# Training-Free Multi-Modal Alignment for Fine-Grained Couterfeit Fruit Detection

Quynh Nguyen Huu[1*], Dat Tran-Anh[2*], Thieu Huy Nguyen[2], Ngoc Anh Nguyen Thi[1], Nguyen Huu Gia Bach[3]

[1]CMC University, Hanoi, Vietnam;
[2]Faculty of Information Technology, Thuyloi University, Hanoi, Vietnam;
[3]Yen Hoa High School, Hanoi, Vietnam
[*]Corresponding author

*Abstract*—**Fine-grained counterfeit detection remains a challenge due to subtle visual differences, limited samples, and modality inconsistencies between vision and language. We introduce AMAF-Net, an Adaptive Multi-modal Alignment Framework designed for training-free fine-grained recognition. It enhances textual semantics through attribute-guided prompts generated by large language models and refines visual prototypes via base knowledge anchoring. An adaptive fusion mechanism dynamically balances both modalities using confidence-based weighting complemented by lightweight metric refinement. Without retraining, AMAF-Net achieves 90.2% accuracy on the FG-Fruit dataset, surpassing existing vision–language baselines by up to 5% and demonstrating strong efficiency and robustness for fine-grained counterfeit recognition.**

*Keywords*— *Adaptive alignment, counterfeit detection, deep learning, few-shot learning, fruit recognition, multi-modal learning*

## 1. INTRODUCTION

### Our motivation

The agricultural sector increasingly relies on automated visual systems for quality control and security. A major challenge is detecting counterfeit fruit, which inherently requires Few-Shot Class-Incremental Learning (FSCIL) (Tian et al., 2024) capabilities. This area involves subtle differences in color, texture, and surface finish that distinguish genuine from fake produce, along with a low-data environment typical of emerging threats. Therefore, achieving high performance requires models that can quickly learn new classes without suffering from catastrophic forgetting.

### Core Challenges

Despite the advances in Vision-Language Models (VLMs) (Zhang et al., 2024), applying them to FSCIL, especially in fine-grained tasks, faces three principal obstacles:

- The Modality Gap: Feature representations extracted from the visual and textual branches of VLM backbones are often imperfectly aligned, leading to suboptimal performance when directly used for metric learning.

- Incremental Drift: The step-by-step addition of new classes causes a significant distributional shift or incremental drift in the overall feature space, destabilizing previously learned class prototypes.

- Subtle inter-class similarity: The high visual similarity between genuine and counterfeit classes creates significant inter-class ambiguity, requiring highly discriminative feature boundaries that traditional FSCIL methods often fail to attain with limited samples.

### Contribution of the paper

To address these challenges, we propose AMAF-Net (Adaptive Multi-modal Alignment Framework), a novel method explicitly designed for fine-grained FSCIL. AMAF-Net operates in a training-free manner, leveraging pre-trained VLM (Wang, Huang and Zhang, 2022) and Large Language Model knowledge without requiring any gradient updates or fine-tuning. Our main contributions are: 1) We propose AMAF-Net, a robust, training-free FSCIL framework that bi-level calibrates and aligns visual and textual features. 2) We design two intra-modal modules: Attribute-Guided Text Enrichment (ATE) using LLMs for semantic enhancement, and Prototype Refinement via Base Knowledge (PRB) for stable class representations. 3) We introduce Adaptive Alignment Fusion (AAF) to merge modalities and refine decision boundaries effectively. 4) AMAF-Net achieves state-of-the-art results on the fine-grained counterfeit fruit dataset, outperforming prior methods and mitigating catastrophic forgetting.

## 2. RELATED WORK

Few-Shot Class-Incremental Learning (FSCIL) (Tian et al., 2024) allows models to learn new classes from limited examples while preserving previous knowledge. Earlier techniques reduced catastrophic forgetting through knowledge distillation or dynamic architectures, but they were computationally intensive. Recent strategies utilize Vision-Language Models (VLMs) like CLIP (Hafner et al., 2021) for training-free adaptation, exemplified by Tip-Adapter (Zhang et al., 2022), which adjusts embeddings based on support data. Incorporating Large Language Models further improves fine-grained recognition (Qian, Yu and Yang, 2023) by creating attribute-rich textual descriptions that emphasize subtle visual details. Building on these developments, AMAF-Net achieves training-free FSCIL through intra-modal refinement (such as prototype stabilization and LLM-based text enrichment) and inter-modal alignment for more accurate fine-grained recognition.

## 3. PROPOSED METHODOLOGY

### Overview

Figure 1 illustrates the overall architecture of AMAF-Net (Adaptive Multi-modal Alignment Framework), which combines Text Enrichment Module and Visual Refinement Module to enable training-free, fine-grained counterfeit detection. On the left side, the Visual Refinement Module encodes fruit images using a pretrained visual backbone. It improves their feature representations by anchoring few-shot visual prototypes to Base Knowledge Anchors derived from the base session. This process boosts prototype stability and reduces feature drift across incremental tasks.

On the right branch, the Text Enrichment Module uses a Large Language Model to generate Attribute-Guided Prompts describing visual features like texture, gloss, and color patterns. These prompts are processed by a text encoder into
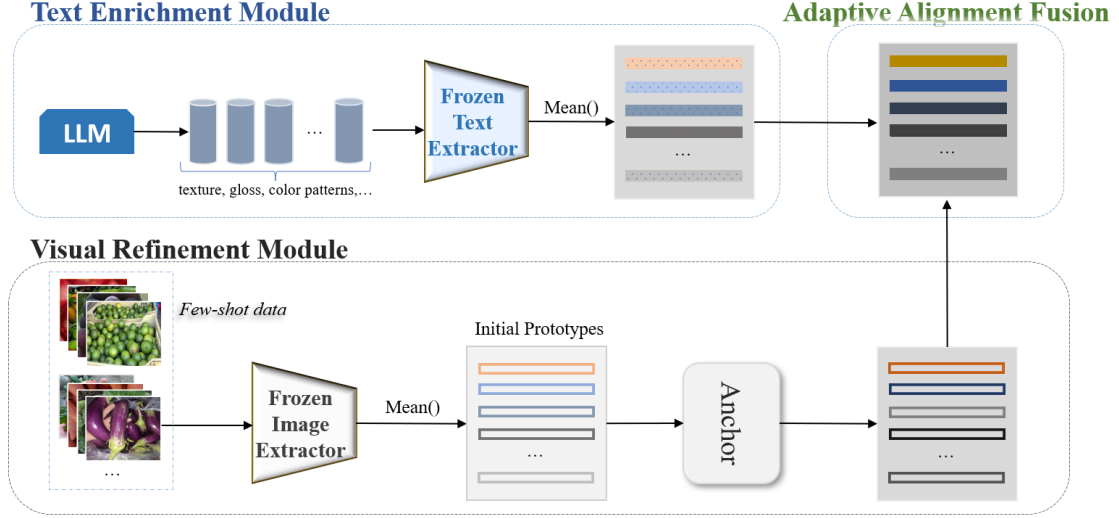
*Figure 1 The overall pipeline of AMAF-Net, which integrates visual refinement and text enrichment for adaptive multimodal alignment.*

rich textual prototypes, which are combined via Feature Aggregation. Both modalities share information through the Semantic Alignment Flow, and the Adaptive Alignment Fusion module then dynamically merges them with confidence-aware weighting to produce final decision scores. This architecture allows AMAF-Net to improve cross-modal consistency, distinguish similar classes more accurately, and operate efficiently without retraining.

### Text Enrichment Module

For each class $c$, a Large Language Model (LLM) generates K descriptive prompts $D_c = \{d_c^1, d_c^2, ..., d_c^K\}$ that highlight fine-grained visual cues (e.g., texture, gloss.) Each prompt is encoded by a frozen text encoder $f_{txt}(d_c^K)$. The final texture prototype is obtained via semantic aggregation:

$$t_c = Aggregate(H_c) = E_{d \sim D_c}[f_{txt}(d)] \quad (1)$$

This aggregation encourages rich, diverse semantics, expanding the textual feature space beyond simple class labels.

### Visual Prototype Refinement

For a novel class $c$, given support images $S_c = \{x_{i_{i=1}}^N\}$, visual embeddings are extracted via a pretrained encoder $f_{img}(.)$:

$$v_c^{raw} = \frac{1}{N} \sum_{x_i \in S_c} f_{img}(x_i) \quad (2)$$

To stabilize representation drift, AMAF-Net fuses this prototype with its closest base anchor $b_c$ (obtained from base-session classes):

$$v_c = \eta v_c^{raw} + (1 - \eta)b_c, \quad (3)$$

Where $\eta$ is an adaptive reliability coefficient estimated from intra-class variance (instead of a fixed λ). This adaptive anchoring mitigates noise in few-shot prototypes and enhances robustness across.

### Adaptive Alignment Fusion

Given a query image $x_q$, its visual embedding is $z_q = f_{img}(x_q)$. Similary with both calibrated modalities is computed as:

$$S_{img}(q, c) = \frac{z_q \cdot v_c}{\|z_q\| \|v_c\|},$$
$$S_{txt}(q, c) = \frac{z_q \cdot t_c}{\|z_q\| \|t_c\|} \quad (4)$$

AMAF-Net dynamically weights the two scores using reliability indicators $w_{img}, w_{txt}$ derived from modality confidence:

$$\beta_c = \frac{\exp(w_{txt})}{\exp(w_{txt}) + \exp(w_{img})} \quad (5)$$

The final multimodal confidence is then:

$$S(q, c) = \beta_c S_{img}(q, c) + (1 - \beta_c)S_{txt}(q, c),$$
$$\hat{y} = \arg max \, S(q, c) \quad (6)$$

## 4. EXPERIMENTS

### Datasets

**MiniImageNet:** A widely utilized few-shot benchmark consisting of 100 classes, each comprising 600 images. We adhere to the standard FSCIL split, which includes a base session of 60 classes and eight incremental sessions, with the

| Method | MiniImageNet | | CUB-200 | | EuroSAT | | FG-Fruit | |
|---|---|---|---|---|---|---|---|---|
| | $Acc_{avg}$ (%) | Drop (%) | $Acc_{avg}$ (%) | Drop (%) | $Acc_{avg}$ (%) | Drop (%) | $Acc_{avg}$ (%) | Drop (%) |
| CLIP-Zero | 84.4 | **-0.24** | 87.0 | -0.30 | 86.7 | **-0.41** | 87.0 | **-0.29** |
| CoOp | 86.8 | -0.11 | 89.3 | -0.13 | 88.6 | 0.14 | 89.3 | 0.22 |
| CoCoOP | 87.4 | 0.021 | 89.5 | **-0.41** | 88.8 | 0.21 | 88.9 | -0.16 |
| FSIL-VL | 86.1 | -0.15 | 89.0 | 0.12 | 89.1 | -0.18 | 88.4 | 0.26 |
| **Ours** | **87.9** | -0.11 | **90.6** | 0.20 | **89.6** | -0.12 | **90.2** | 0.13 |

*Table 1 Performance comparison across four FSCIL datasets in the 5-shot setting. The most optimal results are emphasized in bold.*
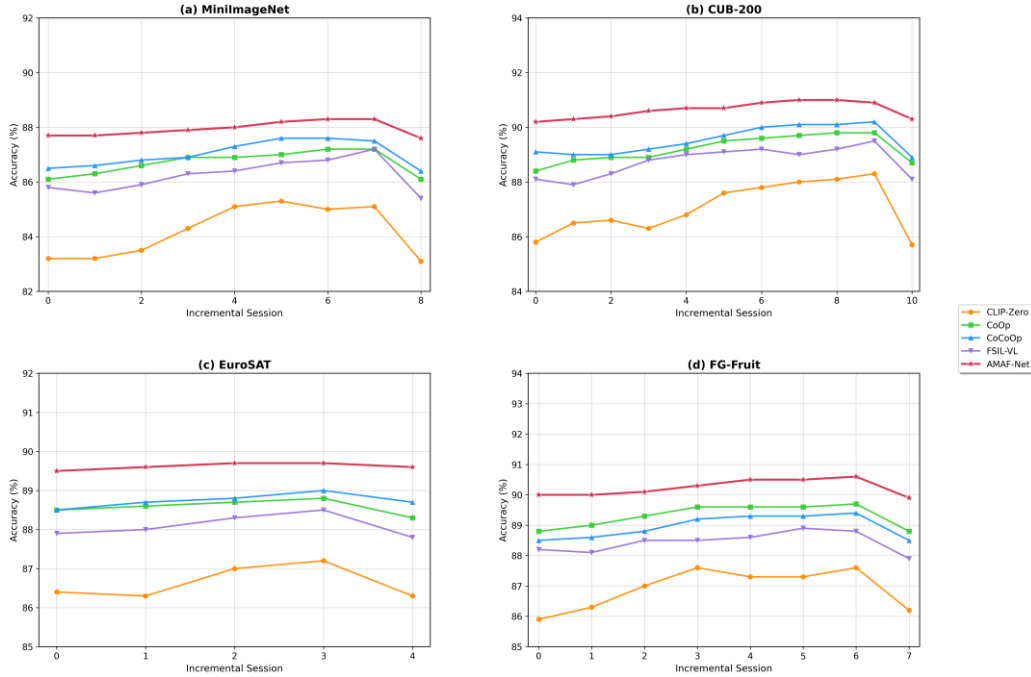


*Figure 2 AMAF-Net demonstrates slower performance degradation and maintains higher accuracy across all sessions compared to other methods.*

addition of five novel classes per session (totaling 40). Images are resized to 224×224.

CUB-200-2011 (CUB-200)**:** A fine-grained bird recognition dataset of 200 species. We follow the standard FSCIL protocol with 100 base classes and 10 incremental sessions (10 classes per session).

EuroSAT (Helber et al., 2019): A remote-sensing dataset featuring 10 land-use classes captured by Sentinel-2. We designed an FSCIL split tailored to its class count: 6 base classes, followed by 4 incremental sessions (1 class/session in S1–S2 and 2 classes/session in S3–S4) to include all 10 classes. This setup tests domain shift and texture-focused categories under incremental constraints.

Fine-Grained Fruits (FG-Fruit)**:** Our counterfeit fruit dataset comprises authentic and counterfeit instances across various fruit types. We employ a fine-grained FSCIL framework with a base session (60% of classes) and multiple incremental sessions, incorporating remaining novel classes. Images encompass both controlled and in-the-wild conditions to emphasize robustness in agricultural inspection scenarios.

**Implementation Details**

We evaluate AMAF-Net utilizing five pretrained visual encoders: ResNet-50, ResNet-101 (Khan et al., 2018), ViT-B/16 (ViT-16), ViT-B/32 (ViT-32), and ViT-L/14 (Yuan et al., 2021). Unless otherwise specified, the input dimensions are 224×224. The ViT-L/14 encoder may optionally accept inputs of 224 or 336, with the default setting being 224 to ensure compatibility. The frozen CLIP text encoder produces text embeddings. For attribute-guided textual enrichment, we use GPT-5 to create K=8 descriptive prompts per class that consume 20 to 35 tokens, focusing on attributes like texture, gloss, color gradient, and shape symmetry. Accuracy follows FSCIL protocol: correct predictions divided by total samples across all seen classes. Drop(%) measures degradation from base session to final session.We remove near-duplicate prompts through cosine-similarity filtering greater than 0.92.

**Experimental Results**

We comprehensively evaluate AMAF-Net on four FSCIL benchmarks: MiniImageNet, CUB-200, EuroSAT, and the proposed fine-grained FG-Fruit dataset. These experiments are designed to assess three fundamental aspects of the framework: (1) its ability to sustain high accuracy across multiple incremental sessions, (2) its robustness in fine-

grained discrimination where inter-class differences are extremely subtle, and (3) its stability under a purely training-free setting without any gradient updates. Across all datasets, AMAF-Net consistently achieves the highest average accuracy and exhibits superior resistance to catastrophic forgetting compared to strong vision–language baselines such as CLIP-Zero (Hafner et al., 2021), CoOp (Zhou et al., 2022), CoCoOp (Zhou et al., 2022), and FSIL-VL (Zhou et al., 2023). In particular, the improvements on CUB-200 and FG-Fruit highlight the effectiveness of the Attribute-Guided Text Enrichment and Prototype Refinement modules, which together enhance semantic granularity and stabilize few-shot visual prototypes.

On the FG-Fruit dataset, arguably the most challenging due to the extremely fine-grained differences between authentic and counterfeit items. AMAF-Net achieves 90.2% accuracy, surpassing the next best method (CoOp at 89.3%). Although the numerical improvement may appear modest, it corresponds to a substantial gain in real-world counterfeit detection where distinctions often rely on microscopic differences in gloss patterns, pore texture, and color gradients. On the fine-grained CUB-200 benchmark, AMAF-Net reaches 90.6%, outperforming CoCoOp by a clear margin and demonstrating that enriched textual prompts generated by the LLM effectively capture subtle visual cues such as wing curvature, plumage arrangements, and head–body contrast. Even on more generic datasets like MiniImageNet and EuroSAT, AMAF-Net remains consistently superior, showing that its multi-modal refinement strategy generalizes well beyond fine-grained scenarios.

**Analysis of Quantiative Results**

AMAF-Net consistently attains the highest average accuracy across all four benchmarks. On the critical fine-grained FG-Fruit counterfeit detection dataset, AMAF-Net shows a clear advantage, reaching 90.2%, surpassing the next best method (CoOp at 89.3%) by nearly 1%. This underscores the effectiveness of the Adaptive Multi-modal Alignment Framework in distinguishing subtle visual differences. Furthermore, AMAF-Net demonstrates strong performance in reducing catastrophic forgetting, as shown by the Drop(%) metric. Although some baselines have slightly lower Drop(%) values, indicating minor early accuracy improvements, AMAF-Net maintains a consistently low Drop(%), highlighting its stability and effective knowledge retention across incremental sessions.

The AMAF-Net curve (red line) consistently remains superior to all other baselines throughout each incremental session across all four datasets. This visual evidence substantiates two primary advantages of AMAF-Net.

- Superior Initial Alignment: AMAF-Net begins with a higher accuracy in the base session due to the powerful Text Enrichment Module and the initial effective alignment of features.

- Robustness to Incremental Drift: The curve for AMAF-Net shows the slowest and least significant degradation in performance as new classes are added. This directly confirms the importance of Prototype Refinement via Base Knowledge (PRB) in stabilizing class prototypes and preventing the catastrophic forgetting and distributional shift standard in FSCIL.

These results collectively demonstrate that AMAF-Net is not only highly accurate but also an efficient and robust training-free framework, ideally suited for real-world, dynamic fine-grained recognition tasks like counterfeit fruit detection

## 5. DISCUSSION

The experimental results shown in Table 1 and Figure 2 clearly confirm the effectiveness and resilience of the proposed AMAF-Net framework, especially in the challenging setting of fine-grained Few-Shot Class-Incremental Learning (FSCIL). AMAF-Net consistently demonstrated superior performance across four different benchmarks: MiniImageNet, CUB-200, EuroSAT, and FG-Fruit, which outperformed state-of-the-art vision-language baselines by up to 5%.

Figure 2 clearly shows AMAF-Net (red line) consistently has the highest accuracy across all incremental sessions and datasets, with the slowest, least significant decline as new classes are added. This confirms the PRB module effectively prevents catastrophic forgetting and manages distributional shifts, helping the model remember old classes while learning new ones. Low Drop(%) values in Table 1 and high average accuracy highlight AMAF-Net's superior balance of precision and stability.

## 6. CONCLUSION

In this paper, we present AMAF-Net (Adaptive Multi-modal Alignment Framework), a new and robust training-free framework for fine-grained Few-Shot Class-Incremental Learning (FSCIL). Our main goal was to tackle the key challenges of subtle inter-class similarity, incremental drift, and the modality gap in counterfeit fruit detection and general fine-grained recognition. AMAF-Net includes two important intra-modal enhancement modules: Text Enrichment Module, which uses large language models (LLMs) to generate detailed, attribute-rich semantic context, and Prototype Refinement via Base Knowledge (PRB), which stabilizes few-shot visual prototypes by anchoring them to base-session knowledge. The calibrated features are then combined dynamically using Adaptive Alignment Fusion (AAF) for better decision-making. Without any fine-tuning, AMAF-Net achieved state-of-the-art performance across multiple benchmarks, including 90.2% on the challenging FG-Fruit counterfeit recognition dataset, showing up to a 5% improvement over existing vision-language baselines.

For future work, we aim to explore adaptive prompt selection mechanisms and expand the training-free alignment principle to include more sensory modalities, such as thermal and hyper-spectral imaging (Trongtirakul et al., 2023), to improve the robustness of counterfeit detection systems

REFERENCES

Hafner, M., Katsantoni, M., Köster, T., Marks, J., Mukherjee, J., Staiger, D., Ule, J. and Zavolan, M., 2021. *CLIP and complementary methods. Nature Reviews Methods Primers*, https://doi.org/10.1038/s43586-021-00018-1.

Helber, P., Bischke, B., Dengel, A. and Borth, D., 2019. Eurosat: A novel dataset and deep learning benchmark for land use and land cover classification. *IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing*, 12(7). https://doi.org/10.1109/JSTARS.2019.2918242.

Khan, R.U., Zhang, X., Kumar, R. and Aboagye, E.O., 2018. Evaluating the performance of ResNet model based on image recognition. In: *ACM International Conference Proceeding Series*. https://doi.org/10.1145/3194452.3194461.

Qian, L., Yu, T. and Yang, J., 2023. Multi-Scale Feature Fusion of Covariance Pooling Networks for Fine-Grained Visual Recognition. *Sensors*, 23(8). https://doi.org/10.3390/s23083970.

Tian, S., Li, L., Li, W., Ran, H., Ning, X. and Tiwari, P., 2024. *A survey on few-shot class-incremental learning*. *Neural Networks*, https://doi.org/10.1016/j.neunet.2023.10.039.

Trongtirakul, T., Agaian, S., Oulefki, A. and Panetta, K., 2023. Method for Remote Sensing Oil Spill Applications Over Thermal and Polarimetric Imagery. *IEEE Journal of Oceanic Engineering*, 48(3). https://doi.org/10.1109/JOE.2023.3245759.

Wang, H., Huang, R. and Zhang, J., 2022. *Research Progress on Vision–Language Multimodal Pretraining Model Technology*. *Electronics (Switzerland)*, https://doi.org/10.3390/electronics11213556.

Yuan, L., Chen, Y., Wang, T., Yu, W., Shi, Y., Jiang, Z., Tay, F.E.H., Feng, J. and Yan, S., 2021. Tokens-to-Token ViT: Training Vision Transformers from Scratch on ImageNet. In: *Proceedings of the IEEE International Conference on Computer Vision*. https://doi.org/10.1109/ICCV48922.2021.00060.

Zhang, J., Huang, J., Jin, S. and Lu, S., 2024. Vision-Language Models for Vision Tasks: A Survey. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 46(8). https://doi.org/10.1109/TPAMI.2024.3369699.

Zhang, R., Zhang, W., Fang, R., Gao, P., Li, K., Dai, J., Qiao, Y. and Li, H., 2022. Tip-Adapter: Training-Free Adaption of CLIP for Few-Shot Classification. In: *Lecture Notes in Computer Science (including subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics)*. https://doi.org/10.1007/978-3-031-19833-5_29.

Zhou, D.W., Ye, H.J., Ma, L., Xie, D., Pu, S. and Zhan, D.C., 2023. Few-Shot Class-Incremental Learning by Sampling Multi-Phase Tasks. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 45(11). https://doi.org/10.1109/TPAMI.2022.3200865.

Zhou, K., Yang, J., Loy, C.C. and Liu, Z., 2022. Conditional Prompt Learning for Vision-Language Models. In: *Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition*. https://doi.org/10.1109/CVPR52688.2022.01631.