set               | 91.6 / 91.9        | NA / 57.8                                 | NA / 35.9                                  |
| $\mathcal{L}_{cls}$      | 94.0 / 94.4        | 46.5 / 54.7                               | 23.3 / 32.8                                |
| + $\mathcal{L}_{im-cst}$ | <b>98.1 / 98.0</b> | 55.3 / <b>70.8</b>                        | <b>23.7 / 47.3</b>                         |
| + Semantic Enrich        | 97.2 / 97.4        | <b>55.6 / 70.8</b>                        | 11.7 / <b>50.7</b>                         |

(b) F1-measure over objects that exist in observations.