**TIME**
Towards Intelligence MEchanism

山东大学
SHANDONG UNIVERSITY

深圳河套学院
Shenzhen Loop Area Institute

NEURAL INFORMATION
PROCESSING SYSTEMS

# VT-FSL: Bridging Vision and Text with LLMs for Few-Shot Learning

**Wenhao Li[1,2], Qiangchang Wang[1]\*, Xianjing Meng[3], Zhibin Wu[1], Yilong Yin[1]\***

[1]School of Software, Shandong University    [2]Shenzhen Loop Area Institute
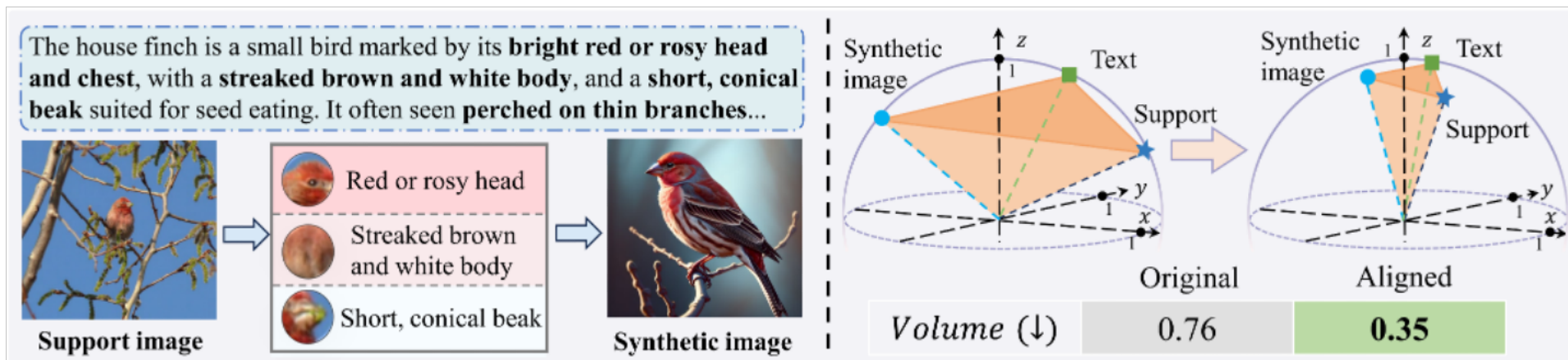[3]School of Computing and Artificial Intelligence, Shandong University of Finance and Economics
{wenhao.li, zhibinwu}@mail.sdu.edu.cn,
rongmengyuan@gmail.com, {qiangchang.wang, ylyin}@sdu.edu.cn

*Corresponding authors
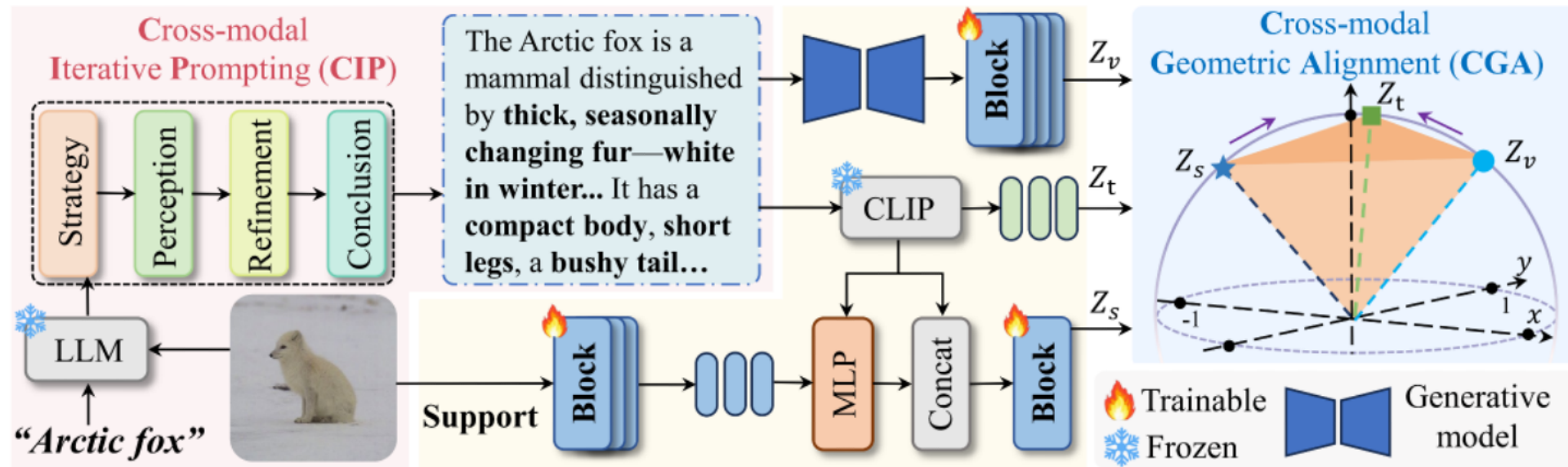
# FSL Task and Motivation Introduction

Imagine the challenge of asking a model to recognize an entirely novel category from only a few sample images. Existing LLM-based methods are prone to **semantic hallucinations**, i.e., generating descriptions that contradict the visual evidence due to the lack of grounding in actual instance,requiring significant effort to correct. So the question arises:

- How can we ensure text descriptions generated by large models truly align with visual features?
- How can we simultaneously capture high-level semantic information and low-level visual diversity with limited samples?
- Compared to traditional CLIP pairwise contrastive learning, how can we integrate multimodal information more comprehensively and structurally?
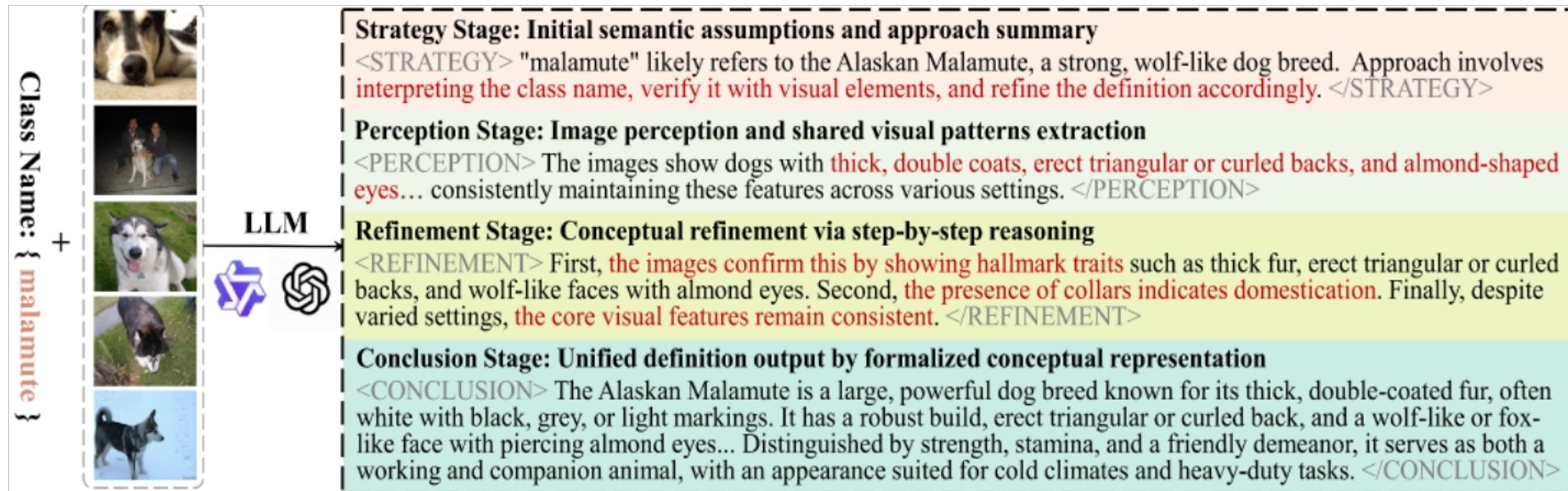
# VT-FSL Framework

- **Cross-modal Iterative Prompting (CIP):** Beyond mere class names, generate precise and visual-grounding descriptions conditioned on LLMs and support images in a single structured iterative reasoning process.
- Semantically Consistent Images Generation: Zero-shot generation based on textual descriptions to enrich limited samples with diverse yet semantically aligned variations.
- **Cross-modal Geometric Alignment (CGA):** Minimize the kernelized volume of a parallelotope to enforce global consistency among textual descriptions, synthetic images, and supports.

# Cross-modal Iterative Prompting (CIP)

> ## Chain-of-Thought (CoT)-based Four Stage Structured Reasoning Paradigm

- **Strategy** Stage: Interpret the and propose initial semantic assumptions.
- **Perception** Stage: Extract visual patterns and shared attributes.
- **Refinement** Stage: Iteratively revise semantics via step-by-step reasoning aligned with visual evidence.
- **Conclusion** Stage: Output a unified, visually grounded class definition with formalized conceptual representation.

# Cross-modal Iterative Prompting (CIP)

The generated description is then fed into a text-to-image generative model to produce synthetic imagesin a zero-shot manner. An LLM-based pairwise comparison strategy is designed to select top-K images per class by ranking them against the descriptions without compromising low-data regime



**Description:** The Alaskan Malamute is a large, powerful dog breed known for its thick, double-coated fur, white with black, grey, or light markings. It has a robust build, erect triangular...

**Synthetic images**     **Support**

Generative model

**Judge Prompt:** Given two images and the description: ... which image best reflects typical visual features of this category?

LLM

+



(a) The effect of the generated number

(b) Visualization of generated images

Figure 7: The effect of the generated number and visualization of generated images.

# Cross-modal Geometric Alignment (CGA)

➢ **Geometric-aware structured and consistent multimodal learning**

Given multiple vectors $\{v_1, \ldots, v_k\}$ to construct matrix $A = [v_1, \ldots, v_k]$, the Gram is:

$$G(v_1, \ldots, v_k) = A^\mathrm{T} A, \quad G_{ij} = \langle \mathbf{v_i}, \mathbf{v_j} \rangle.$$

Gram matrix $G$ reflects the square of the corresponding volume Vol :

$$\mathrm{Vol}(\mathbf{v_1}, \ldots, \mathbf{v_k}) = \sqrt{\det(\mathbf{G})}.$$

To capture nonlinear relations, we compute volume in RKHS via an RBF kernel:

$$\mathrm{Vol}_{\mathcal{H}}(v_1, \ldots, v_k) = \sqrt{\det(\mathbf{K})}, \quad K_{ij} = \kappa(v_i, v_j).$$

Given normalized triplets $\{v_t, v_v, v_s\}$ from textual, vision, and synthetic embeddings, minimizing the kernelized volume through a contrastive objective:

$$\mathcal{L}_{\mathcal{D}2\mathcal{A}} = \frac{1}{B} \sum_{i=1}^{B} \log \frac{\exp\left(-\mathrm{Vol}_{\mathcal{H}}\left(z_t^i, z_s^i, z_v^i\right)/\tau\right)}{\sum_{j=1}^{K} \exp\left(-\mathrm{Vol}_{\mathcal{H}}\left(z_t^i, z_s^i, z_v^j\right)/\tau\right)} \qquad \mathcal{L}_{\mathcal{A}2\mathcal{D}} = -\frac{1}{B} \sum_{i=1}^{B} \log \frac{\exp\left(-\mathrm{Vol}_{\mathcal{H}}\left(z_t^i, z_s^i, z_v^i\right)/\tau\right)}{\sum_{j=1}^{K} \exp\left(-\mathrm{Vol}_{\mathcal{H}}\left(z_t^j, z_s^j, z_v^i\right)/\tau\right)}.$$

# Experimental Results

**On ten standard, cross-domain, and fine-grained datasets, improving accuracy by 4.2% on average. Bold and blue indicates the best and suboptimal results.**

| Model | Venue | Backbone | ≈ # Params | *mini*ImageNet 1-shot | *mini*ImageNet 5-shot | *tiered*ImageNet 1-shot | *tiered*ImageNet 5-shot |
|---|---|---|---|---|---|---|---|
| CPEA [58] | ICCV'23 | ViT-S/16 | 22.0M | $71.97_{\pm0.65}$ | $87.06_{\pm0.38}$ | $76.93_{\pm0.70}$ | $90.12_{\pm0.45}$ |
| FeatWalk [13] | AAAI'24 | ResNet-12 | 12.4M | $70.21_{\pm0.44}$ | $87.38_{\pm0.27}$ | $75.25_{\pm0.48}$ | $89.92_{\pm0.29}$ |
| SemFew [28] | CVPR'24 | Swin-T | 29.0M | $78.94_{\pm0.66}$ | $86.49_{\pm0.50}$ | $82.37_{\pm0.77}$ | $89.89_{\pm0.52}$ |
| UAP [59] | NeurIPS'24 | ResNet-12 | 12.4M | $81.63_{\pm0.28}$ | $79.05_{\pm0.19}$ | $79.68_{\pm0.30}$ | $76.78_{\pm0.21}$ |
| VT-FSL | ours | Visformer-T | 10.0M | $\mathbf{83.66_{\pm0.31}}$ | $\mathbf{88.38_{\pm0.25}}$ | $\mathbf{88.02_{\pm0.34}}$ | $\mathbf{91.71_{\pm0.27}}$ |

| Method | Venue | CUB-200-2011 1-shot | CUB-200-2011 5-shot | Stanford-Dogs 1-shot | Stanford-Dogs 5-shot | Stanford-Cars 1-shot | Stanford-Cars 5-shot |
|---|---|---|---|---|---|---|---|
| SUITED [67] | AAAI'25 | $86.02_{\pm0.47}$ | $94.13_{\pm0.24}$ | $76.55_{\pm0.47}$ | $88.86_{\pm0.27}$ | $89.97_{\pm0.36}$ | $96.53_{\pm0.16}$ |
| VT-FSL | ours | $\mathbf{91.08_{\pm0.28}}$ | $\mathbf{94.63_{\pm0.19}}$ | $\mathbf{86.58_{\pm0.30}}$ | $\mathbf{90.69_{\pm0.25}}$ | $\mathbf{92.95_{\pm0.24}}$ | $\mathbf{96.62_{\pm0.15}}$ |

| Method | Venue | CUB 1-shot | CUB 5-shot | Places 1-shot | Places 5-shot | Plantae 1-shot | Plantae 5-shot |
|---|---|---|---|---|---|---|---|
| StyleAdv [58] | CVPR'23 | $48.49_{\pm0.72}$ | $68.72_{\pm0.67}$ | $58.58_{\pm0.83}$ | $77.73_{\pm0.62}$ | $41.13_{\pm0.67}$ | $61.52_{\pm0.68}$ |
| MEFP [59] | NeurIPS'24 | $51.55_{\pm0.70}$ | $73.61_{\pm0.66}$ | $52.06_{\pm0.69}$ | $73.78_{\pm0.61}$ | $41.55_{\pm0.65}$ | $61.39_{\pm0.67}$ |
| SVasP [70] | AAAI'25 | $49.49_{\pm0.72}$ | $68.95_{\pm0.66}$ | $59.07_{\pm0.81}$ | $77.78_{\pm0.62}$ | $41.22_{\pm0.62}$ | $60.63_{\pm0.64}$ |
| VT-FSL | ours | $\mathbf{66.86_{\pm0.47}}$ | $\mathbf{81.02_{\pm0.36}}$ | $\mathbf{73.68_{\pm0.41}}$ | $\mathbf{81.52_{\pm0.33}}$ | $\mathbf{45.90_{\pm0.40}}$ | $\mathbf{61.54_{\pm0.38}}$ |



(a) T-SNE visualization on novel classes

(b) Visualization of attention maps