PART-LEVEL SEMANTIC-GUIDED CONTRASTIVE LEARNING FOR FINE-GRAINED VISUAL CLASSIFI-CATION

Anonymous authorsPaper under double-blind review

ABSTRACT

Fine-Grained Visual Classification (FGVC) aims to distinguish visually similar subcategories within a broad category, and poses significant challenges due to subtle inter-class differences, large intra-class variations, and data scarcity. Existing methods often struggle to effectively capture both part-level detail and spatial relational features, particularly across rigid and non-rigid object categories. To address these issues, we propose Part-level Semantic-guided Contrastive Learning (PSCL), a novel framework that integrates three key components. (1) The Part Localization Module (PLM) leverages clearCLIP to enable text-controllable region selection, achieving decoupled and semantically guided spatial feature extraction. (2) The Multi-scale Multi-part Branch Progressive Reasoning (MMBPR) module captures discriminative features across multiple parts and scales, while reducing inter-branch redundancy. (3) The Visual-Language Contrastive Learning based on Multi-grained Text Features (VLCL-MG) module introduces intermediategranularity category concepts to improve feature alignment and inter-class separability. Extensive experiments on five publicly available FGVC datasets demonstrate the superior performance and generalization ability of PSCL, validating the effectiveness of its modular design and the synergy between vision and language. Code is available at: https://anonymous.4open.science/r/PSCL-3E1F.

1 Introduction

Fine-Grained Visual Classification (FGVC) aims to accurately distinguish between subcategories that belong to the same high-level category yet exhibit subtle visual differences. Typical applications include the classification of bird species (Wah et al., 2011; Van Horn et al., 2015), car brands (Krause et al., 2013), and aircraft (Maji et al., 2013) models. As FGVC focuses on fine-level distinctions within specific domains, it has demonstrated unique practical value—distinct from general visual classification tasks—in areas such as intelligent transportation, medical image analysis, and ecological environment monitoring. However, FGVC remains a challenging task due to factors such as low inter-class variance, high intra-class variance, a large number of categories, and data scarcity.

We observe that existing models exhibit notable feature preferences when processing rigid and non-rigid objects. We argue that FGVC tasks require the modeling of two key types of features: (1) part-level fine-grained features that capture detailed local differences and (2) spatial relational features that describe inter-class differences in spatial structure. For rigid objects, inter-class differentiation is often affected by external factors such as viewpoint variation and occlusion. In contrast, non-rigid objects tend to exhibit more significant posture variations, leading to greater uncertainty in their spatial structural features. Different model architectures vary considerably in their capacity to capture these two types of features.

Some existing works have consciously incorporated mechanisms for modeling spatial structural features. For example, CAP (Behera et al., 2021) captures spatial relations through region consistency integration, while SFETrans (Yu et al., 2025) extracts spatial features via phase spectrum analysis. These methods have demonstrated effectiveness in improving classification performance for rigid objects. However, the core objective of FGVC lies in accurately modeling subtle inter-class differences. Since spatial relational features often rely on matching shared regions across categories,

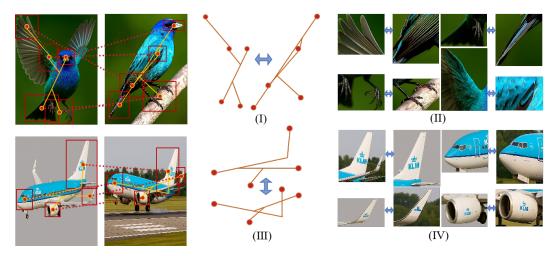


Figure 1: Two types of critical features in rigid and non-rigid objects. (I) Spatial deformation in non-rigid objects (e.g., birds) due to articulated motion; (II) Diverse part-level details in non-rigid objects; (III) Stable spatial structure in rigid objects (e.g., airplanes); (IV) Consistent part-level appearance in rigid objects.

they may conflict with the precise representation of part-level details—particularly for non-rigid objects—potentially weakening the model's ability to focus on critical parts. Furthermore, current models generally adopt a unified strategy for designing part-based branches across all categories, overlooking the homogeneity of part-level details among similar categories. This can lead to misclassifications and redundant representations across branches.

To address these issues, we propose a novel framework called Part-level Semantic-guided Contrastive Learning (PSCL). This model introduces a Part Localization Module (PLM), which leverages clearCLIP (Lan et al., 2024) as an auxiliary component to enable text-guided region selection, thereby achieving effective decoupling of feature region selection and feature representation. Additionally, we design a Multi-scale Multi-part Branch Progressive Reasoning (MMBPR) module, where part-based branches represent fine-grained features of individual parts, while a global branch adaptively integrates features based on spatial relations. Through progressive reasoning, MMBPR enables each branch to refine its feature representations across multiple scales.

During the multi-scale feature fusion stage, we further design the Reverse-key Scale-aware Attention Fusion Module (ReSAF) to suppress the influence of high-level features on semantically similar regions at lower levels, thereby encouraging the model to extract information from less similar areas. This effectively mitigates feature redundancy among branches.

Finally, in the classification phase, we introduce a novel Visual-Language Contrastive Learning based on Multi-grained Text Features (VLCL-MG) module. By incorporating intermediate-level category concepts, this module leverages prior knowledge to aggregate fine-grained categories into semantically coherent mid-level groups, promoting more meaningful clustering of similar subcategories in the feature space.

Our main contributions can be summarized as follows:

- We propose a Part Localization Module (PLM) that enables text-controllable spatial feature extraction via clearCLIP;
- We design a Multi-scale Multi-part Branch Progressive Reasoning (MMBPR) module to reduce feature redundancy and enhance part-level and global representations;
- We introduce a Visual-Language Contrastive Learning module based on Multi-grained Text Features (VLCL-MG) to improve the semantic alignment of visually similar subcategories;
- Extensive experiments on five publicly available FGVC datasets validate the effectiveness and generalization ability of our proposed PSCL framework.

2 RELATED WORK

110111112

113

114

115

116

117

118

108

109

2.1 FINE-GRAINED VISUAL CLASSIFICATION

region segmentation without strong part annotations.

2.2 VISION-LANGUAGE LEARNING

Fine-grained visual classification (FGVC) methods primarily focus on capturing subtle inter-class

differences through refined feature representation and part localization. Early feature representation

approaches relied on high-level features (Lin et al., 2015; Zheng et al., 2019; Sun et al., 2020), later

incorporating multi-scale fusion techniques such as AP-CNN (Ding et al., 2021) and PMG (Du et al.,

2020), as well as attention-based mechanisms like MA-CNN (Zheng et al., 2017), OSME (Sun et al.,

2018), and Transformer-based methods such as TransFG (He et al., 2022) and CAMF (Miao et al.,

2021). MDCM (Zhang et al., 2025) introduces a multi-scale ViT framework that improves fine-

grained bird recognition by activating, selecting, and aggregating discriminative cues across scales.

2014) and Pose Normalized CNN (Branson et al., 2014) relied on strong supervision, recent methods

have shifted to weak supervision for better scalability. Notable examples include MGE-CNN (Zhang

et al., 2019), P2P-Net (Yang et al., 2022), CAP (Behera et al., 2021), TBMSL-Net (Zhang et al.,

2021), and PART (Zhao et al., 2021), which explore part-level semantics via multi-scale learning,

context modeling, or Transformer-based architectures. CSQA-Net (Xu et al., 2025) introduces a

Part Navigator module to assign saliency scores to different image regions, enabling discriminative

Vision-language models (VLMs), particularly CLIP (Radford et al., 2021), have demonstrated strong

potential in open-vocabulary tasks by learning joint representations from large-scale image-text pairs. While early FGVC-related works using CLIP (Li et al., 2023; Wang et al., 2023b) empha-

sized alignment between descriptive text and novel categories, MP-FGVC (Jiang et al., 2024) intro-

duced CLIP to closed-set FGVC by leveraging multimodal prompts to enhance category discrimi-

nation. For region-level tasks, CLIP's utility has been extended to open-vocabulary segmentation.

MaskCLIP (Zhou et al., 2022) revealed that dense patch-level features from CLIP's attention layers

could be aligned with textual representations. Building on this, ClearCLIP (Lan et al., 2024) demon-

strates that by removing residual connections in CLIP, enabling self-attention, and eliminating the

feed-forward network, open-vocabulary semantic segmentation can be achieved directly without ad-

ditional training. We empirically demonstrate that ClearCLIP is also effective for part-level semantic

119 120 121

122

123

124

Part localization methods identify discriminative regions through cropping and scaling strategies. This line of work aims to locate category-relevant regions within the input image by analyzing attention maps generated by the backbone network. The identified regions are then cropped and reprocessed to retain high-resolution, fine-grained details that are critical for classification. This strategy explicitly extracts spatial structural features by emphasizing salient parts, often leading to superior classification performance. While early approaches like Part-based R-CNN (Zhang et al.,

125 126 127

128 129 130

131132

133 134

135 136

144145146

147148149

151

152

153

154

155

156

157

158

159

160

161

149 3 APPROACH

concepts.

The proposed PSCL architecture is illustrated in fig. 2. In the visual pathway, the input image is first processed separately by the backbone and ClearCLIP. ClearCLIP generates part masks by computing matching scores and applying channel selection, while the backbone produces multi-scale features. The two outputs are combined using the Hadamard product to obtain multi-scale part-level features, forming the Part Localization Module. The designed Multi-scale Multi-part Branch Progressive Reasoning module processes the resulting visual features, progressively enhancing the model's confidence in its predictions from low-level to high-level features. This confidence enhancement is achieved through a combination of hyperparameters for contrastive loss weights across different scales and noise parameters. In the text pathway, contrastive loss leverages intermediate-grained

textual priors as input, generating multi-grained textual features for different categories. These features are then rearranged and restructured to produce multi-grained textual representations for each fine-grained label.

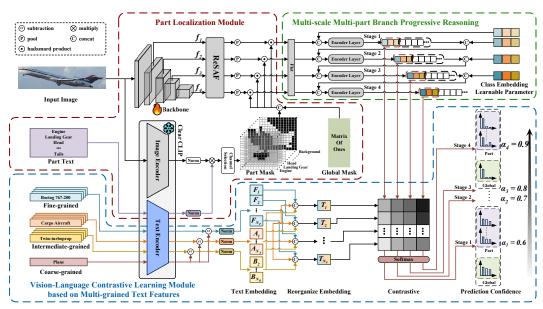


Figure 2: Detailed illustration of Part-level Semantic-guided Contrastive Learning model (PSCL).

3.1 PART LOCALIZATION MODULE

The proposed Part Localization Module (PLM) is designed to address the conflicting requirements of modeling fine-grained part-level features and spatial relational features in FGVC. This conflict is particularly pronounced for non-rigid objects, where posture variation undermines the stability of spatial structures and affects precise part representation. To resolve this, PLM processes the input image $\mathbf{x} \in \mathbb{R}^{C \times H \times W}$ through two separate branches: one for capturing difference-aware features and the other for localizing discriminative parts, enabling more effective and targeted feature learning across object types.

The branch responsible for representing differences processes the input \mathbf{x} to produce multiscale features, with features denoted as $f_s \in \mathbb{R}^{C_s \times H_s \times W_s}$ across multiple stages. When low-level features are less relevant for classification, only higher stages may be selected, such that

$$s \in \{s_{\min}, \dots, 4\}, \quad s_{\min} \ge 1,\tag{1}$$

where s_{\min} denotes the earliest stage used, which can be adjusted based on task requirements.

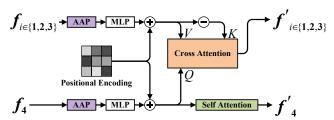


Figure 3: Illustration of ReSAF.

The resulting f_s is then passed into the Reverse-key Scale-aware Attention Fusion Module (ReSAF) to suppress redundant channel representations across scales, as illustrated in fig. 3. In the figure, AAP denotes Adaptive Average Pooling, and the positional encoding is implemented as a learnable parameter. By flipping key vector directions, ReSAF inverts

similarity scores, guiding high-level queries to attend away from similar low-level regions. This contrastive attention promotes the extraction of complementary.

The branch responsible for identifying the parts of interest is built upon the ClearCLIP backbone. The input image \mathbf{x} is encoded by the image encoder f_{img} , producing patch-level image features $\mathbf{F}_{\text{img}} = f_{\text{img}}(\mathbf{x}) \in \mathbb{R}^{H \times W \times d}$. Similarly, the textual prompts corresponding to N parts, denoted as $\mathbf{T} = \{T_1, T_2, \ldots, T_N\}$, are processed by the text encoder f_{text} , yielding part-specific text feature representations $\mathbf{F}_{\text{text}} = \{f_{\text{text}}(T_1), f_{\text{text}}(T_2), \ldots, f_{\text{text}}(T_N)\} \in \mathbb{R}^{N \times d}$. To align image patches with text descriptions, the similarity tensor \mathbf{S} is computed via matrix multiplication.

$$S = F_{\text{img}} F_{\text{text}}^{\top}, S \in R^{H \times W \times N}, \tag{2}$$

To generate the final part mask \mathbf{M} , the indices of the maximum similarity scores across the N channels are first determined as

$$\max_{j \in \{1,\dots,N\}} \mathbf{S}[j], \tag{3}$$

Using these indices, a one-hot-like tensor is constructed:

$$\mathbf{S}[j] = \begin{cases} 1, & j = \text{max_indices}, \\ 0, & \text{otherwise}, \end{cases}$$
 (4)

The one-hot-like tensor undergoes morphological refinement by first applying dilation to expand regions, followed by erosion to refine connectivity and remove noise:

$$\mathbf{M} = \min_{\mathbf{k} \in \mathcal{K}} \max_{\mathbf{k} \in \mathcal{K}} \mathbf{S}(\mathbf{k}),\tag{5}$$

where K represents the kernel window.

The multi-scale multi-part features $G_{s,n'}$ can be expressed as:

$$G_{s,n'} = \operatorname{concat}(f_s' \odot \mathbf{M}_{s,n}, f_s' \odot \mathbf{1}), n \in \{1, 2, \dots, N\}, n' \in \{1, 2, \dots, N+1\},$$
 (6)

where \odot denotes the Hadamard product, concat(\cdot) denotes the concatenation operation, $\mathbf{1}$ is a matrix of ones, representing the global mask, $f_s' \odot \mathbf{1}$, captures the global features. The global features are subsequently processed by the global branch, which adaptively aggregates part-level information according to spatial relationships, aiming to model spatial relational features.

3.2 Branch Progressive Reasoning

Our proposed Multi-scale Multi-part Branch Progressive Reasoning (MMBPR) module extends the multi-scale reasoning framework introduced by PMG (Du et al., 2020) and PART (Zhao et al., 2021). It progressively enhances the constraint on loss from low-level feature branches to high-level feature branches. Unlike previous approaches, our method incorporates PLM to reduce redundant representations across branches. In addition to progressive reasoning along the multi-scale hierarchy, the framework also integrates part-level branches for capturing fine-grained local details and a global branch for modeling spatial relational representations through part feature aggregation.

Following the ViT architecture, we adopt Class Embedding Learnable Parameters to extract category-specific visual representations. TWe utilize three class tokens to comprehensively represent intermediate categories.

The overall input to MMBPR is $G_{s,n'}$ obtained from eq. (6), and the progressive reasoning process begins from the lowest-level features $G_{s_{\min},n'}$, as formulated below:

$$G_{s_{\min},n'} \xrightarrow{\text{flatten}} \{V_{s_{\min},n',m} \mid m \in \{1,2,\dots,P^2\}\},\tag{7}$$

The flattened tokens are concatenated with class tokens:

$$Z_{s_{\min}} = \operatorname{concat}([C_{\operatorname{cls}, s_{\min}, 1}; C_{\operatorname{cls}, s_{\min}, 2}; C_{\operatorname{cls}, s_{\min}}, 3], [V_{s_{\min}, 1, 1}; \dots; V_{s_{\min}, N+1, P^2}]), \tag{8}$$

where $C_{\text{cls},j}$ $(j \in \{1,2,3\})$ are class tokens.

The sequence Z is then processed through an encoder:

$$\begin{split} Z_{s_{\min}}' &= \text{LN}(Z_{s_{\min}} + \text{MHSA}(Z_{s_{\min}})), \\ Z_{s_{\min}}'' &= \text{LN}(Z_{s_{\min}}' + \text{MLP}(Z_{s_{\min}}')), \end{split} \tag{9}$$

where $LN(\cdot)$ is Layer Normalization, $MHSA(\cdot)$ is Multi-Head Self-Attention, $MLP(\cdot)$ is a feedforward neural network. The resulting output is:

$$Z''_{s_{\min}} = [C''_{\text{cls},s_{\min},1}; C''_{\text{cls},s_{\min},2}; C''_{\text{cls},s_{\min},3}; V''_{s_{\min},1,1}; \dots; V''_{s_{\min},N+1,P^2}],$$
(10)

The output $Z''_{s_{\min}}$ is then divided into two parts: Class tokens, representing the visual features of the categories at the stage $\mathbf{I}_{s_{\min}} = [C''_{\text{cls},s_{\min},1}; C''_{\text{cls},s_{\min},2}; C''_{\text{cls},s_{\min},3}]$ are passed to VLCL-MG for contrastive learning. Feature tokens $[V''_{s_{\min},1,1};\ldots;V''_{s_{\min},N+1,P^2}]$ are forwarded to the next stage, where they are concatenated with the flattened high-level features.

By ensuring non-interference between the lower and higher branches, this design enables the higher-level feature branch to acquire stronger discriminative capabilities, thereby leading to more confident category-specific visual representations.

The process is recursively applied to the next stage:

$$Z_{s_{\min+1}} = \operatorname{concat}([C_{\operatorname{cls},s_{\min+1},1}; C_{\operatorname{cls},s_{\min+1},2}; C_{\operatorname{cls},s_{\min+1},3}], \\ [V_{s_{\min+1},1,1}; \dots; V_{s_{\min+1},N+1,P^2}], [V''_{s_{\min},1,1}; \dots; V''_{s_{\min},N+1,P^2}]),$$

$$(11)$$

This procedure is iteratively applied until the highest-level feature branch completes its reasoning process.

3.3 VISION-LANGUAGE CONTRASTIVE LEARNING

Our proposed Vision-Language Contrastive Learning Module based on Multi-grained Text Features (VLCL-MG) constrains the inter-class differences of visual features to align with real-world distinctions by introducing intermediate category constraints, which are primarily implemented through the model structure.

Additionally, in terms of loss computation, since FGVC tasks involve highly similar subcategories, we argue that absolute model outputs can unnecessarily enlarge inter-class distances. If the top-scoring category is incorrect, the actual category score may rank lower due to small score differences among other categories. To mitigate this, we employ label smoothing (Szegedy et al., 2016) for regularization. Meanwhile, given the limited number of samples, models are prone to overfitting, making it crucial to learn more from hard-to-classify samples. Therefore, focal loss (Lin et al., 2017) is also essential. Based on these considerations, we propose the Focal-Smooth Contrastive Loss as a complement to our model structure.

Specifically, we first obtain intermediate categories for fine-grained labels. For example, between the coarse-grained category <code>airplane</code> and the fine-grained category <code>Boeing 737-200</code>, intermediate categories include <code>narrow-body airliner</code> and <code>twinjet</code>. This expert knowledge can be efficiently obtained with minimal effort, requiring only a one-time retrieval per category rather than per-image annotation. For most datasets, we employ <code>ChatGPT-4o</code> for semi-automated knowledge retrieval, whereas for the NABirds dataset, we directly utilize the dataset's built-in class hierarchy.

The multi-grained textual labels corresponding to each fine-grained category are represented as

$$\mathbf{t}_{\text{cls}} = \{a_{n_A}, b_{n_B}, f_{n_F}\} \in \mathbb{R}^C, \tag{12}$$

where a_{n_A} and b_{n_B} represent two types of intermediate categories, and f_{n_F} corresponds to fine-grained categories, C represents the number of fine-grained categories. These labels are processed by the ClearCLIP text encoder, yielding multi-grained text features:

$$\mathbf{T}_{n_F} = f_{\text{text}}(\mathbf{t}_{\text{cls}}),\tag{13}$$

To prevent text features of all grained levels from clustering too closely in the embedding space, we subtract the coarse-grained category feature $f_{\text{text}}(c)$ from the multi-grained text features \mathbf{T}_{n_F} and then apply normalization.

$$\mathbf{T'}_{n_F} = \text{norm}(\mathbf{T}_{n_F} - f_{\text{text}}(c)), \tag{14}$$

This operation ensures a more discriminative distribution of text embeddings across different grained levels. To avoid redundant computations, all category label texts are first processed by the text encoder, and then rearranged according to the intermediate-grained categories corresponding to each fine-grained label, as detailed in fig. 2.

The predicted probability distribution P_s is obtained by applying the softmax function to the similarity score between I_s , extracted from eq. (10) at the s-stage, and T'_{n_F} , formulated as:

$$\mathbf{P_s} = \sigma(\boldsymbol{\tau} \odot (\mathbf{I_s} \mathbf{T'}_{n_F}^T) + \boldsymbol{\beta}), \tag{15}$$

where $\tau \in \mathbb{R}^C$ is a learnable temperature scaling parameter, $\beta \in \mathbb{R}^C$ is a learnable bias term, and $\sigma(\cdot)$ denotes the softmax function for normalization. Therefore, the Focal-Smooth Contrastive Loss $\mathcal{FSL}_s(\mathbf{P_s}, y)$ at the s-stage can be expressed as follows:

$$\mathcal{FSL}_{s}(\mathbf{P_{s}}, y) = -\sum_{n'=1}^{N+1} \sum_{c=1}^{C} (1 - p_{s,n',c})^{\gamma} \cdot \tilde{y}_{s,n',c} \log p_{s,n',c},$$
(16)

where n' represents the n'-th part branch, and s denotes the s-th stage. The parameter γ is the focusing factor that adjusts the impact of misclassified examples. The expression for $\tilde{y}_{s,n',c}$ is given as follows:

$$\tilde{y}_{s,n',c} = \begin{cases} \epsilon_s, & \text{if } c = y_{s,n',c}, \\ 1 - \epsilon_s, & \text{otherwise,} \end{cases}$$
(17)

where ϵ_s is a hyperparameter representing the smoothing noise factor, whose value gradually approaches 1 as the stage s increases. This ensures that the MMBPR becomes more confident during progressive reasoning.

The final loss is computed as $\mathcal{FSL}_s(\mathbf{P_s}, y)$ weighted by a stage-dependent coefficient that increases with s, serving the same role as ϵ_s . It is formulated as:

$$\mathcal{L}_{\text{final}} = \sum_{s=s_{\text{min}}}^{4} \tilde{\epsilon}_s \cdot \mathcal{FSL}_s(\mathbf{P_s}, y), \tag{18}$$

where $\tilde{\epsilon}_s$ is the weight coefficient associated with stage s, which increases as s increases.

4 EXPERIMENTS

4.1 EXPERIMENTAL SETUP

Table 1: Statistics of benchmark datasets.

Dataset	Class	Train	Test
FGVC Aircraft (AIR)	100	6,667	3,333
Stanford Dogs (DOG)	120	12,000	8,580
Stanford Cars (CAR)	196	8,144	8,041
CUB-200-2011 (CUB)	200	5,994	5,794
NABirds (NAB)	555	23,929	24,633

Datasets We comprehensively evaluate PSCL on the FGVC Aircraft (Maji et al., 2013), Stanford Dogs(Khosla et al., 2011), Stanford Cars (Krause et al., 2013), CUB-200-2011 (Wah et al., 2011) and NABirds (Van Horn et al., 2015) datasets, which are widely used FGVC benchmarks. In all experiments, we do not utilize part annotations and follow the same train/test split. The details of the five datasets are presented in table 1.

Implementation Details We adopt ResNet-50 (He et al., 2016), Vision Transformer (Dosovitskiy et al., 2021), and Swin Transformer (Liu et al., 2021) as the backbone architectures. The input image resolutions are set to 448×448 for ResNet-50 (RN50), 518×518 for Vision Transformer (ViT-B), and 384×384 for Swin Transformer (Swin-B).

During training, we apply standard data augmentation techniques, including random cropping, random erasing, horizontal flipping, Gaussian blur, color jittering, and rotation. All models are trained for 100 epochs using the AdamW optimizer with a batch size of 16 and a weight decay of 0.01. The initial learning rate is set to 1×10^{-4} for RN50 and 1×10^{-5} for both ViT-B and Swin-B. A warm-up phase of 10 epochs is applied, and the learning rate follows a cosine annealing schedule.

The focusing parameter γ is set to 4, and the smoothing noise factor ϵ_s follows [0.6, 0.7, 0.8, 0.9], while $\tilde{\epsilon}_s$ is set to [0.1, 0.2, 0.4, 1.0].

4.2 Comparison with Other Methods

We evaluate our method on five benchmark datasets using three backbone architectures and compare it with state-of-the-art models, as summarized in table 2. The results demonstrate the superior performance and strong generalization ability of PSCL across diverse FGVC benchmarks. PSCL consistently achieves state-of-the-art or highly competitive accuracy across all datasets and backbones (RN50, ViT-B, and Swin-B). It delivers substantial improvements under Transformer-based

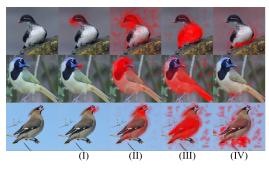
Table 2: Performance comparison on FGVC benchmark datasets (Accuracy %). The best results for each dataset are highlighted in bold.

Method	Backbone	AIR	CAR	CUB	NAB	DOG
CMN (Deng et al., 2022)	RN50	93.8	94.9	88.2	87.8	_
P2P-Net (Yang et al., 2022)	RN50	94.2	95.4	90.2	_	_
GDSMP-Net (Ke et al., 2023)	RN50	94.4	95.3	89.9	89.0	_
SIA-Net (Wang et al., 2023c)	RN50	94.3	95.5	90.7	_	_
PSCL (ours)	RN50	95.1	95.6	89.1	89.0	90.1
TransFG (He et al., 2022)	ViT-B	_	94.8	91.7	90.8	92.3
MpT-Trans (Wang et al., 2023a)	ViT-B	92.2	93.8	92.0	91.3	_
ACC-ViT (Zhang et al., 2024)	ViT-B	_	94.9	91.8	91.4	92.9
MP-FGVC (Jiang et al., 2024)	ViT-B	_	_	91.8	91.0	91.0
PSCL (ours)	ViT-B	96.5	96.4	92.3	93.7	92.3
ViT-NeT (Kim et al., 2022)	Swin-B	_	95.0	91.6	90.9	90.3
TransIFC+ (Liu et al., 2023)	Swin-B	_	_	91.0	90.9	_
HERBS (Chou et al., 2023)	Swin-B	_	_	92.3	93.0	_
CSQA-Net (Xu et al., 2025)	Swin-B	94.7	95.6	92.6	92.3	_
PSCL (ours)	Swin-B	95.3	95.5	93.0	93.8	94.7

backbones, and remains competitive under the CNN-based RN50, particularly on AIR and CAR datasets. These results highlight PSCL's adaptability to different architectures and its effectiveness in capturing both local and structural discriminative cues. Furthermore, its consistent performance across datasets underscores its robustness. Notably, on the large-scale NAB dataset, the availability of an accurate and professionally curated category hierarchy enables precise intermediate-level grouping, further enhancing accuracy and demonstrating the effectiveness of the VLCL-MG module. The strong performance on multiple non-rigid datasets such as DOG and NAB demonstrates that PSCL effectively models the characteristics of non-rigid objects.

4.3 EFFECTIVENESS OF MODULE OPERATION

Locating Relevant Parts We resize eq. (2) to match the original image dimensions for visualizing the part localization results. As observed in fig. 4, the PLM structure effectively identifies the locations of the parts.



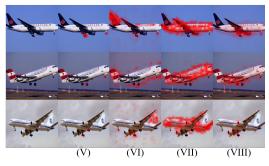


Figure 4: Part score visualization. PLM uses the following textual prompts: (I) mouth; (II) head; (III) body; (IV) foot; (V) landing gear; (VI) tail; (VII) fuselage; (VIII) engine.

Reverse-key Scale-aware Attention Fusion Module To assess the effectiveness of our proposed ReSAF module, we conduct a comparative study with two alternative intermediate mechanisms: a multilayer perceptron (MLP) and cross-attention. All experiments are performed using the RN50 backbone on the AIR dataset. The quantitative results, summarized in table 3, demonstrate that ReSAF consistently outperforms the other two variants, highlighting its superior capability in capturing scale-aware feature interactions.

Table 3: Performance comparison of different intermediate mechanisms on the AIR dataset. Accuracy (%) is reported.

Intermediate Mechanism	Accuracy (%)
MLP	94.71
Cross-Attention	94.99
ReSAF (Ours)	95.14

Hyperparameter Selection All hyperparameters except the learning rate were searched and selected exclusively on the AIR dataset with RN50. Results are shown in table 4, table 5 and fig. 5. The trends of $\tilde{\epsilon}_s$ and γ suggest that progressive inference improves performance, and increasing the focus on hard samples via γ further enhances results. We note that if hyperparameter tuning were performed specifically for a target dataset, our proposed PSCL could potentially achieve even better performance.

Table 4: Accuracy (%) on AIR dataset (RN50) under smoothing noise factor ϵ_s .

ϵ_s	Accuracy (%)
[1.0, 1.0, 1.0, 1.0]	89.92
[0.4, 0.6, 0.8, 1.0]	94.22
[0.3, 0.5, 0.7, 0.9]	94.83
[0.6, 0.7, 0.8, 0.9]	95.14

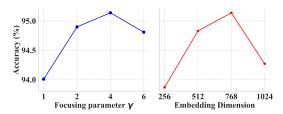


Figure 5: Accuracy (%) on the AIR dataset (RN50) under different settings. Left: focusing parameter γ ; Right: encoder hidden dimension.

Table 5: Accuracy (%) on AIR dataset (RN50) under multi-scale loss weight coefficient $\tilde{\epsilon}_s$.

$ ilde{\epsilon}_s$	Accuracy (%)
[0.0, 0.0, 0.0, 1.0]	94.65
[0.1, 0.2, 0.2, 1.0]	94.77
[0.1, 0.2, 0.4, 1.0]	95.14
[0.2, 0.4, 0.4, 1.0]	94.89

Ablation Studies The ablation results in table 6 across RN50, Swin-B, and ViT-B on CUB, AIR, and CAR confirm the effectiveness of each proposed module. PLM and VLCL-MG alone bring notable gains, highlighting their individual strengths in part localization and semantic alignment. MMBPR further enhances performance via multi-scale reasoning. While their roles differ, all three aim to improve semantic consistency, and may exhibit overlapping effects in low-redundancy settings. Still, the full model consistently achieves the best results, demonstrating the complementary benefits of combining all components.

Table 6: Ablation study on three FGVC datasets using different backbones. Accuracy (%) is reported for each configuration. The best results for each column are highlighted in bold. Features are indicated by a check mark (\checkmark) or a cross (\checkmark).

PLM	MMBPR	VLCL-MG	CUB	RN50 AIR	CAR		Swin-B AIR		CUB	ViT-B AIR	CAR
X	Х	Х	85.09	91.56	91.90	92.32	94.14	94.86	90.23	94.57	95.60
1	X	X	88.82	94.54	95.46	92.68	94.60	94.74	91.99	96.13	96.05
/	✓	X	89.09	94.54	95.54	92.51	95.05	95.04	92.34	96.19	96.17
X	Х	✓	87.90	94.39	95.32	92.65	94.87	95.51	90.94	95.08	96.26
✓	✓	✓	89.13	95.14	95.59	93.01	95.32	95.54	92.34	96.48	96.44

4.4 ADDITIONAL COMPUTATIONAL COST

Overall, although PSCL introduces additional computation, the total computational cost is comparable to that of many modern large-scale Transformer models and can similarly benefit from accelerated computation techniques, and is therefore within an acceptable range.

Table 7: Computational cost.

Component	GFLOPs
ClearCLIP	17.35
MMBPR	$15.72 \times (N+1)$
ReSAF	5.78
others	Negligible

5 CONCLUSION

We introduce Part-level Semantic-guided Contrastive Learning (PSCL), a framework for Fine-Grained Visual Classification that jointly models part-level details and spatial relations for both rigid and non-rigid objects. PSCL employs a Part Localization Module (PLM) with ClearCLIP for semantically guided, interpretable part extraction, a Multi-scale Multi-part Branch Progressive Reasoning (MMBPR) module to fuse fine-grained and global features, and a Visual-Language Contrastive Learning module with Multi-grained Text Features (VLCL-MG) to align subcategories via intermediate-level semantics. Experiments on five FGVC benchmarks demonstrate PSCL's robust performance and strong generalization.

ETHICS STATEMENT

Our work focuses on fine-grained visual classification using publicly available datasets. No human subjects, personally identifiable information, or sensitive content are involved. All datasets employed, including AIR, CAR, NABirds, DOG, and others, are used strictly for research purposes in accordance with their respective licenses. We acknowledge the potential societal impacts of deploying FGVC models, such as reinforcing biases present in the training data, and emphasize that PSCL should be applied responsibly, with consideration for fairness and ethical implications.

REPRODUCIBILITY STATEMENT

To ensure reproducibility, we provide detailed descriptions of all model components, including the Part Localization Module (PLM), Multi-scale Multi-part Branch Progressive Reasoning (MMBPR), and Visual-Language Contrastive Learning with Multi-grained Text Features (VLCL-MG). Hyperparameters, training procedures, and evaluation protocols are specified in the manuscript. The code, pre-trained weights, and scripts for data processing and evaluation are available upon request, enabling independent verification of reported results. All experiments were conducted using standard hardware and frameworks to facilitate replication.

REFERENCES

- Ardhendu Behera, Zachary Wharton, Pradeep RPG Hewage, and Asish Bera. Context-aware attentional pooling (cap) for fine-grained visual classification. In *Proceedings of the AAAI conference on artificial intelligence*, volume 35, pp. 929–937, 2021.
- Steve Branson, Grant Van Horn, Serge Belongie, and Pietro Perona. Bird species categorization using pose normalized deep convolutional nets. *arXiv preprint arXiv:1406.2952*, 2014.
- Po-Yung Chou, Yu-Yung Kao, and Cheng-Hung Lin. Fine-grained visual classification with high-temperature refinement and background suppression. *arXiv preprint arXiv:2303.06442*, 2023.
- Weijian Deng, Joshua Marsh, Stephen Gould, and Liang Zheng. Fine-grained classification via categorical memory networks. *IEEE Transactions on Image Processing*, 31:4186–4196, 2022.
- Yifeng Ding, Zhanyu Ma, Shaoguo Wen, Jiyang Xie, Dongliang Chang, Zhongwei Si, Ming Wu, and Haibin Ling. Ap-cnn: Weakly supervised attention pyramid convolutional neural network for fine-grained visual classification. *IEEE Transactions on Image Processing*, 30:2826–2836, 2021.
- Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale, 2021. URL https://arxiv.org/abs/2010.11929.
- Ruoyi Du, Dongliang Chang, Ayan Kumar Bhunia, Jiyang Xie, Zhanyu Ma, Yi-Zhe Song, and Jun Guo. Fine-grained visual classification via progressive multi-granularity training of jigsaw patches. In *European Conference on Computer Vision*, pp. 153–168. Springer, 2020.
- Ju He, Jie-Neng Chen, Shuai Liu, Adam Kortylewski, Cheng Yang, Yutong Bai, and Changhu Wang. Transfg: A transformer architecture for fine-grained recognition. In *Proceedings of the AAAI conference on artificial intelligence*, volume 36, pp. 852–860, 2022.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 770–778, 2016.
- Xin Jiang, Hao Tang, Junyao Gao, Xiaoyu Du, Shengfeng He, and Zechao Li. Delving into multi-modal prompting for fine-grained visual classification. In *Proceedings of the AAAI conference on artificial intelligence*, volume 38, pp. 2570–2578, 2024.
- Xiao Ke, Yuhang Cai, Baitao Chen, Hao Liu, and Wenzhong Guo. Granularity-aware distillation and structure modeling region proposal network for fine-grained image classification. *Pattern Recognition*, 137:109305, 2023.

- Aditya Khosla, Nityananda Jayadevaprakash, Bangpeng Yao, and Fei-Fei Li. Novel dataset for finegrained image categorization: Stanford dogs. In *Proc. CVPR workshop on fine-grained visual* categorization (FGVC), volume 2, 2011.
 - Sangwon Kim, Jaeyeal Nam, and Byoung Chul Ko. Vit-net: Interpretable vision transformers with neural tree decoder. In *International conference on machine learning*, pp. 11162–11172. PMLR, 2022.
 - Jonathan Krause, Michael Stark, Jia Deng, and Li Fei-Fei. 3d object representations for fine-grained categorization. In *Proceedings of the IEEE international conference on computer vision work-shops*, pp. 554–561, 2013.
 - Mengcheng Lan, Chaofeng Chen, Yiping Ke, Xinjiang Wang, Litong Feng, and Wayne Zhang. Clearclip: Decomposing clip representations for dense vision-language inference. In *European Conference on Computer Vision*, pp. 143–160. Springer, 2024.
 - Siyuan Li, Li Sun, and Qingli Li. Clip-reid: exploiting vision-language model for image reidentification without concrete text labels. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 37, pp. 1405–1413, 2023.
 - Tsung-Yi Lin, Priya Goyal, Ross Girshick, Kaiming He, and Piotr Dollár. Focal loss for dense object detection. In 2017 IEEE International Conference on Computer Vision (ICCV), pp. 2999–3007, 2017. doi: 10.1109/ICCV.2017.324.
 - Tsung-Yu Lin, Aruni RoyChowdhury, and Subhransu Maji. Bilinear cnn models for fine-grained visual recognition. In *Proceedings of the IEEE international conference on computer vision*, pp. 1449–1457, 2015.
 - Hai Liu, Cheng Zhang, Yongjian Deng, Bochen Xie, Tingting Liu, and You-Fu Li. Transifc: Invariant cues-aware feature concentration learning for efficient fine-grained bird image classification. *IEEE Transactions on Multimedia*, 27:1677–1690, 2023.
 - Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo. Swin transformer: Hierarchical vision transformer using shifted windows. In *Proceedings of the IEEE/CVF international conference on computer vision*, pp. 10012–10022, 2021.
 - Subhransu Maji, Esa Rahtu, Juho Kannala, Matthew Blaschko, and Andrea Vedaldi. Fine-grained visual classification of aircraft. *arXiv preprint arXiv:1306.5151*, 2013.
 - Zhuang Miao, Xun Zhao, Jiabao Wang, Yang Li, and Hang Li. Complemental attention multi-feature fusion network for fine-grained classification. *IEEE Signal Processing Letters*, 28:1983–1987, 2021.
 - Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pp. 8748–8763. PMLR, 2021.
 - Guolei Sun, Hisham Cholakkal, Salman Khan, Fahad Khan, and Ling Shao. Fine-grained recognition: Accounting for subtle differences between similar classes. In *Proceedings of the AAAI conference on artificial intelligence*, volume 34, pp. 12047–12054, 2020.
 - Ming Sun, Yuchen Yuan, Feng Zhou, and Errui Ding. Multi-attention multi-class constraint for fine-grained image recognition. In *Proceedings of the european conference on computer vision (ECCV)*, pp. 805–821, 2018.
 - Christian Szegedy, Vincent Vanhoucke, Sergey Ioffe, Jon Shlens, and Zbigniew Wojna. Rethinking the inception architecture for computer vision. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 2818–2826, 2016.
 - Grant Van Horn, Steve Branson, Ryan Farrell, Scott Haber, Jessie Barry, Panos Ipeirotis, Pietro Perona, and Serge Belongie. Building a bird recognition app and large scale dataset with citizen scientists: The fine print in fine-grained dataset collection. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 595–604, 2015.

- C. Wah, S. Branson, P. Welinder, P. Perona, and S. Belongie. The caltech-ucsd birds-200-2011 dataset. Technical Report CNS-TR-2011-001, California Institute of Technology, 2011.
 - Chuanming Wang, Huiyuan Fu, and Huadong Ma. Multi-part token transformer with dual contrastive learning for fine-grained image classification. In *Proceedings of the 31st ACM International Conference on Multimedia*, pp. 7648–7656, 2023a.
 - Shijie Wang, Jianlong Chang, Haojie Li, Zhihui Wang, Wanli Ouyang, and Qi Tian. Open-set fine-grained retrieval via prompting vision-language evaluator. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 19381–19391, 2023b.
 - Shijie Wang, Zhihui Wang, Haojie Li, Jianlong Chang, Wanli Ouyang, and Qi Tian. Semantic-guided information alignment network for fine-grained image recognition. *IEEE Transactions on Circuits and Systems for Video Technology*, 33(11):6558–6570, 2023c.
 - Qin Xu, Sitong Li, Jiahui Wang, Bo Jiang, Bin Luo, and Jinhui Tang. Context-semantic quality awareness network for fine-grained visual categorization. *Pattern Recognition*, pp. 112033, 2025.
 - Xuhui Yang, Yaowei Wang, Ke Chen, Yong Xu, and Yonghong Tian. Fine-grained object classification via self-supervised pose alignment. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 7399–7408, 2022.
 - Ying Yu, Wei Wei, Cairong Zhao, Jin Qian, and Enhong Chen. Structural feature enhanced transformer for fine-grained image recognition. *Pattern Recognition*, pp. 111955, 2025.
 - Fan Zhang, Meng Li, Guisheng Zhai, and Yizhao Liu. Multi-branch and multi-scale attention learning for fine-grained visual categorization. In *MultiMedia Modeling: 27th International Conference, MMM 2021, Prague, Czech Republic, June 22–24, 2021, Proceedings, Part I 27*, pp. 136–147. Springer, 2021.
 - Lianbo Zhang, Shaoli Huang, Wei Liu, and Dacheng Tao. Learning a mixture of granularity-specific experts for fine-grained categorization. In *Proceedings of the IEEE/CVF international conference* on computer vision, pp. 8331–8340, 2019.
 - Ning Zhang, Jeff Donahue, Ross Girshick, and Trevor Darrell. Part-based r-cnns for fine-grained category detection. In *Computer Vision–ECCV 2014: 13th European Conference, Zurich, Switzerland, September 6-12, 2014, Proceedings, Part I 13*, pp. 834–849. Springer, 2014.
 - Zhicheng Zhang, Hao Tang, and Jinhui Tang. Multi-scale activation, selection, and aggregation: Exploring diverse cues for fine-grained bird recognition. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 39, pp. 10385–10393, 2025.
 - Zi-Chao Zhang, Zhen-Duo Chen, Yongxin Wang, Xin Luo, and Xin-Shun Xu. A vision transformer for fine-grained classification by reducing noise and enhancing discriminative information. *Pattern Recognition*, 145:109979, 2024.
 - Yifan Zhao, Jia Li, Xiaowu Chen, and Yonghong Tian. Part-guided relational transformers for fine-grained visual recognition. *IEEE Transactions on Image Processing*, 30:9470–9481, 2021.
 - Heliang Zheng, Jianlong Fu, Tao Mei, and Jiebo Luo. Learning multi-attention convolutional neural network for fine-grained image recognition. In *Proceedings of the IEEE international conference on computer vision*, pp. 5209–5217, 2017.
 - Heliang Zheng, Jianlong Fu, Zheng-Jun Zha, and Jiebo Luo. Learning deep bilinear transformation for fine-grained image representation. *Advances in Neural Information Processing Systems*, 32, 2019.
 - Chong Zhou, Chen Change Loy, and Bo Dai. Extract free dense labels from clip. In *European Conference on Computer Vision*, pp. 696–712. Springer, 2022.

A THE USE OF LARGE LANGUAGE MODELS

In this work, we leverage large language models (LLMs), specifically ChatGPT-40, to retrieve domain-specific semantic knowledge for guiding part-level selection in fine-grained visual classification. Beyond constructing intermediate category hierarchies, LLMs are also employed to assist in code development and manuscript refinement. All outputs from the LLM are carefully verified to ensure accuracy.

B TRAINING AND HYPERPARAMETER SETTINGS FOR DIFFERENT BACKBONES

Hyperparameter Setting	ResNet-50	ViT-B	Swin-B
Input resolution	448×448	518×518	384×384
Batch size	16	16	16
Weight decay	0.01	0.01	0.01
Optimizer	AdamW	AdamW	AdamW
Optimizer β	(0.9, 0.95)	(0.9, 0.95)	(0.9, 0.95)
Optimizer ϵ	1e-8	1e-8	1e-8
Initial learning rate	1e-4	1e-5	1e-5
Learning rate schedule	Cosine annealing	Cosine annealing	Cosine annealing
Warm-up epochs	10	10	10
Epochs	100	100	100
Focusing parameter γ	4	4	4
Smoothing noise factor ϵ_s	[0.6, 0.7, 0.8, 0.9]	[0.6, 0.7, 0.8, 0.9]	[0.6, 0.7, 0.8, 0.9]
Multi-scale loss weight $\tilde{\epsilon}_s$	[0.1, 0.2, 0.4, 1.0]	[0.1, 0.2, 0.4, 1.0]	[0.1, 0.2, 0.4, 1.0]

C PART TEXT FOR DIFFERENT DATASETS

Dataset	Part Text
AIR	background of a plane, tail of a plane, logo of a plane, engine of a plane, landing gear of a plane, fuselage of a plane
CUB	background of a bird, head of a bird, foot of a bird, body of a bird, mouth of a bird
CAR	background of a car, head of a car, body of a car, back of a car
NAB DOG	background of a bird, head of a bird, foot of a bird, body of a bird, mouth of a bird background of a dog, head of a dog, foot of a dog, body of a dog

D OTHER RESULTS.

Table 8: Classification accuracy (%) on fine-grained datasets using different embedding/masking strategies.

Method		ResNet50			ViT-B		
Troutou	CUB	Aircraft	Car	CUB	Aircraft	Car	
Random text embeddings (F_{text})	88.93	94.75	95.38	92.04	95.45	96.36	
Random masking $(S(\mathbf{k}))$	88.12	94.93	95.25	91.76	95.26	96.26	
Part text embeddings	89.13	95.14	95.59	92.34	96.48	96.44	

We conducted experiments in which either F_{text} or $S(\mathbf{k})$ was randomized. Both random strategies can be viewed as mutually exclusive data augmentation methods based on random masking.

However, Random Text Embeddings tend to occlude semantically similar regions, whereas Random Masking hides regions randomly. Our proposed PSCL architecture demonstrates considerable robustness: thanks to the MMBPR and VLCL-MG modules, the model can still learn to focus on relevant regions autonomously. Nevertheless, providing targeted human guidance could further improve the efficiency of this process.

Table 9: Effect of intermediate-category text annotations on NAB classification performance using ViT-B. Accuracy (%) is reported.

Intermediate-Category Text	Accuracy (%)
Expert annotations (precise hierarchy)	93.74
Generated via ChatGPT-4o (semi-automatic)	93.48
Random-text control group	92.83

We posit that the NAB dataset benefits substantially from its inherent, precise hierarchical category structure, resulting in a significant performance boost. Accordingly, we employed ChatGPT-40 in a semi-automatic manner to generate intermediate-category text annotations, while also creating a random-text control group, and conducted comparative experiments using the ViT-B backbone. The results demonstrate that accurate expert annotations effectively activate the VLCL-MG module, yet even the generated intermediate-category text can improve classification accuracy to a certain extent.

E EXAMPLES OF SOME INTERMEDIATE CATEGORIES

Table 10: Intermediate classes for the AIR dataset

Fine-grained	Intermediate-grained 1	Intermediate-grained 2
737-900	narrow-body airliner	twinjet
747-100	wide-body airliner	four-engined jet aircraft
A330-300	wide-body airliner	twinjet
A340-200	wide-body airliner	four-engined jet aircraft
Cessna 525	business jet	twinjet
Challenger 600	business jet	twinjet
DC-10	wide-body airliner	trijet
DC-3	cargo aircraft	twin-turboprop
Gulfstream V	business jet	twinjet
Hawk T1	light aircraft	single-engine jet
Il-76	cargo aircraft	four-engined jet aircraft
L-1011	wide-body airliner	trijet
MD-11	wide-body airliner	trijet

Table 11: Intermediate classes for the CUB dataset

Fine-grained	Intermediate-grained 1	Intermediate-grained 2
Frigatebird	Seabirds	Waterbirds
Gadwall	Ducks	Waterbirds
American Goldfinch	Finches	Songbirds
Boat-tailed Grackle	Grackles	Songbirds
American Crow	Crows	Corvids
Fish Crow	Crows	Corvids
Black-billed Cuckoo	Cuckoos	Songbirds
Rusty Blackbird	Blackbirds	Songbirds
Yellow-headed Blackbird	Blackbirds	Songbirds
Indigo Bunting	Buntings	Songbirds

Table 12: Intermediate classes for the CAR dataset

Fine-grained	Intermediate-grained 1	Intermediate-grained 2
Audi S4 Sedan 2007	Sedan	Performance Vehicle
Audi TT RS Coupe 2012	Coupe	Performance Vehicle
BMW ActiveHybrid 5 Sedan 2012	Sedan	Hybrid Vehicle
BMW 1 Series Convertible 2012	Convertible	Luxury Vehicle
BMW 1 Series Coupe 2012	Coupe	Luxury Vehicle
Acura Integra Type R 2001	Coupe	Performance Vehicle
Acura ZDX Hatchback 2012	Hatchback	Luxury Vehicle
Aston Martin V8 Vantage Convertible 2012	Convertible	Luxury Vehicle
Chrysler Crossfire Convertible 2008	Convertible	Performance Vehicle
Chrysler PT Cruiser Convertible 2008	Convertible	Family Car
Daewoo Nubira Wagon 2002	Wagon	Family Car

Table 13: Intermediate classes for the DOG dataset

Fine-grained	Intermediate-grained 1	Intermediate-grained 2
Blenheim Spaniel	Sporting	Spaniel
Papillon	Toy	Toy-group
Toy Terrier	Toy	Terrier-toy
Rhodesian Ridgeback	Hound	Sighthound
Afghan Hound	Hound	Sighthound
Weimaraner	Sporting	Pointer
Staffordshire Bullterrier	Terrier	Bull-type
Cocker Spaniel	Sporting	Spaniel
Pug	Toy	Toy-group
Great Pyrenees	Working	Working-group
Irish Water Spaniel	Sporting	Spaniel
Kuvasz	Working	Working-group
Groenendael	Herding	Herding-group

Table 14: Intermediate classes for the NAB dataset

Fine-grained	Intermediate-grained 1	Intermediate-grained 2
Black-bellied Whistling-Duck	Black-bellied Whistling-Duck	Ducks, Geese, and Swans
Semipalmated Plover	Semipalmated Plover	Plovers, Sandpipers, and Allies
American White Pelican	American White Pelican	Pelicans, Herons, Ibises, and Allies
Killdeer	Killdeer	Plovers, Sandpipers, and Allies
Chimney Swift	Chimney Swift	Swifts and Hummingbirds
American Oystercatcher	American Oystercatcher	Plovers, Sandpipers, and Allies
Ross's Goose	Ross's Goose	Ducks, Geese, and Swans
Barn Owl	Barn Owl	Owls
Turkey Vulture	Turkey Vulture	Hawks, Kites, Eagles, and Allies
Brown Pelican	Brown Pelican	Pelicans, Herons, Ibises, and Allies
Scaled Quail	Scaled Quail	Grouse, Quail, and Allies
Rock Pigeon	Rock Pigeon	Pigeons and Doves
Black-necked Stilt	Black-necked Stilt	Plovers, Sandpipers, and Allies