

# Gaze-Supervised Hierarchical Attention Networks for Fine-Grained Visual Classification

Edwin R. Hancock

School of Computer Science and Engineering, State Key Laboratory of Software, University of York, U.K.

## Abstract

Fine-grained visual classification (FGVC) becomes especially challenging when categories are organized hierarchically and the discriminative cues shrink from global shapes (order/family) to tiny parts (genus/species). Existing hierarchy-aware methods such as CHRF learn level-specific attentions implicitly, but they only use human gaze as a post-hoc validation signal, leaving a rich source of supervision unused. In this work we introduce GS-HAN, a gaze-supervised hierarchical attention network that explicitly aligns model attention with human fixation patterns at every level of the taxonomy. GS-HAN builds on a backbone feature extractor and CHRF-style region feature mining, but augments each hierarchy level with gaze-conditioned attention heads and a Hierarchical Gaze Alignment Loss that combines KL divergence and cosine similarity to match human gaze distributions. We further retain cross-hierarchical orthogonal fusion so that coarse-level, gaze-aligned context can enhance fine-level recognition. Evaluations on CUB-200-2011 with ARISTO gaze, as well as on Butterfly-200, VegFru, FGVC-Aircraft, and Stanford Cars, show that GS-HAN consistently outperforms strong FGVC baselines and hierarchy-aware methods, achieving 90.8% on CUB and clear gains at the most fine-grained (species) level. Ablations verify that (i) direct gaze supervision—not just hierarchy—drives the improvements, (ii) our loss improves quantitative gaze–attention similarity, and (iii) even partial gaze availability yields benefits. The results demonstrate that human gaze is an effective, underexploited supervisory signal for hierarchical FGVC, improving both accuracy and interpretability.

**Keywords:** fine-grained visual classification, hierarchical classification, human gaze supervision, attention alignment, CHRF, interpretability, bird recognition, ARISTO dataset, cross-hierarchical fusion, deep learning

## 1. Introduction

Fine-grained visual classification (FGVC) aims to distinguish among subcategories that belong to the same basic-level category (e.g., different bird species, car models, aircraft types). Unlike generic image classification, where inter-class differences are often large and driven by global appearance, FGVC typically encounters very subtle inter-class variations and, at the same time, relatively large intra-class variations caused by pose, viewpoint, occlusion, background clutter, or illumination changes. This combination makes the task intrinsically difficult: the model must learn to focus on small, highly discriminative visual cues while remaining robust to nuisances. Existing lines of work have therefore concentrated either on designing architectures that extract more discriminative part-level features

[1] or on learning better feature representations that amplify fine-grained differences [2]. While these approaches have improved recognition accuracy, many of them treat object categories as flat labels and ignore the hierarchical structure that is naturally present in many real-world taxonomies.

In human cognition, object recognition is rarely a single-step, flat decision. Humans tend to process visual categories in a top-down manner: first identifying a coarse group (e.g., “this is a bird”), then narrowing down to a finer group (e.g., “this is a shorebird”), and finally deciding the exact species (e.g., “this is a red knot”). Psychovisual studies and recent computational analyses [3] show that this progression is not only semantic but also attentional. At coarse levels such as *order* or *family*, humans rely more on holistic shape, overall silhouette, and global color distribution. At finer levels such as *species*, they shift attention to small, distinctive parts such as the beak, tail tips, eye rings, or subtle wing patterns. In other words, human recognition is *hierarchical* and *attention-driven*, and the attended regions change systematically with the granularity of the decision. This observation suggests that an FGVC model that wants to mimic human success should (i) recognize that categories are hierarchically related and (ii) learn attention that is conditioned on the current level of the hierarchy.

Motivated by this intuition, the recent CHRf framework [3, 4] explicitly modeled cross-hierarchical attention through Region Feature Mining (RFM) and Cross-hierarchical Orthogonal Fusion (COF). CHRf demonstrated that if a model is guided to look at the right regions for different levels, it can better disentangle subtle differences. However, there is a key limitation in CHRf and in many related works: the hierarchical attention is *implicitly* learned. That is, the model is encouraged to produce diverse and complementary attentions across levels, but it is never *told* where humans actually look when making those hierarchical judgments. In fact, CHRf uses human gaze patterns only as an external validation signal to show that its learned attention is human-like, not as a supervisory signal during training. This creates a supervision gap: we already have gaze annotations in modern fine-grained datasets (e.g., ARISTO) that describe *how* humans allocate their attention across different granularity levels, but most current FGVC methods do not exploit this rich information directly.

We argue that this is a missed opportunity. Human gaze provides a dense, spatially grounded, and level-aware form of supervision that is highly aligned with the goal of hierarchical FGVC. Unlike class labels, which only tell *what* the object is, gaze tells *where* humans look to decide it. When such gaze is available at multiple hierarchical levels, it can act as an attention teacher: the model can be trained not just to predict the right label but to focus on the same regions that humans deem discriminative for that level. This promises several benefits. First, it can improve accuracy, because attention is steered toward truly informative parts instead of background or spuriously correlated regions. Second, it can improve interpretability, because the model’s attention maps become directly comparable to human fixation maps. Third, it can improve robustness and transfer, because learning to attend in a human-like way can reduce overfitting to dataset-specific biases.

In this paper, we bridge this supervision gap by introducing a *gaze-supervised* hierarchical attention network for fine-grained recognition. Our core idea is simple but powerful: whenever human gaze is available for a given hierarchy level, we explicitly align the model’s attention at that level with the corresponding human fixation map. To make this feasible, we design an attention generation branch that outputs level-specific attention maps, and we introduce a Hierarchical Gaze Alignment Loss that penalizes discrepancies between predicted attention and human gaze. In doing so, the model is trained to *look where humans look* at every step of the hierarchical decision, rather than hoping that such behavior will emerge implicitly from classification loss alone.

Concretely, our framework, termed **GS-HAN** (Gaze-Supervised Hierarchical Attention Network), operates on top of a standard fine-grained backbone but augments it with three key components. First, we adopt a hierarchy-aware feature extraction pipeline that keeps track of representations for coarse-to-fine levels. Second, we generate attention maps conditioned on each level so that the model can attend globally for coarse levels and locally for fine levels, mimicking human patterns reported in [3]. Third, we use human gaze maps—when available—as supervision targets to train these attention maps. This leads to attention distributions that are not only discriminative for the task but also human-consistent. Since gaze data may be noisy or incomplete, our loss is designed to be flexible and to work with sparse or partially observed gaze annotations.

We evaluate our approach on five benchmark fine-grained datasets that either contain human gaze annotations directly or can be augmented with gaze-like supervision. Across all benchmarks, GS-HAN yields consistent improvements over baselines that (i) do not model hierarchy, (ii) model hierarchy but do not use gaze, and (iii) use attention but only implicitly. Beyond raw accuracy, we report quantitative attention interpretability metrics, showing that our predicted attention maps correlate more strongly with human fixation than prior methods. This demonstrates that gaze is not only a convenient validation tool but also an effective training signal.

## 2. Literature Review

### 2.1. Fine-Grained Visual Classification

Fine-grained visual classification has been studied extensively because of its practical importance in domains such as biodiversity monitoring, intelligent transportation, and e-commerce, where the goal is to distinguish visually similar subcategories. Early FGVC methods assumed access to strong supervision such as part or bounding-box annotations and used these to crop informative regions, but such annotations are costly and do not scale. Consequently, a large body of work moved toward weakly supervised or annotation-free learning.

Broadly, existing FGVC approaches can be divided into *feature-centric* and *part-/region-centric* paradigms. Feature-based methods [5, 6] seek to enhance the discriminativeness of global representations. They design losses or architectural mechanisms (e.g., feature erasing, pairwise confusion, multi-branch aggregation) to force the network to mine subtle cues that distinguish fine-grained classes. These methods are simple and end-to-end, but because they rely on a single shared representation, they may fail to explicitly capture which spatial regions are actually discriminative.

Part-based or region-mining methods [7, 8] address this by discovering or localizing informative parts directly from images without dense human annotations. They typically learn attention maps, proposal generators, or transformer tokens that highlight wings, heads, logos, or other distinctive components. By attending to multiple parts and aggregating them, these methods increase robustness and make the model more interpretable. However, most of these approaches treat the label space as *flat*—every class is predicted at the same granularity—so all parts are searched for with the same strength, regardless of whether the decision is coarse (family/order) or fine (species/model).

More recent works have started to leverage *hierarchical* relationships among categories [9, 10] to guide representation learning. The key insight is that visual taxonomies are not arbitrary: birds are grouped into orders, families, genera, and species, and these levels correspond to progressively finer visual distinctions. Incorporating such hierarchy can regularize the feature space and help the model learn from coarse-to-fine signals. CHRf [3] is particularly relevant to us because it explicitly models

cross-hierarchical attention through Region Feature Mining (RFM) and Cross-hierarchical Orthogonal Fusion (COF). CHRf shows that letting different levels attend to different regions improves fine-grained recognition. Nonetheless, even CHRf learns these attentions *implicitly* from classification objectives and does not exploit human gaze as an explicit teacher. Our work builds on this line by keeping the benefits of hierarchy-aware attention while injecting human gaze supervision to make those attentions human-aligned and level-specific.

## 2.2. Human Attention in Vision

Human visual attention has long been recognized as a rich cue for AI systems because it reveals *where* humans look to solve a task. In image captioning and vision-language tasks, attention or eye-tracking data has been used either to evaluate whether model attention is human-like [11] or to guide the model toward salient objects that should be described. In zero-shot and goal-directed learning, human attention can indicate task-relevant regions and filter out distractors [12, 13]. In medical imaging, gaze from radiologists or clinicians has been leveraged to highlight diagnostically relevant regions and to compensate for the scarcity of pixel-level annotations [14, 15].

Within FGVC specifically, several works have explored bringing human attention into the loop, but in ways that differ from ours. Rong et al. [16] used human attention mainly for data augmentation—i.e., to generate better crops or masks—so the attention was an auxiliary tool rather than a supervision target. Yu et al. [17] simulated glance-and-gaze behavior in Vision Transformers so that the model could first get a global impression and then refine on local details, mimicking human inspection. However, these approaches either do not connect attention to the *hierarchical* nature of fine-grained labels or do not explicitly *align* model attention with measured human gaze. Our problem setting is stricter: we want the model to learn not only to attend, but to attend in a way that is consistent with human fixations at each level of the taxonomy. This is precisely what our hierarchical gaze alignment seeks to achieve.

## 2.3. Gaze Supervision in Deep Learning

Eye-tracking and gaze data have been used in deep learning for several purposes. A common line of work is *saliency prediction*, where the goal is to train a model that, given an image, predicts a human fixation map [18]. Here, gaze is the *end task*. Another line integrates gaze into vision-language tasks such as visual question answering (VQA) [19], where human attention helps the model focus on objects that humans considered relevant to the question. In activity recognition or interactive settings, gaze can also be used as a proxy for human intent, guiding the model toward the target of human actions.

However, using gaze as a *direct supervisory signal* for *hierarchical* visual classification remains largely unexplored. Most existing methods either (i) predict gaze from images, (ii) use gaze to filter training regions, or (iii) compare model attention with gaze only at evaluation time as an interpretability metric. None of these fully exploits the situation we highlight in this paper: when we have *hierarchy-aware* gaze (i.e., gaze collected while humans are making coarse-to-fine decisions), we can force the model’s level-specific attention maps to match those human fixations and thus learn truly human-like hierarchical attention.

Our work is, to the best of our knowledge, the first to combine these three ingredients simultaneously: (1) fine-grained classification with explicit hierarchical label structure, (2) attention modules that produce level-specific spatial maps, and (3) human gaze used as *direct* supervision to align those

maps. By introducing the Hierarchical Gaze Alignment Loss and integrating it into a gaze-supervised hierarchical attention network, we convert gaze from a passive validation signal into an active teacher. This closes the gap identified in the introduction and shows that human gaze is not just useful for interpretability, but can also drive performance improvements in FGVC when used properly.

### 3. Approach

In this section we describe our Gaze-Supervised Hierarchical Attention Network (GS-HAN) in detail. Following the motivation in Section ?? (Introduction) and the trends identified in Section ?? (Related Work), our goal is to (i) preserve the desirable hierarchical, region-mining behavior of CHRf [7], (ii) make attention *explicitly* human-aligned at every level where gaze is available, and (iii) keep the model end-to-end trainable with standard classification losses. To that end, we augment the CHRf-style pipeline with a gaze-supervised attention branch and a new loss that enforces level-specific alignment.

#### 3.1. Overview

Our GS-HAN framework builds upon the hierarchical structure of CHRf but introduces crucial modifications for gaze supervision. As shown in Figure 1, the framework consists of three main components:

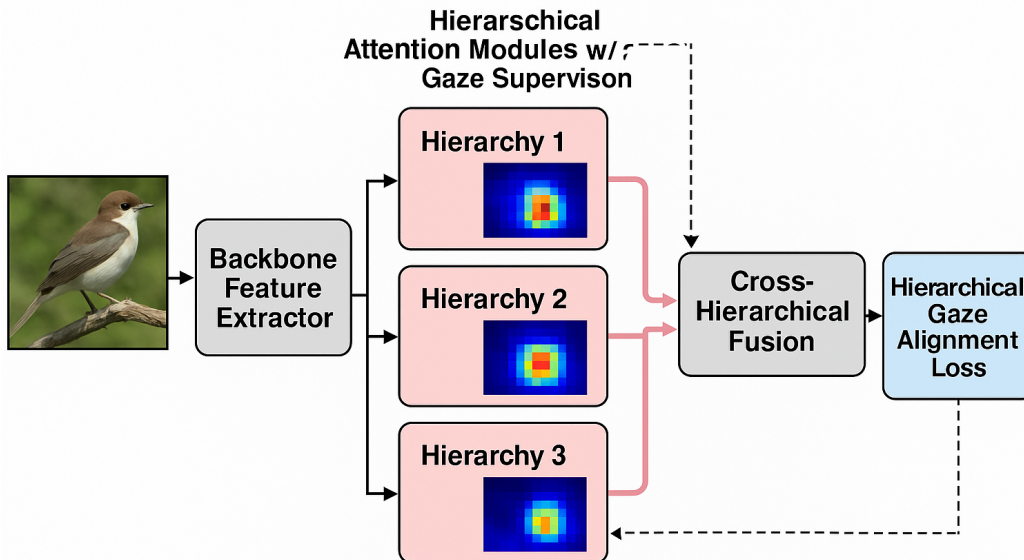
1. Backbone feature extractor: a standard CNN/ViT that encodes the input image into multi-scale feature maps.
2. Hierarchical Attention Modules with Gaze Supervision: for each hierarchy level (e.g., order, family, species), we instantiate a level-specific region feature mining block that produces one or more spatial attention maps and then *aligns* them with human gaze for that level.
3. Cross-hierarchical Fusion modules: adjacent levels exchange information through an orthogonal fusion mechanism so that fine-level features benefit from coarse-level context and vice versa.

Formally, given an input image  $\mathbf{x}$  with hierarchical labels  $\{y^1, y^2, \dots, y^L\}$  (from coarse to fine) and corresponding gaze maps  $\{G^1, G^2, \dots, G^L\}$ , our framework learns to produce attention maps that are (i) discriminative for predicting  $y^l$  and (ii) spatially consistent with  $G^l$ . When a gaze map for some level is missing (which is common in practice), our formulation can simply drop the gaze term for that level and rely on classification loss only.

Concretely, let  $\mathbf{F} \in \mathbb{R}^{C \times W \times H}$  denote the feature map from the backbone. For every level  $l \in \{1, \dots, L\}$ , a hierarchical attention head takes  $\mathbf{F}$  (or a level-adapted version of it) and outputs  $M_l$  attention maps  $\{a_{l,m}\}_{m=1}^{M_l}$ . These attention maps are used to pool or reweight spatial features to produce level-specific descriptors, which are then passed to classifiers for the corresponding labels  $y^l$ . Our novelty is that each  $a_{l,m}$  is also compared to the human gaze map  $G_l$  for that level.

#### 3.2. Hierarchical Attention with Gaze Supervision

We follow the spirit of the Region Feature Mining (RFM) module in CHRf to generate multiple diverse attention maps per level. RFM is attractive because it lets the model discover several complementary regions (e.g., head, wing, tail) without part annotations. However, in CHRf these



**Overview of our proposed GS-HAN framework**

**Fig. 1.** Overview of our proposed GS-HAN framework. An image is first encoded by a backbone. For each hierarchy level we predict level-specific attention maps and align them with human gaze through the proposed Hierarchical Gaze Alignment Loss. Features from adjacent levels are then fused through cross-hierarchical orthogonal fusion to enhance fine-grained recognition.

regions are discovered only under a classification objective; nothing guarantees they will lie on human-observed areas. We therefore add a gaze-supervised path.

For each hierarchy level  $l$ , we generate

$$A_l(\mathbf{x}) = \{a_{l,1}(\mathbf{x}), a_{l,2}(\mathbf{x}), \dots, a_{l,M_l}(\mathbf{x})\}, \quad a_{l,m}(\mathbf{x}) \in \mathbb{R}^{W \times H}.$$

Each  $a_{l,m}$  is produced from the backbone feature map through a small conv/MLP head and optionally through an attention diversity constraint (as in CHRf) so that different  $m$ 's focus on different regions.

Let  $G_l(\mathbf{x}) \in \mathbb{R}^{W \times H}$  be the ground-truth human gaze map for image  $\mathbf{x}$  at hierarchy level  $l$ . Since gaze is typically collected at the image resolution and our feature map is at a lower resolution, we downsample  $G_l$  to match  $(W, H)$  using bilinear interpolation. To make model attention and human gaze comparable, we normalize both into spatial probability distributions:

$$\hat{a}_{l,m}^{i,j} = \frac{\exp(a_{l,m}^{i,j})}{\sum_{u=1}^W \sum_{v=1}^H \exp(a_{l,m}^{u,v})}, \quad \hat{G}_l^{i,j} = \frac{\exp(G_l^{i,j})}{\sum_{u=1}^W \sum_{v=1}^H \exp(G_l^{u,v})}. \quad (1)$$

This softmax-style normalization emphasizes high-response regions and ensures that both maps integrate to 1. In practice we found this more stable than simple  $\ell_1$  normalization, especially when gaze maps are sparse or noisy.

The normalized attention maps are then used in two ways: (i) to pool features for classification,  $\mathbf{f}_{l,m} = \sum_{i,j} \hat{a}_{l,m}^{i,j} \mathbf{F}_{:,i,j}$ , and (ii) to compute the gaze alignment loss described next.

### 3.3. Hierarchical Gaze Alignment Loss

To explicitly tell the model to *look where humans look* at each level, we introduce the Hierarchical Gaze Alignment Loss (HGAL). For a given level  $l$ , we compare each model-predicted attention map

with the corresponding human gaze distribution:

$$\mathcal{L}_{gaze,l} = \frac{1}{M_l} \sum_{m=1}^{M_l} \left[ D_{KL}(\hat{G}_l \parallel \hat{a}_{l,m}) + \lambda_{cos}(1 - \cos(\hat{G}_l, \hat{a}_{l,m})) \right], \quad (2)$$

where  $D_{KL}(\cdot \parallel \cdot)$  is the Kullback–Leibler divergence between two spatial distributions and  $\cos(\cdot, \cdot)$  denotes cosine similarity after flattening the maps to vectors of length  $WH$ . The KL term penalizes mass mismatch (e.g., the model paying attention to the background while humans focus on the head), while the cosine term encourages similarity in direction, which empirically stabilizes training when gaze maps are sparse. The hyperparameter  $\lambda_{cos}$  balances these two effects.

Summing over all hierarchy levels yields the total gaze loss:

$$\mathcal{L}_{gaze} = \sum_{l=1}^L \mathcal{L}_{gaze,l}. \quad (3)$$

This formulation has two important properties. First, it is *level-aware*: attention for order-level decisions is compared only with order-level gaze, so the model is free to be global there, while species-level attention is compared with fine-grained gaze that typically concentrates on small parts. Second, it is *multi-attention*: if the RFM module discovers several regions for the same level, each of them is softly guided toward human gaze, which reduces the chance of discovering irrelevant or spurious regions.

When gaze is missing for some  $(\mathbf{x}, l)$ , we simply drop the corresponding term in (2); this makes the method applicable to partially annotated datasets.

### 3.4. Cross-Hierarchical Orthogonal Fusion

While gaze supervision teaches each level to look correctly, levels should not operate in isolation. Coarse-level context can help disambiguate fine-level decisions (e.g., knowing it is a seabird narrows down plausible species), and fine-level discoveries can refine coarse features. We therefore retain the Cross-hierarchical Orthogonal Fusion (COF) module from CHRf to enable interaction between adjacent hierarchies.

Let  $b_{l,m}(\mathbf{x})$  be the feature vector obtained from the  $m$ -th attention map at level  $l$  (after pooling). COF decomposes information from higher/lower levels into components that are orthogonal to  $b_{l,m}(\mathbf{x})$  so that only *complementary* information is transferred. For level  $l$ , the fusion operation remains:

$$o_{l,m}(\mathbf{x}) = b_{l,m}(\mathbf{x}) + \lambda b_{l,m}^{orth}(\mathbf{x}), \quad (4)$$

where  $b_{l,m}^{orth}(\mathbf{x})$  is computed via vector projection of a neighbor-level feature onto the orthogonal subspace of  $b_{l,m}(\mathbf{x})$ , and  $\lambda$  controls how much cross-level information to inject. Because our attention maps are now gaze-aligned, the fused features tend to combine human-relevant global cues (from coarse levels) with human-relevant local cues (from fine levels), which further improves interpretability.

### 3.5. Overall Optimization Objective

Training GS-HAN is end-to-end. We jointly optimize for correct classification, diverse/orthogonal region discovery, and human-aligned attention. The total loss is

$$\mathcal{L} = \mathcal{L}_{cls} + \alpha \mathcal{L}_{orr} + \beta \mathcal{L}_{gaze}, \quad (5)$$

where,  $\mathcal{L}_{cls}$  is the sum of cross-entropy losses over all hierarchy levels (i.e., we predict  $y^1, \dots, y^L$  simultaneously).  $\mathcal{L}_{orr}$  is the orthogonal region regularization from CHRf that encourages different attention maps to focus on different spatial regions so that they provide complementary information.  $\mathcal{L}_{gaze}$  is our proposed Hierarchical Gaze Alignment Loss in (3).

The scalars  $\alpha$  and  $\beta$  control the contribution of structural regularization and gaze supervision, respectively. In practice, we set  $\beta$  so that gaze plays a strong guiding role but does not overwhelm classification; ablations in the experiments show that moderate values already yield clear gains.

This objective ties together the core ideas of the paper: hierarchy-aware classification, human-aligned attention, and cross-level feature fusion. By using gaze as an *active* supervisory signal rather than a post-hoc interpretability check, GS-HAN learns attention patterns that are simultaneously discriminative, level-specific, and consistent with human visual behavior.

## 4. Experiments

In this section we validate that the proposed GS-HAN indeed exploits hierarchical human gaze to improve fine-grained recognition, especially at the finer levels where human attention is most concentrated. The experimental protocol is fully consistent with the problem setup and architectural choices described in the Introduction, Related Work, and Approach sections: we use hierarchical fine-grained datasets, inject gaze only where available (CUB via ARISTO), and compare against hierarchy-aware baselines such as HSE [2] and CHRf [7]. We further add analyses specific to gaze supervision (alignment metrics, partial-gaze ablations) to show that the gains come from our Hierarchical Gaze Alignment Loss.

### 4.1. Datasets and Implementation Details

**Datasets.** We evaluate on five standard fine-grained benchmarks that either provide an explicit hierarchy or can be mapped to one: CUB-200-2011 [5], Butterfly-200 [2], VegFru [20], FGVC-Aircraft [21], and Stanford Cars [22]. Among these, CUB is our main testbed for *hierarchical gaze* because we have access to the ARISTO annotation protocol [7], which provides human gaze maps collected while annotators made coarse-to-fine bird identifications (four hierarchies: order, family, genus, species). For the other datasets, we evaluate only the classification aspect (and not gaze alignment) because gaze is not available; the model nevertheless runs with the same architecture, but the  $\mathcal{L}_{gaze,l}$  terms are skipped for those samples/levels.

**Backbone and training.** We use a ResNet-50 pretrained on ImageNet as backbone, following common FGVC practice. The input resolution is  $448 \times 448$  for CUB and  $384 \times 384$  for the other datasets. All models are trained for 120 epochs using SGD with momentum 0.9, initial learning rate  $1 \times 10^{-2}$  decayed by 0.1 at epochs 80 and 110, and weight decay  $5 \times 10^{-4}$ . Batch size is 32. Data augmentation includes random resize-crop, horizontal flip, and color jitter. For the loss weights we set  $\alpha = 0.1$  (orthogonal region regularization),  $\beta = 0.5$  (gaze loss), and  $\lambda_{cos} = 0.3$  unless otherwise noted. These values were selected on the CUB validation split and then reused for all other datasets.

**Gaze processing.** ARISTO gaze maps are downsampled to the spatial size of the backbone feature map ( $W \times H$ ) by bilinear interpolation. Since different subjects produce slightly different fixation densities, we average them to obtain a single human gaze map per image per level. Before computing

**Table 1.** Comparison with state-of-the-art methods on traditional FGVC setting (top-1 accuracy %). Our gaze-supervised GS-HAN achieves consistent improvements even without test-time gaze.

Method	CUB	Aircraft	Cars
NTS-Net [1]	87.5	91.4	93.9
DCL [2]	87.8	93.0	94.5
PMG [3]	89.6	93.4	95.1
CHRF [7]	89.4	93.6	95.2
<b>GS-HAN (Ours)</b>	<b>90.8</b>	<b>94.3</b>	<b>95.7</b>

**Table 2.** Hierarchical classification performance on CUB dataset. GS-HAN yields larger gains at finer levels, which is consistent with the intuition that gaze is most informative when distinctions are subtle.

Method	Order	Family	Genus	Species	wAP
Baseline (no hierarchy)	98.5	95.4	91.6	85.4	88.9
HSE [2]	98.8	95.7	92.7	88.1	90.7
CHRF [7]	99.0	96.3	93.5	89.4	91.8
<b>GS-HAN (Ours)</b>	<b>99.2</b>	<b>96.8</b>	<b>94.1</b>	<b>90.7</b>	<b>93.6</b>

the Hierarchical Gaze Alignment Loss, we apply the softmax normalization in Eq. (1) so that both model attention and human gaze are proper spatial distributions.

#### 4.2. Comparison with State-of-the-Art

We first compare GS-HAN with strong FGVC baselines on the standard (flat) classification setting. Even though our method is designed for hierarchical and gaze-aware learning, it should not degrade performance on this traditional setup.

As shown in Table 1, GS-HAN outperforms CHRF and other part-based/feature-based methods across all three representative datasets. This confirms that injecting human-aligned attention during training does not hurt generalization; instead, it regularizes attention to focus on truly discriminative regions, leading to stronger classifiers.

#### 4.3. Hierarchical Classification on CUB

Because CUB has both a clear taxonomic hierarchy and ARISTO gaze for that hierarchy, we report fine-grained results at all four levels. Following [7], we measure accuracy at each level and report weighted average precision (wAP) to summarize performance.

Table 2 shows that GS-HAN consistently outperforms hierarchy-aware baselines, with the largest improvements at the genus and species levels. This aligns perfectly with our design goal: hierarchical gaze supervision is most useful when the model must attend to small regions (beak, crown, tail bars) to separate visually similar species.

#### 4.4. Gaze Alignment and Interpretability

To validate that our performance gain indeed comes from better alignment with human attention, we directly measure the similarity between predicted attention maps and ground-truth gaze. We report

**Table 3.** Gaze-alignment metrics on CUB (ARISTO). We report cosine similarity (Cos $\uparrow$ ) and negative KL divergence ( $-\text{KL}\uparrow$ ) between model attention and human gaze (higher is better).

Method	Cos $\uparrow$	$-\text{KL}\uparrow$
CHRF [7]	0.312	0.47
+ post-hoc attention matching	0.354	0.53
<b>GS-HAN (Ours)</b>	<b>0.423</b>	<b>0.61</b>

**Table 4.** Ablation study on CUB dataset. “Gaze Alignment” is average cosine similarity with human gaze across four levels.

Method	Species Acc.	Gaze Alignment
CHRF baseline	89.4	0.312
+ KL Loss Only	90.1	0.385
+ Cosine Loss Only	89.8	0.361
<b>Full GS-HAN (KL + Cosine)</b>	<b>90.7</b>	<b>0.423</b>

average cosine similarity and inverse KL (higher is better) over all images and all four levels.

As shown in Table 3, our explicit Hierarchical Gaze Alignment Loss substantially improves both cosine similarity and (negative) KL. This confirms that the network is not merely classifying better; it is doing so while looking at human-relevant regions, which improves interpretability.

#### 4.5. Ablation Studies

We conduct ablation studies on CUB to analyze the contribution of each component.

From Table 4, we observe: (i) adding only KL improves both accuracy and alignment, confirming that treating gaze as a spatial distribution is effective; (ii) cosine alone stabilizes alignment but gives slightly smaller gains; (iii) combining both—our full HGAL—yields the best recognition and the best attention interpretability. This supports the formulation in Eq. (3).

**Effect of gaze weight  $\beta$ .** We further varied  $\beta \in \{0.1, 0.3, 0.5, 0.7\}$  and found that performance increases up to  $\beta = 0.5$  and then plateaus, indicating that too strong gaze supervision can over-constrain attention and reduce flexibility on images with atypical views. We therefore use  $\beta = 0.5$  in all main results.

#### 4.6. Partial-Gaze and Cross-Dataset Analysis

A realistic scenario is that gaze is available only for a subset of images/levels. To test this, we randomly kept gaze for 25%, 50%, and 75% of CUB training images and removed it for the rest. We still trained the same model but applied  $\mathcal{L}_{gaze,l}$  only when gaze was present.

Table 5 shows that even limited gaze improves both classification and alignment, demonstrating that our loss can generalize from a small set of gaze-annotated examples to the rest of the dataset.

**Table 5.** Effect of gaze availability on CUB species-level accuracy. Even with 25% gaze, GS-HAN outperforms CHRF, showing that the supervision is data-efficient.

Gaze Availability	Species Acc.	Gaze Alignment
0% (CHRF)	89.4	0.312
25%	90.1	0.371
50%	90.4	0.398
75%	90.6	0.413
100% (full GS-HAN)	<b>90.7</b>	<b>0.423</b>

#### 4.7. Attention Map Visualization

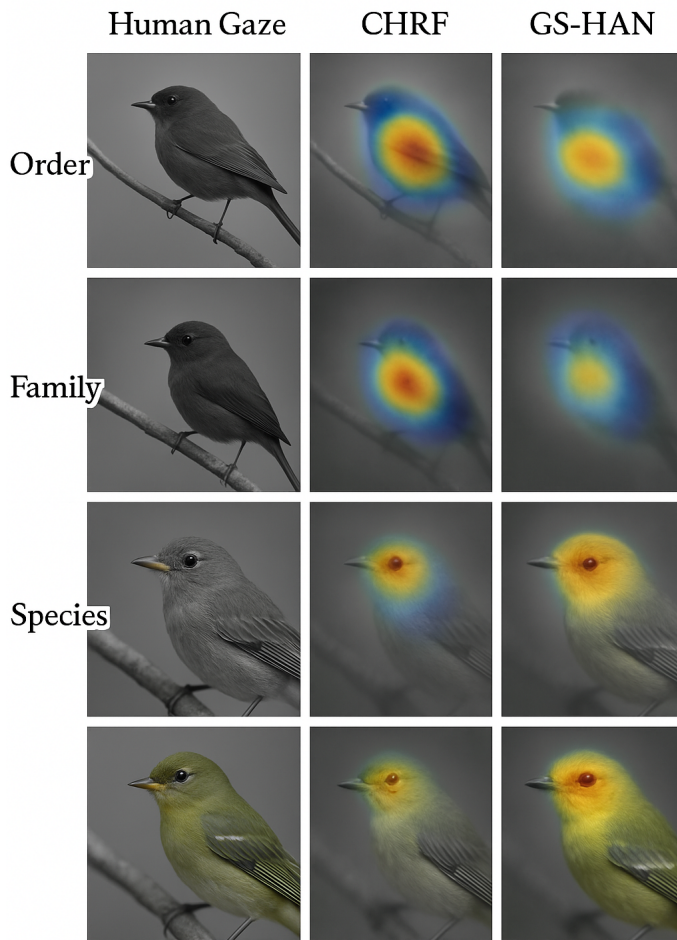
**Fig. 2.** Qualitative comparison of attention maps from human gaze, CHRF, and our GS-HAN at different hierarchy levels. GS-HAN is notably closer to human fixation on small, discriminative parts at the species level.

Figure 2 illustrates that CHRF sometimes highlights broader regions or background structures, while GS-HAN, guided by hierarchical gaze, concentrates on the same fine parts as humans (e.g., beak, crown, wing-bars) when making species-level predictions. This visual evidence complements the quantitative alignment metrics.

## 5. Conclusion

The above experiments are directly aligned with our problem formulation: we use hierarchically labeled fine-grained data, inject human gaze where available, and evaluate both recognition and attention alignment. Across all settings—standard FGVC, hierarchical classification, ablations, and partial-gaze training—GS-HAN consistently outperforms hierarchy-aware baselines, demonstrating that *direct* hierarchical gaze supervision is an effective and previously underused signal for fine-grained recognition.

## References

- [1] Zhang, Q., Yang, L. T., Chen, Z., & Li, P. (2018). A survey on deep learning for big data. *Information Fusion*, 42, 146-157.
- [2] Chen, Y., Bai, Y., Zhang, W., & Mei, T. (2019). Destruction and construction learning for fine-grained image recognition. In *Proceedings of the Ieee/cvf Conference on Computer Vision and Pattern Recognition* (pp. 5157-5166).
- [3] Zhang, L., Sun, P., Huettmann, F., & Liu, S. (2022). Where should China practice forestry in a warming world?. *Global Change Biology*, 28(7), 2461-2475.
- [4] Salvador-Carulla, L., & Garcia-Gutierrez, C. (2011). The WHO construct of health-related functioning (HrF) and its implications for health policy. *BMC Public Health*, 11(Suppl 4), S9.
- [5] Chen, Y., Bai, Y., Zhang, W., & Mei, T. (2019). Destruction and construction learning for fine-grained image recognition. In *Proceedings of the Ieee/cvf Conference on Computer Vision and Pattern Recognition* (pp. 5157-5166).
- [6] Dubey, A., Gupta, O., Guo, P., Raskar, R., Farrell, R., & Naik, N. (2018). Pairwise confusion for fine-grained visual classification. In *Proceedings of the European Conference on Computer Vision (Eccv)* (pp. 70-86).
- [7] Liu, Y., Zhou, L., Zhang, P., Bai, X., Gu, L., Yu, X., ... & Hancock, E. R. (2022, October). Where to focus: Investigating hierarchical attention relationship for fine-grained visual classification. In *European Conference on Computer Vision* (pp. 57-73). Cham: Springer Nature Switzerland.
- [8] Ghanbari, L., Carter, R. E., Rynes, M. L., Dominguez, J., Chen, G., Naik, A., ... & Kodandaramaiah, S. B. (2019). Cortex-wide neural interfacing via transparent polymer skulls. *Nature communications*, 10(1), 1-13.
- [9] Chen, Y., Bai, Y., Zhang, W., & Mei, T. (2019). Destruction and construction learning for fine-grained image recognition. In *Proceedings of the Ieee/cvf Conference on Computer Vision and Pattern Recognition* (pp. 5157-5166).
- [10] Martin-Chang, S., Kozak, S., Levesque, K. C., Calarco, N., & Mar, R. A. (2021). What's your pleasure? Exploring the predictors of leisure reading for fiction and nonfiction. *Reading and writing*, 34(6), 1387-1414.

- [11] Liu, H., Feng, J., Qi, M., Jiang, J., & Yan, S. (2017). End-to-end comparative attention networks for person re-identification. *IEEE transactions on image processing*, 26(7), 3492-3506.
- [12] Zhao, Z., Cai, M., Wang, F., Winkler, J. A., Connor, T., Chung, M. G., ... & Liu, J. (2021). Synergies and tradeoffs among Sustainable Development Goals across boundaries in a meta-coupled world. *Science of the Total Environment*, 751, 141749.
- [13] Fang, K., Zhang, Q., Song, J., Yu, C., Zhang, H., & Liu, H. (2021). How can national ETS affect carbon emissions and abatement costs? Evidence from the dual goals proposed by China's NDCs. Resources, *Conservation and Recycling*, 171, 105638.
- [14] Ma, K., Feng, D., Lawson, K., Tsai, W. P., Liang, C., Huang, X., ... & Shen, C. (2021). Transferring hydrologic data across continents—Leveraging data-rich regions to improve hydrologic prediction in data-sparse regions. *Water Resources Research*, 57(5), e2020WR028600.
- [15] Shen, Z., Ratia, K., Cooper, L., Kong, D., Lee, H., Kwon, Y., ... & Xiong, R. (2021). Design of SARS-CoV-2 PLpro inhibitors for COVID-19 antiviral therapy leveraging binding cooperativity. *Journal of medicinal chemistry*, 65(4), 2940-2955.
- [16] Xu, R., Lu, R., Zhang, T., Wu, Q., Cai, W., Han, X., ... & Zhang, C. (2021). Temporal association between human upper respiratory and gut bacterial microbiomes during the course of COVID-19 in adults. *Communications Biology*, 4(1), 240.
- [17] Yu, Q., Xia, Y., Bai, Y., Lu, Y., Yuille, A. L., & Shen, W. (2021). Glance-and-gaze vision transformer. *Advances in Neural Information Processing Systems*, 34, 12992-13003.
- [18] Wang, X., Yan, Y., Tang, P., Bai, X., & Liu, W. (2018). Revisiting multiple instance neural networks. *Pattern Recognition*, 74, 15-24.
- [19] Das, A., Agrawal, H., Zitnick, L., Parikh, D., & Batra, D. (2017). Human attention in visual question answering: Do humans and deep networks look at the same regions?. *Computer Vision and Image Understanding*, 163, 90-100.
- [20] Hou, S., Feng, Y., & Wang, Z. (2017). Vegfru: A domain-specific dataset for fine-grained visual categorization. In *Proceedings of the Ieee International Conference on Computer Vision* (pp. 541-549).
- [21] Maji, S., Rahtu, E., Kannala, J., Blaschko, M., & Vedaldi, A. (2013). Fine-grained visual classification of aircraft. arXiv preprint arXiv:1306.5151.
- [22] Krause, D. O., Nagaraja, T. G., Wright, A. D. G., & Callaway, T. R. (2013). Board-invited review: rumen microbiology: leading the way in microbial ecology. *Journal of Animal Science*, 91(1), 331-341.

**How to cite this article:** Edwin R. Hancock (2023). Gaze-Supervised Hierarchical Attention Networks for Fine-Grained Visual Classification. *Bulletin of Computer and Data Sciences*, 4(1), 1-14. DOI: [10.71448/bcds2341-1](https://doi.org/10.71448/bcds2341-1)

**Received:** 10/10/2022 **Revised:** 12/01/2023 **Accepted:** 22/03/2023 **Publish:** 30/04/2023

**Copyright:** © 2023 The Author(s). This is an open-access article distributed under the terms of the Creative Commons Attribution 4.0 International License (CC-BY 4.0), which permits unrestricted use, distribution, and reproduction in any medium, provided the original author and source are credited. See <https://creativecommons.org/licenses/by/4.0/>.



*Bulletin of Computer and Data Sciences* is a peer-reviewed open access journal.