

# From Generalized zero-shot learning to long-tail with class descriptors

Dvir Samuel<sup>1</sup> Yuval Atzmon<sup>2</sup> Gal Chechik<sup>1,2</sup>

<sup>1</sup>Bar-Ilan University, Ramat Gan, Israel

<sup>2</sup>NVIDIA Research, Tel Aviv, Israel

dvirsamuel@gmail.com, yatzmon@nvidia.com, gal.chechik@biu.ac.il

## Abstract

Real-world data is predominantly unbalanced and long-tailed, but deep models struggle to recognize rare classes in the presence of frequent classes. Often, classes can be accompanied by side information like textual descriptions, but it is not fully clear how to use them for learning with unbalanced long-tail data. Such descriptions have been mostly used in (Generalized) Zero-shot learning (ZSL), suggesting that ZSL with class descriptions may also be useful for long-tail distributions.

We describe DRAGON, a late-fusion architecture for long-tail learning with class descriptors. It learns to (1) correct the bias towards head classes on a sample-by-sample basis; and (2) fuse information from class-descriptors to improve the tail-class accuracy. We also introduce new benchmarks CUB-LT, SUN-LT, AWA-LT for long-tail learning with class-descriptors, building on existing learning-with-attributes datasets and a version of Imagenet-LT with class descriptors. DRAGON outperforms state-of-the-art models on the new benchmark. It is also a new SoTA on existing benchmarks for GFSL with class descriptors (GFSL-d) and standard (vision-only) long-tailed learning ImageNet-LT, CIFAR-10, 100, and Places365-LT.

## 1. Introduction

Real-world data is predominantly unbalanced, typically following a long-tail distribution. From text data (Zipf’s law), through acoustic noise (the 1-over- $f$  rule) to the long-tail distribution of classes in object recognition [45], few classes are frequently observed, while the many remaining ones are rarely encountered.

Long-tail data poses two major challenges to learning: *data paucity* and *data imbalance*. First, at the tail of the distribution, classes are poorly sampled and one has to use few-shot and zero-shot learning techniques. Second, when training a single model for both richly-sampled classes and poorly-sampled classes, the common classes dominate

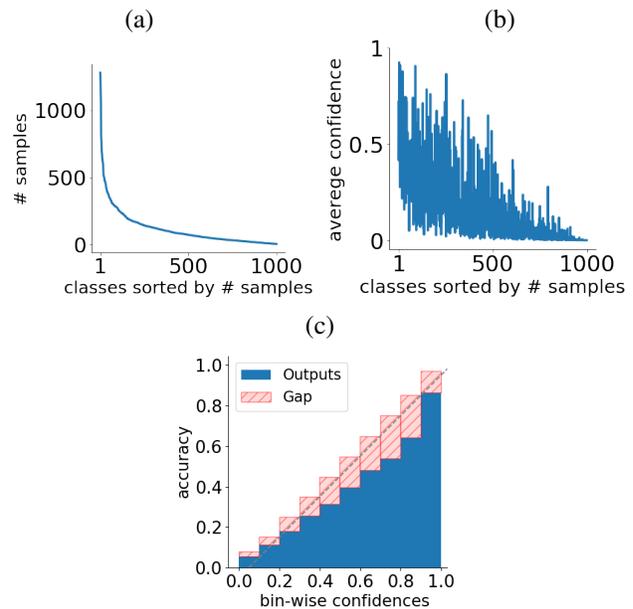


Figure 1: Training with unbalanced data leads to a “familiarity bias”, where models are more confident and more over-confident about frequent classes [48, 47, 6]. (a) Class distribution of a long-tailed ImageNet [12]. Classes are ordered from left to right by decreasing number of samples. (b) When training a ResNet-10 on ImageNet-LT, validation (and test) predictions tend to have low confidence for tail classes. We show the mean output of softmax for each class, conditioned on samples from that class. (c) A reliability graph for the model in b. Predictions are grouped based on confidence. The model has larger confidence gaps (pink boxes) for more confident predictions, which usually come from head classes. This result suggests that overconfidence is strongly affected by class frequency, and we can learn to correct it if the number of samples is known.

training, and as we show below, this skews prediction confidence towards rich-sampled classes.

To address *data paucity* of tail classes, note that visual

examples can very often be augmented with class descriptors. Namely, semantic information about classes given as text or attributes [26, 36, 4, 52]. This approach, *learning with class-descriptors* has been studied mostly for zero-shot and generalized zero-shot learning [38, 32, 43, 52, 53]. Here, we propose to adapt it to generalized few-shot learning (GFSL) by fusing information from two modalities. A visual classifier, expected to classify correctly head classes, and a semantic classifier, trained with per-class descriptors and is expected to classify correctly tail classes. We explain the subtleties of GZSL and GFSL in Section 2 (Related work).

To address *data imbalance*, we first note that learning with unbalanced data leads to a **familiarity effect**, where models become biased to favor the more familiar, rich-sampled classes [48, 6]. Since deep models tend to be overly confident about high-confidence sample predictions [17, 25], they become over-confident about head classes [48, 47, 6]. Figure 1 illustrates this phenomenon. It shows that a model trained on unbalanced data (Figure 1a), has higher confidence for head classes (Figure 1b). It is also an over-estimate of the true accuracy (Figure 1c), especially for head classes. See further analysis in Appendix A.

As an interesting side note, studies of human decision making and preference learning show a similar bias towards familiar classes. This effect is widely observed and has been connected to the availability heuristic studied by Tversky and Kahneman [44].

A natural way to correct the familiarity bias would be to penalize high-frequency classes, either during training using a balanced loss, or post-training [23]. However, it would be a grave mistake to penalize all samples from rich classes, because confidence is sometimes justified, as in the case of "easy" prototypical examples of a class. Indeed, we show below that addressing the familiarity bias benefits from *per-sample debiasing*, going beyond class-based debiasing. To summarize, **model overconfidence is affected by class frequency. It can be estimated by observing the full vector of predictions to correct for overconfidence. Not all samples of a class should be penalized for belonging to a frequent class.**

Importantly, the familiarity effect caused by data imbalance has a crippling effect on model accuracy and on aggregating predictions from multiple modalities. Several approaches attempted calibrating predictions of deep networks to remedy the above biases (see e.g. a survey in [17]) and some became common practice. Unfortunately, the problem is still far from being solved.

We propose to address both the *data-imbalance* and the *data paucity* learning challenges, using a single late-fusion architecture. We describe an easy-to-implement debiasing module that offsets the familiarity effect by learning to predict the magnitude of the bias for any given sample. It fur-

ther improves learning at the tail by learning to fuse and balance information from visual and semantic modalities. It can easily be reduced to address long-tail learning with a single modality (vision only), where it improves over current baselines.

The paper has four main novel contributions:

- (1) **A new late-fusion architecture** (DRAGON) that learns to fuse predictions based on vision with predictions based class descriptors.
- (2) **A module that rebalances class predictions across classes on a sample-by-sample basis.**
- (3) **New benchmarks** CUB-LT, SUN-LT, and AWA-LT for evaluating long-tail learning with textual class descriptors (LT-d). DRAGON is SoTA on these datasets, and also on ImageNet-LT augmented with class descriptors and on existing two-level benchmarks.
- (4) **A new SoTA** on existing (vision-only) long-tail learning benchmarks: CIFAR-10, CIFAR-100, ImageNet-LT and Places365-LT.

## 2. Related methods

**Long-tail learning:** Learning with unbalanced data causes models to favor head classes [6]. Previous efforts to address this effect can be viewed as either algorithmic or data-manipulations approaches.

Algorithmic approaches encourage learning of tail classes using a non-uniform cost per misclassification function. A natural approach is to rescale the loss based on class frequency [19]. [27] proposed to down-weight the loss of well-classified examples, preventing easy negatives from dominating the loss. [37] dynamically rescaled the cross-entropy loss based on the difficulty to classify a sample. [7] proposed a loss that encourages larger margins for rare classes. [23] decoupled the learning procedure into representation learning and classification and studied four approaches. Among them, LWS  $L_2$ -normalizes the last-layer, since the weight magnitude correlates with class cardinality. The effect of this approach is similar to that presented in this paper, but here we apply recalibration dynamically on a sample-by-sample basis.

Data-manipulation approaches aim to flatten long-tail datasets to correct the bias towards majority classes. Popular techniques employ over-sampling of minority classes (more likely to overfit)[10, 18], under-sampling the majority classes (wastes samples)[13], or generating samples from the minority classes (can be costly to develop)[5].

Another approach is to transfer meta-level-knowledge from data-rich classes to data-poor classes. [49] gradually transfer hyperparameters from rich classes to poor classes by representing knowledge as trajectories in model space that capture the evolution of parameters with increasing training samples. [29] first learns representations on the

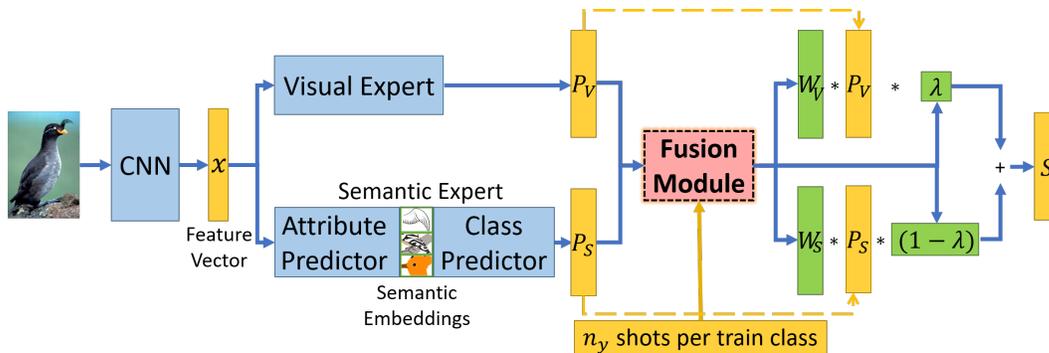


Figure 2: The DRAGON architecture for long-tail learning with class-descriptors. The visual-expert and attribute-expert each outputs a prediction vector fed to a fusion module. The fusion module combines expert predictions and debias them. Blue, network components. Yellow, input to the fusion module. Green, the outputs of the fusion module.

unbalanced data and then fine-tunes them using a class-balanced sampling and a memory module.

**Learning with class descriptors:** Learning with class descriptors is usually applied to zero-shot learning (ZSL) [52, 26, 4], where a classifier is trained to recognize (new) unseen classes based on their semantic description, which can include a natural-language textual description or predefined attributes. In several ZSL studies, attributes detected in a test image are matched with ground-truth attributes of each class, and several studies focused on this matching [26, 4, 42, 56, 55, 8].

A series of papers proposed to learn a shared representation of visual and text features (class-descriptors). As one example, [43] learns such a shared latent manifold using autoencoders and then minimizes the MMD loss between the two domains. Another recent line of work synthesizes feature vectors of unseen classes using generative models like VAE or GAN, and then use them in training a conventional classifier [32, 51, 14, 1, 31, 59, 38, 53]. The major baseline we compare our approach with is CADA-VAE [38], the current SoTA for Generalized FSL with class descriptors. CADA-VAE uses a variational autoencoder that aligns the distributions of image features and semantic (attribute) class embedding in a shared latent space. A recent work, [53], uses a mixture of VAEs and GANs. We could not directly compare with [53] because their FSL protocol deviates from the standard benchmark of [38, 52] by fine-tuning the CNN features. *Without* fine-tuning, their reported metrics for GZSL are similar to CADA-VAE.

Some studies fused information from vision and *per sample* descriptors (e.g., [58]). This is outside the scope of this paper because it may require extensive labeling.

**Generalized ZSL (GZSL) and Generalized FSL:** GZSL extends ZSL to the scenario where the test data contains both seen and unseen classes [9, 52, 40]. Recently, GZSL extended to Generalized Few-Shot-Learning

with class descriptors (GFSL-d), where the unseen classes are augmented with a fixed number of few training samples [38, 43]. Namely, the distribution of samples across classes is a 2-level distribution, with *many* “head” classes and a smaller set of “tail” classes all having the same (small) number of samples per class. Both GZSL and GFSL-d can be viewed as special cases of long-tail learning with class-descriptors, but with a *short-tailed* unnatural distribution.

Most related GZSL approaches are [3, 40, 54]. They use a gating mechanism to weigh the decisions of seen-classes experts and a ZSL expert. The gating module is modeled as an out-of-distribution estimator. The current paper differs from their work by (1) The problem setup is different. Here, all samples are *in-distribution* and the distribution of classes is smooth and long-tail with a much smaller number of head classes. (2) DRAGON architecture first quantifies and corrects the (smooth) familiarity effect. Then it learns how to fuse the debiased decision of the two experts.

**Early vs late fusion:** When learning from multiple modalities, one often distinguishes between early and late fusion models [28]. Early fusion models combine features from multiple modalities to form a joint representation. Late fusion methods combine decisions of per-modality models [22, 2, 35]. Our approach addresses the long-tail setup, by leveraging the information in the familiarity bias to debias experts predictions.

### 3. Long-tail learning with class descriptors

We start with a formal definition of the problem of learning over unbalanced distributions with class descriptors.

We are given a training set of  $n$  labeled (image) samples:  $\{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)\}$ , where each  $\mathbf{x}_i$  is a feature vector and  $y_i$  is a label in  $\{1, 2, \dots, k\}$ . Samples are drawn from a distribution  $\mathcal{D} = p(x, y)$  such that the marginal distribution over the classes  $p(y)$  is strongly non uniform. For example,  $p(y)$  may be exponential  $p(y) \sim \exp(-ky)$ .

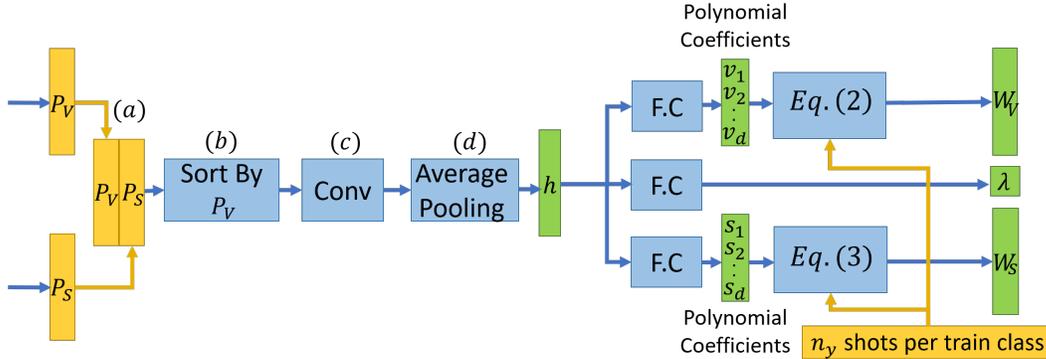


Figure 3: Architecture of the fusion-module for long-tail learning with class-descriptors. In blue, network components. In yellow, inputs to the fusion-module and in green, activations or outputs of the fusion-module. The inputs  $P_V$  denote the softmax prediction vector of the *Visual Expert*, and  $P_S$  that of the *Semantic Expert*. The outputs  $W_V$ ,  $W_S$  and  $\lambda$  are used in Eq. (1) for re-weighting the inputs. See Section 4.2 for more details.

As a second supervision signal, each class  $y$  is also accompanied with a *class-description* vector  $\mathbf{a}_j$ ,  $j = 1, \dots, k$ , in the form of semantic attributes [26] or natural-language embedding [36, 59, 40]. For example, classes in CUB [46] are annotated with attributes like `Head-color:red`.

At test time, a new set of  $m$  test samples  $\{\mathbf{x}_{n+1}, \dots, \mathbf{x}_{n+m}\}$  is given from the same distribution  $\mathcal{D}$ . We wish to predict their correct classes.

## 4. Our approach

Our approach is based on two observations: (1) Semantic descriptions of classes are easy to collect and can be very useful for tail (low-shot) classes, because they allow models to recognize classes even with few or no training samples [26, 52, 4] (and Appendix B). (2) The average prediction confidence over samples of a class is correlated with the number of training samples of that class (Figure 1).

Our architecture leverages these observations and learns to (1) Combine predictions of two expert classifiers: A conventional visual expert which is more accurate at head classes and a semantic expert which excels at tail classes; (2) Reweight the scores of each class prediction, taking into account the number of training samples for that class.

### 4.1. The architecture

The DRAGON architecture<sup>1</sup> follows two design considerations: *modularity* and *low-complexity*. First, modularity; DRAGON allows to plug-in existing multi-modal experts, each trained for its own modality. Below we show experiments with language-based experts and a visual expert, but

<sup>1</sup>Dragons, like many distributions, have long-tails and are cool. For acronym lovers, DRAGON also stands for “a moDulaR Approach for long-tail classificatiON”.

other modalities can be considered (e.g., mesh, depth, motion, or multi-spectral information). Second, limiting the model to have a small number of parameters is important because tail classes only have few training samples and the model must perform well at the tail.

Our general architecture (Figure 2) takes a late fusion approach. It consists of two experts modules: A visual expert and a semantic expert. Each expert outputs a prediction vector which is fed to a fusion module. The fusion module combines the expert predictions and learns to debias the familiarity effect, by weighing the experts and re-scaling their class predictions.

### 4.2. A fusion-module

The fusion module takes as input the prediction vectors of two experts  $p_V$ ,  $p_S$  for a given image, and a vector containing the number of training samples per class. It has three outputs:  $\lambda \in (0, 1)$  is a scalar to trade-off the visual expert against the semantic expert.  $\mathbf{w}_V \in (0, 1)^k$  is a vector that weighs the predictions of the visual expert. Similarly,  $\mathbf{w}_S$  weighs the semantic expert predictions. Given these three outputs, a debiased score is computed for a class  $y$ :

$$S(y) = \lambda \mathbf{w}_V(y) p_V(y) + (1 - \lambda) \mathbf{w}_S(y) p_S(y). \quad (1)$$

Figure 3 describes the architecture of the fusion-module. It has two main parts. The first part maps the prediction scores to a meaningful joint space, by first aligning the prediction of both classifiers, and then sorting according to confidence.

In more detail, the first part has four steps. (a) Stacking together the predictions of two experts to a  $\mathcal{Y} \times 2$  vector. This makes the following convolution meaningful across the 2 experts axis. (b) To make convolution meaningful also along the classes axis, and since classes are categorical, we reorder classes by their prediction score according

to one of the experts. Section 4.3 explains the rationale of this reordering. (c) Now that predictions are sorted, feed the sorted scores to a  $N_{filters} \times 2 \times 2$  convolutional network. (d) Follow with an average-pooling layer (per class), yielding a  $(\mathcal{Y} - 1)$ -dimensional vector  $\mathbf{h}$ .

The goal of the second part is simply to predict a debiasing coefficient for each prediction, namely, learn a function from  $n_y$  to  $(0, 1)^k$ . We know from Figure 1 that the bias is inversely related to the number of samples. Aiming for a simple model, we train a polynomial regression that takes as input the number of samples and outputs a debiasing weight ( $w_V(y)$ ). The coefficients of this polynomial  $v_0, \dots, v_{d-1}$  are learned as a deep function over  $\mathbf{h}$ . Similarly for  $f_S(y)$ .

More formally, let  $n_y$  be the number of training samples of class  $y$ , and let  $\bar{n}_y = n_y / \max_y n_y$  be the normalized counts, then we have

$$\mathbf{w}_V(y) = \sigma \left[ \sum_{j=0}^{d-1} v_j(\mathbf{h}) \bar{n}_y^j \right], \quad \mathbf{w}_S(y) = \sigma \left[ \sum_{j=0}^{d-1} s_j(\mathbf{h}) \bar{n}_y^j \right], \quad (2)$$

where  $d$  is the polynomial degree and  $\sigma$  denotes a sigmoid that ensures that the resulting scale is in  $[0, 1]$ .

Finally, the fusion module also predicts the trade-off scalar  $\lambda$  to control the relative weight of the visual and semantic experts. This is achieved using a fully connected layer over  $\mathbf{h}$ . Section 9 analyzes the contribution of each component of the approach with an ablation study.

### 4.3. Architecture design decision

DRAGON is designed with a small number of model parameters so it can improve predictions at the tail, where very few samples are available.

**Debiasing based on all expert predictions.** To achieve the above goals, DRAGON implicitly learns how frequent is a class of a given sample. Namely, when the model receives a new test sample, it predicts if it is from a head class, a tail class, or somewhere between, and adjusts the confidence of the experts accordingly. To do this, it has been shown in the context of zero-shot learning that profile of confidence values are a good predictor if a sample comes from a seen or unseen class [3]. DRAGON generalizes this idea to long-tail distributions. To do this it takes as input all class predictions from both experts (Figure 3a).

**Using order statistics over predictions.** To process expert prediction, we point out that order statistics over the prediction vector – the maximum confidence,  $2^{nd}$  max, etc. – provides a strong signal about confidence calibration. Using the maximum of a vector is a very common operator in deep learning, known as max pooling. Here however, there is additional important information the gap between subsequent order statistics, like  $\max - 2^{nd} \max$ . As an intuitive example, a maximal prediction of 0.6 should be interpreted differently if the  $2^{nd}$  max is 0.4 or 0.1.

Order statistics can be easily computed by sorting the vector of predictions (Figure 3b). Sorting also increases the sample efficiency for learning, because later layers have each order statistic located at a fixed position in their input regardless of class. The function learned over order-statistics gaps is therefore shared across all classes.

The  $2 \times 2$  convolution (Figure 3c) works well with the sorted expert predictions. Its filters capture two signals: (1) the confidence gaps between the two experts for each class; and (2) the confidence gaps between order-statistics for each expert alone.

## 5. Experiments

We evaluate DRAGON in three unbalanced benchmark scenarios. (1) “*Smooth-Tail*”, the long-tailed distribution of classes decays smoothly (Figure 4) and each class is accompanied with textual class descriptors. (2) “*Two-Level*”, the distribution has a step-shape as in [38]; Most classes have many samples and the rest have few samples (Figure 5). (3) “*Vision-only*”, a long-tail setup, as in Smooth-Tail, but without class descriptors.

We compare DRAGON with SoTA approaches on standard benchmarks for each of these three scenarios. See Appendix D for implementation details.

Code available at <https://github.com/dvirsamuel/DRAGON>

### 5.1. Overview of main results

**For Smooth-Tail distributed data,** we evaluated DRAGON on four benchmarks that we created from existing datasets. We added textual descriptors to ImageNet-LT [29], and generated long-tail versions of CUB, SUN and AWA. Dragon outperforms all baselines on various metrics.

**For Two-Level distributed data,** DRAGON surpasses the current SoTA [38], tested using their experimental setup.

**For Vision-only long-tail data,** we tested the calibration component of DRAGON (without fusion). It achieves a new SoTA on ImageNet-LT [29], Places365-LT [29], Unbalanced CIFAR-10/100 [7] and comparable results on iNaturalist2018 [20].

## 6. Smooth-Tail distribution

### 6.1. Datasets

To evaluate long-tail learning with class descriptors, we created benchmark datasets in two ways. First, we created long-tail versions of existing learning-with-class-descriptors benchmarks. Second, we augmented existing long-tail benchmark (ImageNet-LT) with class descriptors.

Specifically, we created new long-tail variants of the 3 main learning-with-class-descriptors benchmarks: CUB [46], SUN [33] and AWA [26], illustrated in Figure 4. We ranked classes by the number of samples in each class

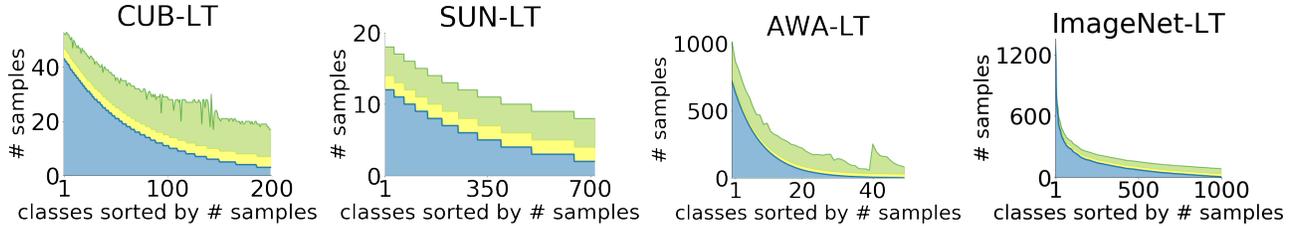


Figure 4: Long-tailed versions of CUB, SUN, AWA and ImageNet. Number of samples for the training, validation and test sets are shown respectively by blue, yellow and green.

	CUB-LT	SUN-LT	AWA-LT
Train # samples	2,945	4,084	6,713
Val # samples	600	1,434	1,250
Test # samples	2,348	2,868	6,092
<b>Train Set Properties</b>			
Max # samples	43	12	720
Min # samples	3	2	2
Mean # samples	14.725	5.696	123.460
Median # samples	11	5	35

Table 1: Properties of CUB-LT, SUN-LT and AWA-LT.

after assigning tail classes to be consistent with those in the Two-Level benchmark [51, 38] (See Appendix D.2 for more details). We then computed a frequency level for each class following an exponentially decaying function of the form  $f(class) = ab^{-rank(class)}$ .  $a$  and  $b$  were selected such that the first class has the maximum number of samples, and the last class has 2 or 3 samples depending on the dataset. We then drew a random subset of samples from each class based on their assigned frequency  $f(class)$ . To create the validation set, we randomly drew a constant number of samples per class, while keeping an overall size of 20% of the training set. See dataset statistics in Table 1.

As a second type of benchmark, we used the existing long-tailed ImageNet [29] and augmented it with class descriptors. Specifically, we used the word2vec embeddings provided by [8], which are widely used in the literature. Their word embeddings were created by training a skip-gram language model on a Wikipedia corpus to extract a 500-dimensional word vector for each class. See [8] Figure 4 for visualization.

Together, this process yielded the following datasets:

1. **CUB-LT**, based on [46], consists of 2,945 training visual images of 200 bird species. Each species is described by 312 attributes (like tail-pattern:solid, wing-color:black). Classes have between 43 and 3 images per class.
2. **SUN-LT**, based on [33], consists of 4,084 training images, from 717 visual scene types and 102 attributes (like material:rock, function:eating, surface:glossy). Classes

have between 12 and 2 images per class.

3. **AWA-LT**, based on [26], consists of 6,713 training images of 50 animal classes and 85 attributes (like texture:furry, or color:black). Classes have between 720 and 2 images per class.

4. **ImageNet-LT-d**, based on [29], consists of 115.8K images from 1000 categories with 1280 to 5 images per class. We use Word2Vec [30] class embeddings features provided by [8], as textual descriptors.

## 6.2. Training scheme

The familiarity effect is substantial in the validation and test data, but not in the training data, where models may actually become more confident on rare classes. We observed this effect in CUB-LT, SUN-LT, and AWA-LT. Since we wish to train the fusion module using data that exhibits the familiarity bias, we hold-out a subset of the training data and use it to simulate the response of experts to test samples. Note that in large-scale datasets, like ImageNet-LT-d, no hold-out set is needed and DRAGON is trained on the training set. There, the familiarity bias is also present in the training data, as the models did not overfit the tail classes. Appendix C illustrates this effect in more detail.

## 6.3. Baselines and variants

We compared DRAGON with long-tail learning and unbalanced data approaches: **Focal Loss** [27], **Anchor Loss** [37], **Range Loss** [57], and **LDAM Loss** [7] are loss manipulation approaches for long-tail distributions. **FSLwF** [16], **OLTR** [29] and **Classifier-Balancing (CB)** [23] are algorithmic approaches in the long-tail learning benchmarks. **Mixture** and **Class Balanced Experts** [39] are late fusion approaches. **Mixture** resembles mixture-of-experts (MoE) [21] without EM optimization. It fuses the raw outputs of the two experts by a gating module. As with standard MoE models, the gating module is trained with visual-features as inputs.

For CUB-LT, SUN-LT, and AWA-LT, the visual expert was a linear layer over a pre-trained ResNet from [52, 50], which we trained with a balanced xent loss (CE Loss). The semantic expert was LAGO [4]. For ImageNet-LT-d, we

(a) Method	CUB-LT		SUN-LT		AWA-LT	
	$Acc_{PC}$	$Acc_{LT}$	$Acc_{PC}$	$Acc_{LT}$	$Acc_{PC}$	$Acc_{LT}$
<b>Vision-only</b>						
CE Loss* (VE)	53.0	65.5	33.7	40.0	73.7	93.4
Focal Loss [27]*	46.4	61.3	30.1	37.4	73.5	91.8
Anchor Loss [37]*	48.3	64.7	28.2	36.2	69.1	93.2
Range Loss [57]*	48.5	65.3	27.9	36.0	68.9	93.5
LDAM Loss [7]*	50.1	64.1	29.8	36.4	69.1	93.5
CB LWS [23]*	53.1	65.7	33.9	40.2	73.4	93.6
<b>Multi-modal</b>						
LAGO [4]* (SE)	54.8	66.0	18.2	18.3	74.0	93.0
CADA-VAE [38]*	48.3	57.4	32.8	35.1	73.5	89.5
<b>Late Fusion</b>						
Mixture	54.8	66.0	34.0	40.3	74.0	93.7
<b>DRAGON (ours)</b>	<b>57.8</b>	<b>67.7</b>	<b>34.8</b>	<b>40.4</b>	<b>74.1</b>	<b>94.1</b>
<b>DRAGON + Bal’Loss (ours)</b>	<b>60.1</b>	<b>66.5</b>	<b>36.1</b>	38.5	<b>76.2</b>	92.2

(b) ImageNet-LT Method	ResNet-10	ResNeXt-50			
	$Acc_{PC}$	$Acc_{MS}$	$Acc_{MED}$	$Acc_{FS}$	$Acc_{PC}$
<b>Vision-only</b>					
CE Loss* (VE)	34.8	65.9	37.5	7.7	44.4
FSLwF [16]	28.4	-	-	-	-
Focal Loss [27]	30.5	-	-	-	-
Range Loss [57]	30.7	-	-	-	-
OLTR [29]	35.6	-	-	-	37.7
CB LWS [23]	41.4	60.2	47.2	30.3	49.9
<b>Multi-modal</b>					
DEM [56]* (SE)	18.1	16.5	13.0	<b>50.8</b>	19.5
CADA-VAE [38]*	42.1	57.4	43.7	27.0	49.3
<b>Late Fusion</b>					
Mixture	40.7	63.8	36.3	23.9	45.1
Sharma et al. [39]	39.2	-	-	-	-
<b>DRAGON (ours)</b>	<b>43.1</b>	<b>66.0</b>	<b>38.3</b>	47.6	<b>51.2</b>
<b>DRAGON + Bal’Loss (ours)</b>	<b>46.5</b>	62.0	<b>47.4</b>	50.2	<b>53.5</b>

Table 2: **Smooth-tail distribution:** Rows with \* denote results reproduced by us. The rest were taken from [23, 29]. VE and SE refer to the *visual-expert* and *semantic-expert* that were used to train DRAGON. Bal. refers to training DRAGON with a balanced loss. (a) Comparing DRAGON with baselines on the long-tailed benchmark datasets. We report Per-Class Accuracy  $Acc_{PC}$  and Long-Tailed Accuracy  $Acc_{LT}$ . (b) Comparing DRAGON with baselines on the long-tailed ImageNet with word embeddings.

followed [23] and set the visual expert to be ResNet-10 or ResNeXt-50. The semantic expert was DEM [56].

#### 6.4. Evaluation metrics and protocol

**Evaluation Protocol:** The experiments for CUB-LT, SUN-LT, and AWA-LT follow the standard protocols set by [38, 52, 50], including their ResNet-101 features. Their split ensures that **none of the test classes appear in the training data used to train the ResNet-101 model.** For

ImageNet-LT-d we used the protocols in [29, 23] with the pre-trained ResNet-10 and ResNeXt-50 provided by [23].

**Evaluation metrics:** We evaluated DRAGON on the Smooth-Tail benchmark with the following metrics:

(a) **Per-Class Accuracy ( $Acc_{PC}$ ):** Balanced accuracy metric that uniformly averages the accuracy of each class  $\frac{1}{k} \sum_{y=1}^k Acc(y)$ , where  $Acc(y)$  is the accuracy of class  $y$ .

(b) **Long-Tailed Accuracy ( $Acc_{LT}$ ):** Test accuracy, where

CUB-LT	$Acc_{PC}$	$Acc_{LT}$
Max.	57.3	70.8
Avg.	57.0	70.0
Product	56.9	70.3
Mixture	53.7	66.5
<b>DRAGON (ours)</b>	<b>60.0</b>	<b>70.9</b>

Table 3: Comparing DRAGON against common late-fusion approaches on the validation set of CUB-LT.

the distribution over test classes is long-tailed like the training distribution. This is expected to be the typical case in real-world scenarios. See Appendix D.1 for more details.

(c) **Many-Shot, Medium-shot and Few-Shot accuracies:** For ImageNet-LT-d, we follow [29] and report accuracy for:  $Acc_{MS}$  (>100 training images),  $Acc_{MED}$  (20-100 images) and  $Acc_{FS}$  (< 20 images).

### 6.5. Results with smooth-tail distribution

Table 2 (a) provides the test accuracy for three long-tail benchmark datasets and compares DRAGON to baselines and individual components of the DRAGON model. DRAGON achieves higher accuracy compared with all competing methods, both with respect to class-balanced accuracy ( $Acc_{PC}$ ) and to test-distribution accuracy ( $Acc_{LT}$ ). Improving  $Acc_{LT}$  indicates that DRAGON effectively classifies head classes, which are heavily weighted in  $Acc_{LT}$ . At the same time, improving  $Acc_{PC}$  indicates that DRAGON also effectively classifies tail classes, which are up-weighted in  $Acc_{PC}$ .

Table 2 (b) provides the  $Acc_{PC}$  accuracy for ImageNet-LT-d. We can directly see the benefit of fusion information between modalities - the visual expert excels on many-shot classes,  $Acc_{MS}$ , while the semantic expert excels only on few-shot classes,  $Acc_{FS}$ . DRAGON recalibrate and fuse both experts to excel in all classes.

We further trained DRAGON with a balanced cross-entropy loss (Bal.). This strategy has a synergistic effect with DRAGON (last row in Table 2): It improves tail accuracy  $Acc_{PC}$  for all benchmarks while only marginally hurting head accuracy  $Acc_{LT}$ .

Table 3 compares DRAGON against common late fusion strategies, on the validation set of CUB-LT: *AVG* (averaging expert predictions), *Max* (taking the largest prediction), *Product* (multiplying expert predictions) and *Mixture*. We show that those approaches, which late fuse predictions of the two experts, are usually better at head classes ( $Acc_{LT}$ ) while giving less accurate results for tail classes ( $Acc_{PC}$ ). DRAGON achieves better results on both metrics because it also calibrates expert predictions. In the ablation study (Table 8) we compare our fusion module to more ablated fusion

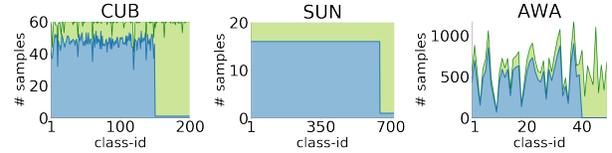


Figure 5: Two-level variants of CUB, SUN and AWA as in [38]. Blue: training set, green: test set.

components.

## 7. Two-Level (GFSL-d) benchmark

We follow the protocol of [38] on the original CUB, SUN, and AWA (Figure 5), to compare DRAGON in a Two-Level setting.

For those datasets, many-shot classes are kept as in the original train-set, while few-shot classes have an increasing number of shots: 1,2,5,10 and 20 (in SUN up to 10 shots).

### 7.1. Baselines and variants

We compared DRAGON with **LDAM Loss** [7] and with SoTA multi-modal GFSL-d approaches: **ReViSE** [43], **CA-VAE** [38], **DA-VAE** [38] and **CADA-VAE** [38]. Their results were obtained from the authors of [38], while LDAM results were reproduced by us.

The visual expert was a linear layer over a pre-trained ResNet from [52, 50], which we trained with a balanced cross-entropy loss. The semantic expert was LAGO [4].

### 7.2. Evaluation metrics

Following [38], we evaluated the Two-Level benchmark with the Harmonic mean metric ( $Acc_H$ ): It quantifies the overall performance of head and tail classes, by  $Acc_H = 2(Acc_{ms}Acc_{fs})/(Acc_{ms} + Acc_{fs})$ . Where,  $Acc_{ms}$  is the per-class accuracy over many-shot classes and  $Acc_{fs}$  is the per-class accuracy over few-shot classes.

### 7.3. Results with two-level distribution

Table 4 compares DRAGON with SoTA baselines on the Two-Level setup. Our model wins in CUB and SUN on all shots but loses on AWA for fewer than 10 samples. Furthermore, DRAGON gains better results when the number of shots increases in contrast to complex generative models like CADA-VAE [38]. Appendix E.1 provides results for  $Acc_{fs}$  and  $Acc_{ms}$  with 1,2,5,10,20 shots.

## 8. Vision-only long-tail learning

The approach presented in this paper focuses on learning from two modalities, vision and language. To understand the effect of re-calibrating we now study a simpler variant of DRAGON that can be applied to the more common vision-only long-tail learning. We name it smDRAGON for

Two-Level # shots	CUB					SUN				AWA				
	1	2	5	10	20	1	2	5	10	1	2	5	10	20
LDAM [7]*	2.4	10.9	36.0	52.2	61.5	4.3	11.5	26.6	37.0	12.4	24.8	41.1	57.0	68.6
REVISE [43]	36.3	41.1	44.6	50.9	-	27.4	33.4	37.4	40.8	56.1	60.3	64.1	67.8	-
CA-VAE [38]	50.6	54.4	59.6	62.2	-	37.8	41.4	44.2	45.8	64.0	71.3	76.6	79.0	-
DA-VAE [38]	49.2	54.6	58.8	60.8	-	37.8	40.8	43.6	45.1	68.0	73.0	75.6	76.8	-
CADA-VAE [38]	55.2	59.2	63.0	64.9	66.0	40.6	43.0	46.0	47.6	<b>69.6</b>	<b>73.7</b>	<b>78.1</b>	80.2	80.9
CE Loss* (VE)	1.2	6.9	30.2	50.2	60.9	1.8	8.9	25.1	38.3	11.0	20.0	47.8	69.9	73.9
LAGO* (SE)	23.0	33.2	49.0	58.6	64.8	19.5	23.2	25.6	27.8	20.2	33.0	59.0	68.7	75.8
<b>DRAGON (ours)</b>	<b>55.3</b>	<b>59.2</b>	<b>63.5</b>	<b>67.8</b>	<b>69.9</b>	<b>41.0</b>	<b>43.8</b>	<b>46.7</b>	<b>48.2</b>	67.1	69.1	76.7	<b>81.9</b>	<b>83.3</b>

Table 4: **Two-Level distributions:** Comparing DRAGON on Two-Level CUB, SUN and AWA with SoTA GFSL models and baselines and with increasing number of few-shot training samples. Values denote the Harmonic mean  $Acc_H$ . VE and SE refer to the *visual-expert* and *semantic-expert* that were used to train DRAGON.

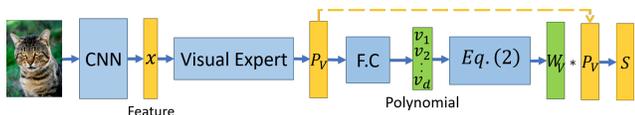


Figure 6: Architecture of vision-only smDRAGON.

*single-modality*-DRAGON, and show that it achieves new state-of-the-art results, compared to uni-modals baselines, on ImageNet-LT, Places365-LT, CIFAR-10 and CIFAR-100. On iNaturalist smDRAGON is comparable to SoTA. To adapt to single-modality, we train smDRAGON only on the predictions of the visual-expert. It outputs a single set of coefficients  $\{\mathbf{w}_V(y)\}_{y \in \mathcal{Y}}$  to rescale the predictions of the visual expert, instead of two sets of coefficients. Subsequently, Eq. (1) reduces to  $S(y) = \mathbf{w}_V(y)p_V(y)$ .

In other words, smDRAGON is a simplified version of DRAGON that is trained on the predictions of the **visual-expert only** (no class descriptors being used). smDRAGON takes the predictions of a frozen visual-expert and rescale it by learning a single set of polynomial coefficients. During inference, smDRAGON balances the visual-expert predictions in a sample-by-sample basis.

Tables 5 and 6 compare smDRAGON against approaches in the unbalanced CIFAR-10 and CIFAR-100 benchmarks, as presented in [7]: **CE Loss**, **Resample [11]**, **Reweight [11]**, **Focal [11]** and **LDAM Loss [7]**. *DRW* [7] denotes models that were trained with the training schedule proposed by [7].

Table 7 compares smDRAGON with popular baselines and recent long-tail learning approaches in the ImageNet-LT and Places-LT benchmarks. Those are the same baselines as in the Smooth-Tail setup (Section 6).

The results demonstrate that (1) smDRAGON outperforms all baselines and (2) combining smDRAGON with SoTA approaches (LDAM or DRW) has a synergistic effect.

Table 10 compares DRAGON with smDRAGON on CUB-

LT. It shows that fusing information between modalities (third row) gives better results than re-scaling expert predictions alone (first and second rows). Finally, on iNaturalist 2018, smDRAGON is comparable to SoTA, reaching 69.1% compared to 69.5% (CB-LWS [23]).

## 9. Ablation study

To understand the contribution of individual components of DRAGON, we carried ablation experiments. We report results on the validation set, which were consistent with the test set (Appendix E.2).

**Fusion-Module Architecture:** Table 8 compares the performance of various components of the fusion-module on CUB-LT. (1) *F.C.*: predicts  $\lambda$  using a fully-connected layer over  $\mathbf{p}_V, \mathbf{p}_S$ , no re-balancing ( $\mathbf{w}_V(y) = \mathbf{w}_S(y) = 1, \forall y$ ). (2) *F.C. &  $1/n_y$  rescale*: learns  $\lambda$  as in *F.C.*, rescales experts predictions by  $n_y$ . (3) *F.C. & non-parametric rescale*: learns  $\lambda$  as *F.C.* and rescales both experts predictions by a learned non-parametric weight for each class instead of a polynomial. (4) *Conv. & non-parametric rescale*: like (3), then applies sorting and convolution (Section 4.2). (5) *Conv. & single parametric rescale* replaces the non-parametric re-scaling weights by a single polynomial of parametrized weights. (6) DRAGON is our full approach described in Section 4. The comparison shows that rescaling expert predictions significantly improves  $Acc_{PC}$  and that reducing the number of parameters using the convolutional layer is important.

To quantify the contribution of *per-sample* weighting, Table 9 compares it against *per-class* weighting on three long-tail benchmarks: ImageNet-LT, Places365-LT and CIFAR100-LT. To keep the comparison fair, this was done using vision-only ( $\lambda = 1$ ), and sweeping over the same set of hyper parameters. To gain more intuition on how per-sample weighting helps, Fig. 7 plots the per-sample weights of four images from a ImageNet-LT head class (mouse-

Vision-Only	Unb. CIFAR-10				Unb. CIFAR-100			
	long-tail		two-level		long-tail		two-level	
Imbalance Type								
Imbalance Ratio	100	10	100	10	100	10	100	10
[11] ReSample	29.45	13.21	38.14	15.41	66.56	44.94	66.23	46.92
[11] ReWeight	27.63	13.46	38.06	16.20	66.01	42.88	78.69	47.52
[11] Focal	25.43	12.90	39.73	16.54	63.98	42.01	80.24	49.98
CE [7]	29.64	13.61	36.70	17.50	61.68	44.30	61.45	45.37
CE* (VE)	29.81	13.12	36.61	17.78	61.72	43.77	61.59	45.75
Focal [27]	29.62	13.34	36.09	16.36	61.59	44.22	61.43	46.54
LDAM [7]	26.65	13.04	33.42	15.00	60.40	43.09	60.42	43.73
<b>smDRAGON (ours)</b>	<b>22.08</b>	<b>12.17</b>	<b>27.10</b>	<b>12.38</b>	<b>58.01</b>	<b>42.22</b>	<b>54.43</b>	<b>40.97</b>

Table 5: **Vision-only long-tail.** Error rate of ResNet32 on unbalanced CIFAR-10 and CIFAR-100 [7], comparing smDRAGON and SoTA techniques. smDRAGON was trained over predictions of the cross-entropy model (CE\*). Reported values are the top-1 validation error. Asterisks \* denote results reproduced using published code.

Vision-Only	Unb. CIFAR-10				Unb. CIFAR-100			
	long-tail		two-level		long-tail		two-level	
Imbalance Type								
Imbalance Ratio	100	10	100	10	100	10	100	10
CE-DRW* (VE1)	24.73	13.52	28.65	13.90	59.23	42.19	58.93	45.00
M-DRW [15]	24.94	13.57	26.67	13.17	59.49	43.49	58.91	44.72
LDAM-DRW [7]	22.97	11.84	23.08	12.19	57.96	41.29	54.64	40.54
LDAM-DRW* [7] (VE2)	22.96	11.84	23.41	12.20	57.89	41.61	54.65	43.48
<b>VE1 + smDRAGON</b>	<b>20.37</b>	12.06	21.54	<b>11.94</b>	<b>56.50</b>	42.11	<b>53.32</b>	40.66
<b>VE2 + smDRAGON</b>	21.22	<b>11.84</b>	<b>20.64</b>	12.37	56.70	<b>41.23</b>	54.07	<b>40.35</b>

Table 6: **Vision-only long-tail:** smDRAGON was trained on top models trained with DRW (VE1) or LDAM-DRW (VE2) [7]. Similar to Table 5 except that all models were trained with DRW schedule [7]. Reported values are top-1 validation error. Asterisks \* denote results that we reproduced using code published by the authors of [7].

Vision-Only	Places365-LT		ImageNet-LT		Vision-Only	iNaturalist ResNet-50
	ResNet-50	ResNet-10	ResNet-50	ResNeXt-50		
CE Loss* (VE)	30.2	34.8	44.4			
Bal' Loss	32.4	33.1	-			
Lifted Loss [41]	35.2	30.8	-		[11] Focal	61.1
Focal Loss [27]	34.6	30.5	-		LDAM [7]	64.6
Range Loss [57]	35.1	30.7	-		LDAM-DRW [7]	68.0
FSLwF [16]	34.9	28.4	-		CB $\tau$ -norm [23]	69.3
OLTR [29]	35.9	34.1	37.7		CB LWS [23]	69.5
CB $\tau$ -norm [23]	37.9	40.6	49.4			
CB LWS [23]	37.6	41.4	49.9		<b>smDRAGON (ours)</b>	69.1
<b>smDRAGON (ours)</b>	<b>38.1</b>	<b>42.0</b>	<b>50.1</b>			

Table 7: **Vision-only long-tail:** Baseline results where copied directly from [7] and [23]. **Left:** smDRAGON achieves better  $Acc_{PC}$  on Places365-LT and ImageNet-LT. **Right:** Comparing smDRAGON on long-tailed iNaturalist. smDRAGON achieve comparable results compared to SoTA baselines.

trap). Per-sample weighs more strongly "easy" samples (low entropy) than non-typical samples. This illustrates that

per-sample weighting does not penalize "justified" high-confidence predictions if they happen to arrive from a head

	$Acc_{PC}$	$Acc_{LT}$	$\#params$
F.C.	54.0	67.1	403
F.C. & $1/n_y$ RESCALE	56.7	60.3	403
F.C. & NON-PARAMETRIC RESCALE	58.2	68.0	81,406
CONV. & NON-PARAMETRIC RESCALE	58.7	68.2	40,612
CONV. & SINGLE PARAMETRIC RESCALE	59.0	67.5	613
<b>DRAGON (OURS)</b>	<b>60.0</b>	<b>70.9</b>	<b>1,015</b>

Table 8: Ablation study, comparing different fusion and re-scaling approaches. The results show the contribution of the convolutional backbone and the re-scaling method for the two experts (validation set, CUB-LT).

	Places365-LT	ImageNet-LT		CIFAR100-LT
	ResNet-50	ResNet-10	ResNeXt-50	ResNet-32
CE Loss (VE)	30.2	34.8	44.4	38.3
Per-class	36.9	40.0	49.2	40.5
Per-sample	<b>38.1</b>	<b>42.0</b>	<b>50.1</b>	<b>42.0</b>

Table 9: Ablation of per-sample weighting on vision-only benchmarks. VE refers to *visual-expert*.

Ground Truth: **Mousetrap (Head Class)**



Prediction confidence:	0.9	0.9	0.8	0.7
Per-class weight:	0.7	0.7	0.7	0.7
Per-sample weight:	0.9	0.7	0.3	0.1

Figure 7: Using per-sample weighting, images typical for the class Mousetrap (softmax has low-entropy) are weighed more strongly than non-typical images (softmax has high-entropy). Per-class weighting reweighs all samples for that class the same (0.7), hurting recognition of typical images.

	$Acc_{PC}$	$Acc_{LT}$
VISUAL EXPERT + SMDRAGON	55.8	66.0
SEMANTIC EXPERT + SMDRAGON	57.7	63.4
<b>DRAGON (OURS)</b>	<b>60.1</b>	<b>67.7</b>

Table 10: Ablation study, comparing smDRAGON to DRAGON on CUB-LT: Fusing information between modalities improves performance (test set, CUB-LT).

Sorting	$Acc_{PC}$	$Acc_{LT}$
No Sorting	58.7	68.2
<b>Sorting By Visual Expert</b>	<b>60.0</b>	<b>70.9</b>
Sorting By Semantic Expert	60.0	70.8

Table 11: Ablation study, quantifying the contribution of sorting the fusion-module inputs (validation set, CUB-LT).

class. At the same time, per-sample weighting gives more chance to tail classes, reducing the familiarity bias.

**Sharing order statistics (sorting):** Table 11 quantifies the benefit of sorting expert predictions. As discussed in Section 4.3, sorting enables sharing of information across classes by fixing the input location of each order statistic (max,  $2^{nd}$  max etc.).

## 10. Conclusion

This paper discussed two key challenges for learning with long-tail unbalanced data: A “familiarity bias”, where models favor head classes, and low accuracy over tail classes due to lack of samples. We address these challenges, with DRAGON, a late-fusion architecture for visual recognition that learns with per-class semantic information. It outperforms existing methods on new long-tailed versions of ImageNet, CUB, SUN, and AWA. It further sets new SoTA on a Two-Level benchmark [38]. Finally, a single-modality variant of DRAGON improves accuracy over standard long-tail learning benchmarks, including ImageNet-LT, Places365-LT, and unbalanced CIFAR-10/100. These results show that information about the number of samples per-class can be effectively used to reduce prediction biases.

Strongly unbalanced data with a long-tail is ubiquitous in numerous domains and problems. The results in this paper show that a light-weight late-fusion model can be used to address many of the challenges posed by class imbalance.

## Acknowledgments

DS was funded by a grant from the Israeli innovation authority, through the AVATAR consortium and by a grant from the Israel Science Foundation (ISF 737/2018). Study was also funded by an equipment grant to GC and Bar-Ilan University from the Israel Science Foundation (ISF 2332/18).

## References

- [1] G. Arora, V-K. Verma, A. Mishra, and P. Rai. Generalized zero-shot learning via synthesized examples. In *CVPR*, 2018. 3
- [2] A. Asvadi, L. Garrote, C. Premebida, P. Peixoto, and U. Nunes. Multimodal vehicle detection: fusing 3d-lidar and color camera data. *Pattern Recognition Letters*, 2018. 3
- [3] Y. Atzmon and G. Chechik. Adaptive confidence smoothing for generalized zero-shot learning. *CVPR*, 2018. 3, 5
- [4] Y. Atzmon and G. Chechik. Probabilistic and-or attribute grouping for zero-shot learning. In *UAI*, 2018. 2, 3, 4, 6, 7, 8, 14, 17
- [5] S. Beery, Y. Liu, D. Morris, J. Piavis, A. Kapoor, M. Meister, and P. Perona. Synthetic examples improve generalization for rare classes. *Preprint arXiv:1904.05916*, 2019. 2
- [6] M. Buda, A. Maki, and M. Mazurowski. A systematic study of the class imbalance problem in convolutional neural networks. *Neural Networks*, 2018. 1, 2
- [7] K. Cao, C. Wei, A. Gaidon, N. Arechiga, and T. Ma. Learning imbalanced datasets with label-distribution-aware margin loss. *NIPS*, 2019. 2, 5, 6, 7, 8, 9, 10, 17
- [8] S. Changpinyo, W. L. Chao, B. Gong, and F. Sha. Synthesized classifiers for zero-shot learning. In *CVPR*, 2016. 3, 6
- [9] R. Chao, S. Changpinyo, B. Gong, and Sha F. An empirical study and analysis of generalized zero-shot learning for object recognition in the wild. In *ICCV*, 2016. 3
- [10] Nitesh V. Chawla, K. Bowyer, L. Hall, and W. P. Kegelmeyer. Smote: Synthetic minority over-sampling technique. *ArXiv*, 2002. 2
- [11] Y. Cui, M. Jia, T. Lin, Y. Song, and S. Belongie. Class-balanced loss based on effective number of samples. *CVPR*, 2019. 9, 10
- [12] J. Deng, W. Dong, R. Socher, L. Li, K. Li, and F. Li. Imagenet: A large-scale hierarchical image database. *CVPR*, 2009. 1
- [13] Chris Drummond. C 4 . 5 , class imbalance , and cost sensitivity : Why under-sampling beats oversampling. 2003. 2
- [14] R. Felix, V. Kumar, I. Reid, and G. Carneiro. Multi-modal cycle-consistent generalized zero-shot learning. In *ECCV*, 2018. 3
- [15] W Feng, j. Cheng, W. Liu, and H. Liu. Additive margin softmax for face verification. *IEEE Signal Processing Letters*, 2018. 10
- [16] S. Gidaris and Komodakis.N. Dynamic few-shot visual learning without forgetting. *CVPR*, 2018. 6, 7, 10
- [17] C. Guo, G. Pleiss, Y. Sun, and K. Weinberger. On calibration of modern neural networks. In *ICML*, 2017. 2, 14, 15
- [18] H. Han, W. Wang, and B. Mao. Borderline-smote: A new over-sampling method in imbalanced data sets learning. In *ICIC*, 2005. 2
- [19] H. He and E. A. Garcia. Learning from imbalanced data. *IEEE Transactions on Knowledge and Data Engineering*, 2009. 2
- [20] G. Horn, O. Aodha, Y. Song, A. Shepard, H. Adam, P. Perona, and S. Belongie. The inaturalist challenge 2017 dataset. *ArXiv*, 2017. 5
- [21] R. Jacobs, S. Nowlan, M. Jordan, and G. Hinton. Adaptive mixtures of local experts. *Neural computation*, 1991. 6
- [22] S. et al. Kahou. Combining modality specific deep neural networks for emotion recognition in video. In *ICMI '13*, 2013. 3
- [23] B. Kang, S. Xie, M. Rohrbach, M. Yan, A. Gordo, J. Feng, and Y. Kalantidis. Decoupling representation and classifier for long-tailed recognition. *ICLR*, 2020. 2, 6, 7, 9, 10
- [24] D. Kingma and J. Ba. Adam: A Method for Stochastic Optimization. *ICLR*, 2015. 15
- [25] M Kull, M. Perelló-Nieto, M. Kängsepp, T. Filho, Hao. Song, and P. Flach. Beyond temperature scaling: Obtaining well-calibrated multiclass probabilities with dirichlet calibration. *ArXiv*, 2019. 2, 14, 15
- [26] C.H. Lampert, H. Nickisch, and S. Harmeling. Learning to detect unseen object classes by between-class attribute transfer. In *CVPR*, 2009. 2, 3, 4, 5, 6, 14
- [27] T. Lin, P. Goyal, R. Girshick, K. He, and P. Dollár. Focal loss for dense object detection. *ICCV*, 2017. 2, 6, 7, 10, 16
- [28] K. Liu, Y. Li, N. Xu, and P. Natarajan. Learn to combine modalities in multimodal deep learning. *arXiv preprint arXiv:1805.11730*, 2018. 3
- [29] Z. Liu, Z. Miao, X. Zhan, J. Wang, B. Gong, and S. Yu. Large-scale long-tailed recognition in an open world. In *CVPR*, 2019. 2, 5, 6, 7, 8, 10
- [30] T. Mikolov, I. Sutskever, K. Chen, G. Corrado, and J. Dean. Distributed representations of words and phrases and their compositionality. *ArXiv*, 2013. 6
- [31] A. Mishra, M. Reddy, A. Mittal, and H. A. Murthy. A generative model for zero shot learning using conditional variational autoencoders. In *WACV*, 2018. 3
- [32] A. Pambala, T. Dutta, and S. Biswas. Generative model with semantic embedding and integrated classifier for generalized zero-shot learning. In *WACV*, 2020. 2, 3
- [33] G. Patterson and J. Hays. Sun attribute database: Discovering, annotating, and recognizing scene attributes. In *CVPR*, 2012. 5, 6
- [34] J. Platt. Probabilistic outputs for support vector machines and comparisons to regularized likelihood methods. In *Advances in large margin classifiers*, 1999. 15
- [35] S. Pouyanfar, T. Wang, and S. Chen. A multi-label multimodal deep learning framework for imbalanced data classification. In *MIPR*, 2019. 3
- [36] S. Reed, Z. Akata, H. Lee, and B. Schiele. Learning deep representations of fine-grained visual descriptions. In *CVPR*, 2016. 2, 4
- [37] S. Ryou, S. Jeong, and P. Perona. Anchor loss: Modulating loss scale based on prediction difficulty. In *ICCV*, 2019. 2, 6, 7
- [38] E. Schönfeld, S. Ebrahimi, S. Sinha, T. Darrell, and Z. Akata. Generalized zero-shot learning via aligned variational autoencoders. In *CVPR*, 2019. 2, 3, 5, 6, 7, 8, 9, 11, 15, 16, 17
- [39] S. Sharma, N. Yu, M. Fritz, and B. Schiele. Long-tailed recognition using class-balanced experts. *arXiv preprint arXiv:2004.03706*, 2020. 6, 7

- [40] R. Socher, M. Ganjoo, C.D. Manning, and A. Ng. Zero-shot learning through cross-modal transfer. In *NIPS*, 2013. 3, 4
- [41] H. Song, Y. Xiang, S. Jegelka, and S. Savarese. Deep metric learning via lifted structured feature embedding. *CVPR*, 2016. 10
- [42] F. Sung, Y. Yang, L. Zhang, T. Xiang, P. HS Torr, and T. M Hospedales. Learning to compare: Relation network for few-shot learning. In *CVPR*, 2018. 3
- [43] Y-H. Tsai, L-K. Huang, and R. Salakhutdinov. Learning robust visual-semantic embeddings. In *ICCV*, 2017. 2, 3, 8, 9, 17
- [44] A. Tversky and D. Kahneman. Availability: A heuristic for judging frequency and probability. *Cognitive psychology*, 1973. 2
- [45] G. Van Horn and P. Perona. The devil is in the tails: Fine-grained classification in the wild. *arXiv preprint arXiv:1709.01450*, 2017. 1
- [46] 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. 4, 5, 6
- [47] B. Wallace and I. Dahabreh. Improving class probability estimates for imbalanced data. *Knowledge and Information Systems*, 2013. 1, 2
- [48] B. C. Wallace and I. J. Dahabreh. Class probability estimates are unreliable for imbalanced data (and how to fix them). In *2012 IEEE 12th International Conference on Data Mining*. 1, 2
- [49] Yu-Xiong Wang, Deva Ramanan, and Martial Hebert. Learning to model the tail. In *NIPS*, 2017. 2
- [50] Y. Xian, C.H. Lampert, B. Schiele, and Z. Akata. Zero-shot learning - A comprehensive evaluation of the good, the bad and the ugly. *TPAMI*, 2018. 6, 7, 8
- [51] Y. Xian, T. Lorenz, B. Schiele, and Z. Akata. Feature generating networks for zero-shot learning. In *CVPR*, 2018. 3, 6
- [52] Y. Xian, B. Schiele, and Z. Akata. Zero-shot learning - the good, the bad and the ugly. In *CVPR*, 2017. 2, 3, 4, 6, 7, 8, 14
- [53] Yongqin Xian, Saurabh Sharma, Bernt Schiele, and Zeynep Akata. F-vaegan-d2: A feature generating framework for any-shot learning. *CVPR*, 2019. 2, 3
- [54] H. Zhang and P. Koniusz. Model selection for generalized zero-shot learning. In *ECCV*, 2018. 3
- [55] Hongguang Zhang and Piotr Koniusz. Zero-shot kernel learning. In *CVPR*, 2018. 3
- [56] L. Zhang, T. Xiang, and S. Gong. Learning a deep embedding model for zero-shot learning. In *CVPR*, 2017. 3, 7
- [57] X. Zhang, Z. Fang, Y. Wen, Z. Li, and Y. Qiao. Range loss for deep face recognition with long-tailed training data. *ICCV*, 2017. 6, 7, 10
- [58] Z. Zhang, P. Luo, C. Loy, and X. Tang. Learning deep representation for face alignment with auxiliary attributes. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 38, 2016. 3
- [59] Y. Zhu, M. Elhoseiny, B. Liu, X. Peng, and A. Elgammal. A generative adversarial approach for zero-shot learning from noisy texts. In *CVPR*, 2018. 3, 4

## A. Additional analysis of the familiarity effect

Here we provide a deeper analysis showing that DRAGON effectively addresses the “familiarity bias”.

**The familiarity bias causes models to incorrectly favor head classes:** Figure S8(a) shows the confusion matrix of a standard ResNet-101 trained on CUB-LT, as computed on the validation set. Classes, of the confusion matrix, are ordered by a decreasing number of training samples, with class #1 having many samples and class #200 have few samples. Black dots denote count larger than 15.

It illustrates two effects. First, the trained model correctly classifies head classes, based on the fact that the top rows have no incorrect (off-diagonal) predictions. Second, for mid and tail classes, predictions are clearly biased towards the head, since there are many more off-diagonal predictions to the left (head class predictions).

**DRAGON corrects for the familiarity bias:** Figures S8(b) and S8(c) demonstrate that DRAGON learns to offset the familiarity bias. The left panel (b) shows the familiarity effect on CUB-LT before recalibration. The right panel (c) shows that DRAGON corrects the familiarity bias and produces a more balanced average confidence across the head and tail classes.

**DRAGON re-calibrate predictions:** In the main paper (Figure 1(c)) we showed that a model that is trained on unbalanced data has higher confidence for head classes and it over-estimate them. **By reversing the familiarity bias, smDRAGON, implicitly, also re-calibrate experts predictions.** Figure 9 compares the reliability diagrams for smDRAGON against raw *ResNeXt-152*, *Temp Scaling* [17] and *Dirichlet Calibration* [25] (common and SoTA calibration approaches). We report both per-class-accuracy (ACC) and expected-calibration-error (ECE) for each model.

## B. Visual experts are better at the head, semantic experts excel at the tail

Here we provide supporting evidence to our observation from Section 4 of the main paper that semantic experts are better at the tail: “*Semantic descriptions of classes can be very useful for tail (low-shot) classes, because they allow models to recognize classes even with few or no training samples [26, 52, 4]*”. Additionally, we demonstrate that the visual expert is better for the many-shot regime.

We focus on the Two-Level CUB distribution, and evaluate the accuracy for the many shot classes when restricting predictions to these classes (many-among-many), and separately the accuracy for the few-shot classes when predictions are restricted to these tail classes (few-among-few).

Figure S10(a) shows the accuracy over few-shot classes of both experts in few-among-few setting. The semantic expert outperforms the visual one, and this effect stronger with fewer samples. For example, with 1-shot learning, the se-

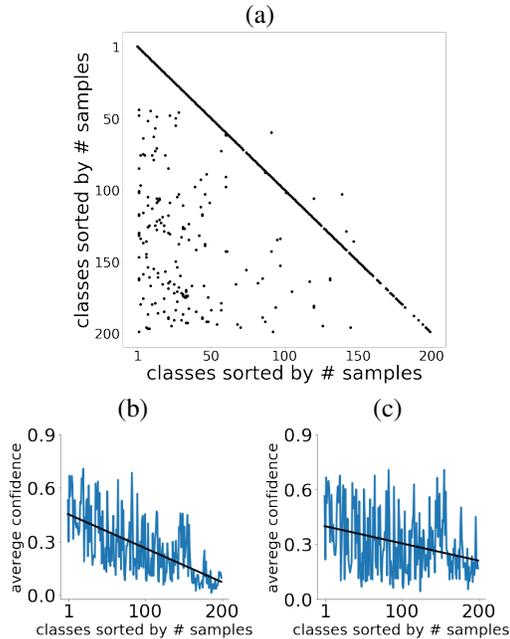


Fig. S 8: DRAGON learns to offset the familiarity bias. (a) A confusion matrix of a ResNet101 trained on CUB-LT as a function of the number of samples per class. The matrix shows markers for pairs of (gt, predicted) whose count is larger than 15. (b) Average-confidence per-class of the classifier. (c) Similar curve as (b) but for DRAGON. Black lines depict a linear regression line. DRAGON per-class confidence has smaller dependence on the number of samples in the train.

semantic expert is almost 100% better than the Visual Expert. Additionally, when we measure the accuracy of the many-shot classes in the many-among-many setup (accuracy at the head), the visual expert is better than the semantic expert Figure S10(b).

## C. Training the fusion-module in small scale datasets

Our goal is to have the fusion-module learn to capture the correlations between the number of training samples and the output confidence (the familiarity bias), so it can adjust for it. Unfortunately, while the familiarity effect is substantial in the validation data and the test data, it may not present in the training data in small scale datasets. The reason is: Models tend to overfit and become overconfident over rare classes in the training set. This effect is illustrated in Figure S11(a) (compare Train versus Validation curves) for CUB-LT. We observed the effect in also in SUN-LT and AWA-LT.

To address this mismatch, we hold-out 50% of the samples of the tail classes and 20% of the samples of the head classes of the training data and use it to simulate the re-

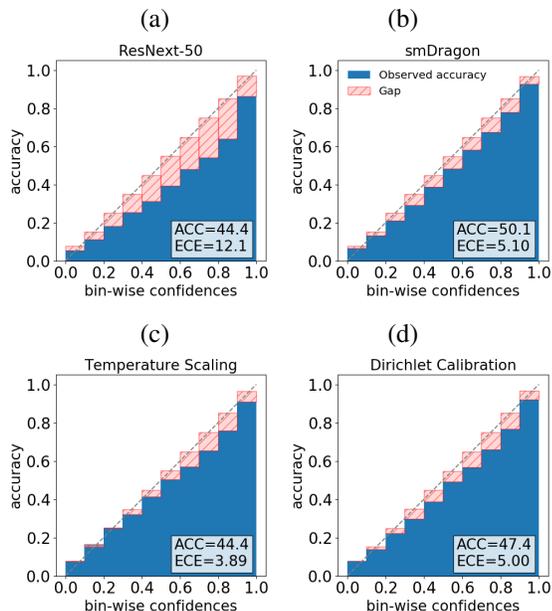


Fig. S 9: Reliability Diagrams on ImageNet-LT, for (a) raw ResNext-50, (b) smDRAGON, (c) Temperature-Scaling [17] and (d) Dirichlet-Calibration [25]. We report expected-calibration-error (ECE) and per-class accuracy (ACC)

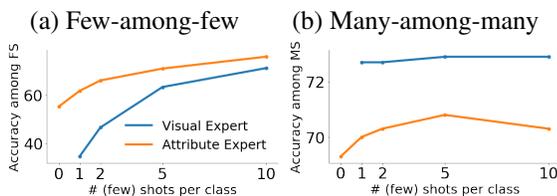


Fig. S 10: Accuracy as a function of number of samples at the tail, of the *Visual Expert* and the *Semantic Expert* used in our study. (a) Accuracy among few-shot classes; The Semantic expert outperforms the visual expert. (b) Accuracy among many-shot classes; The Visual expert outperforms, regardless of the number of samples at the tail.

sponse of experts to test samples. This set is used for training the fusion-module.

Note, that after training the fusion module, we re-train the experts on all the training data (including the hold-out set), in order to use all data available. (See Section D for more details).

In large-scale datasets, like ImageNet-LT, no hold-out set is needed and DRAGON is trained on the training set. There, the familiarity bias is also present on the training data (Figure S11(b)), as the models did not overfit the tail classes.

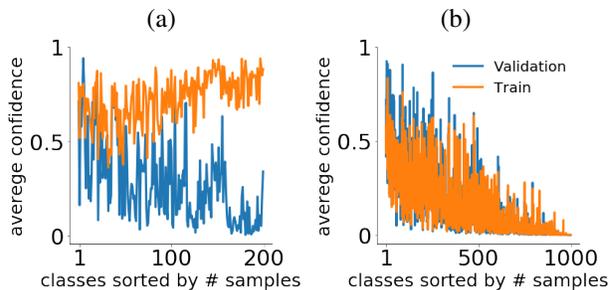


Fig. S 11: The familiarity bias effect: (a) On CUB-LT the effect is strong on validation samples (blue) but not on training samples (orange). (b) On ImageNet-LT the effect is prominent on both train and validation samples.

## D. Implementation details

**Training:** Considering the observation from Section C, regarding CUB-LT, SUN-LT and AWA-LT, we train the architecture in three steps: First, we train each expert on the training data excluding the hold-out set. Second, we freeze the expert weights and train the *fusion-module* on all the training set. Finally, we re-train the experts on all the training data in order to use all data available. For the hold-out set, we randomly draw half the samples of the tail classes and 20% of the samples of the head classes. For inference, we use the fusion-module trained at the second step with the experts trained at the third step.

**Platt-scaling:** We used Platt-scaling [34] to tune the combination coefficient  $\lambda$  by adding constant bias  $\beta$  and applying a sigmoid on top of its scores:  $\lambda = \sigma[f_0 - \beta]$ , where  $\beta$  is a hyperparameter selected with cross validation.

**Fusion-module:** We trained the fusion-module using ADAM [24] optimizer. For large-scale datasets, like ImageNet-LT, Places-LT and iNaturalist, we used  $L_2$  regularization, selected by hyperparameter optimization using grid search  $\in \{10^{-5}, 10^{-4}, 10^{-3}\}$ , to avoid overfitting.

**Hyper-parameter tuning:** We determined the number of training epochs (early-stopping), selected architecture alternatives, and tuned hyperparameters using the validation set, using  $Acc_{LT}$  for *Smooth-Tail* and *Vision-only*, and  $Acc_{PC}$  for *Two-Level*.

For *DRAGON*: We optimized the following hyperparameters: (1) Number of filters in the convolution layer  $\in \{1, \dots, 4\}$ . (2) Degree of polynomial in Eq.2  $\in \{2, 3, 4\}$ . (3) Learning rate  $\in \{10^{-5}, 10^{-4}, 10^{-3}\}$ . (4) Bias term of Platts rescaling  $\beta \in [-2, 2]$ .

For *CADA-VAE* [38]: We applied a grid search for the latent embedding space  $\in [12, 25, 50, 64, 100, 200, 250]$ , variational-autoencoder learning rate  $\in [0.0001, 0.00015, \dots, 0.015]$  and classifier learning rate  $\in [0.0001, 0.0005, \dots, 0.1]$ . We used a batch size of 64.

Sorting	$Acc_{PC}$	$Acc_{LT}$
No-Sorting	58.5	57.0
<b>Sorting-By-Visual-Expert</b>	<b>60.1</b>	<b>67.7</b>
Sorting-By-Semantic-Expert	59.8	67.7

Table S 12: Ablation study, quantifying the contribution of sorting the fusion-module inputs (test set, CUB-LT).

For *Focal Loss* [27]: We applied a grid-search for gamma  $\in [1, 2, \dots, 15]$  and alpha  $\in [0.1, 0.2, 0.5, 0.75, 0.9, 1]$ .

For *Range Loss* [27]: We applied a grid-search for alpha  $\in [0.1, 0.2, 0.5, 0.75, 0.9, 1]$  and beta  $\in [0.1, 0.2, 0.5, 0.75, 0.9, 1]$ .

For *Anchor Loss* [27]: We applied a grid-search for gamma  $\in [0.1, 0.5, 1, \dots, 15]$  and slack  $\in [0.001, 0.005, 0.01, \dots, 0.5]$ .

For *LDAM Loss* [27] we applied a grid-search for C  $\in [0.1, 0.2, \dots, 0.9]$ .

### D.1. Computing $Acc_{LT}$ :

$Acc_{LT}$  measures the accuracy over a test distribution that resembles the training distribution. However, the test and validation samples of CUB-LT, SUN-LT and AWA-LT have a different distribution because they were originated from an approximate uniform distribution. Thus, to compute  $Acc_{LT}$  we measure the accuracy for each individual class, and then take a weighted sum according to the class frequencies in the training set. Specifically, for each class, we assign a weight  $P_{train}(y)$  according to the train-set distribution such that  $0 < P_{train}(y) < 1$  and  $\sum_y P_{train}(y) = 1$ . Then we compute the accuracy per class and report the weighted average across all classes:  $Acc_{LT} = \sum_{y=1}^k P_{train}(y) acc(y)$ . This is equivalent to transforming the test set to have the same distribution as the train set.

### D.2. A clarification about the Smooth-Tail benchmark

In this section, we explain how the long-tail benchmark was aligned with the two-level benchmark, as was mentioned in the paragraph that describes the long-tailed datasets (Section 6 of the main paper).

To align the long-tail benchmark with the two-level benchmark, we first ordered the classes according to their number of samples in the *two-level* distribution. Then we calculated the number of samples for each class according to the required long-tail distribution, and accordingly drew samples to construct the training set.

Architecture	$Acc_{PC}$	$Acc_{LT}$
F.C.	56.4	66.3
F.C. & $1/n_y$ re-scale	56.4	56.2
F.C. & non-parametric re-scale	58.2	64.3
Conv. & non-parametric re-scale	58.4	67.1
Conv. & single parametric re-scale	59.3	64.3
<b>DRAGON (ours)</b>	<b>60.1</b>	<b>67.7</b>

Table S 13: Ablation study, comparing different fusion and re-scaling approaches. The results show the contribution of the convolutional backbone and the re-scaling method for the two experts (test set, CUB-LT).

Training Process	$Acc_{PC}$
All-Train	56.6
End-To-End	46.4
<b>Three-Stage-Training</b>	<b>60.1</b>

Table S 14: Ablation study, quantifying the contribution the effect of three-stage training as proposed in Section D. (test-set, CUB)

### D.3. Training CADA on Smooth-Tail benchmark

In this section, we explain how we trained CADA-VAE [38] for the long-tail benchmark.

To evaluate CADA-VAE [38] on long-tail benchmarks we used the code published by the authors and followed the training protocol exactly as they used for the two-level distribution. Since the protocol relies on a hard distinction between head classes and tail classes, we had to choose where to partition the smooth long-tail distribution to head and tail. Our solution is simple. It is based on the fact that we aligned the order of classes in the long-tail distribution to be the same order as in the two-level split (Section D.2). The alignment allowed us to use the same partition to head and tail as used for the two-level benchmark.

## E. Additional metrics

### E.1. $Acc_{fs}$ and $Acc_{ms}$ on Two-Level benchmarks

Table S15 provides the results of  $Acc_{fs}$  and  $Acc_{ms}$  (described in section 7) for the Two-Level benchmark. We show results for 1,2,5 and 10-shots. At the main paper we reported the results for the  $Acc_H$  metrics, which is derived from  $Acc_{fs}$  and  $Acc_{ms}$  reported here.

### E.2. Ablation results on the test set

In the main paper (9) we described results for ablation study on the validation set. Here we report results for the

Model	$Acc_{ms}$	$Acc_{fs}$	$Acc_H$
Most Common Class*	0.7, 0.7, 0.7, 0.7	0, 0, 0, 0	0, 0, 0, 0
LDAM [7]*	71.5, 71.9, 71.6, 71.5	1.2, 5.9, 24.1, 41.2	2.4, 10.9, 36.0, 52.2
REVISE [43]	-	-	36.3, 41.1, 44.6, 50.9
CA-VAE [38]	58.2, 57.6, 60.0, 62.2	44.8, 51.6, 59.4, 62.3	50.6, 54.4, 59.6, 62.2
DA-VAE [38]	50.6, 56.0, 56.8, 56.8	47.9, 53.2, 61.0, 65.4	49.2, 54.6, 58.8, 60.8
CADA-VAE [38]	59.6, 60.9, 62.3, 63.1	51.4, 57.5, 63.6, 68.8	55.2, 59.2, 63.0, 64.9
CE Loss* (VE)	72.7, 72.9, 72.7, 72.0	0.6, 3.7, 19.1, 38.6	1.2, 6.9, 30.2, 50.2
LAGO [4]* (SE)	69.2, 69.0, 69.0, 68.1	13.8, 21.9, 38.1, 51.5	23.0, 33.2, 49.0, 58.6
<b>DRAGON (ours)</b>	58.0, 62.9, 63.3, 66.1	52.8, 55.9, 63.8, 69.6	<b>55.3, 59.2, 63.5, 67.8</b>

(a) Two-Level CUB

Model	$Acc_{ms}$	$Acc_{fs}$	$Acc_H$
Most Common Class*	0.2, 0.2, 0.2, 0.2	0, 0, 0, 0	0, 0, 0, 0
LDAM [7]*	43.7, 44.0, 44.2, 44.3	2.2, 6.6, 19.0, 31.8	4.3, 11.5, 26.6, 37.0
REVISE [43]	-	-	27.4, 33.4, 37.4, 40.8
CA-VAE [38]	35.8, 37.5, 37.5, 39.0	40.0, 46.5, 53.8, 55.7	37.8, 41.4, 44.2, 45.8
DA-VAE [38]	34.8, 37.3, 38.6, 38.2	41.4, 45.1, 50.2, 54.8	37.8, 40.8, 43.6, 45.1
CADA-VAE [38]	37.6, 38.2, 39.4, 41.9	44.1, 49.0, 55.3, 55.1	40.6, 43.0, 46.0, 47.6
CE Loss* (VE)	46.3, 46.3, 46.2, 45.6	0.9, 4.9, 17.2, 33.0	1.8, 8.9, 25.1, 38.3
LAGO [4]* (SE)	30.6, 30.4, 30.7, 31.0	14.4, 18.7, 21.9, 25.2	19.6, 23.2, 25.6, 27.8
<b>DRAGON (ours)</b>	37.2, 39.2, 40.5, 41.6	45.5, 49.6, 55.1, 57.2	<b>41.0, 43.8, 46.7, 48.2</b>

(b) Two-Level SUN

Model	$Acc_{ms}$	$Acc_{fs}$	$Acc_H$
Most Common Class*	2.5, 2.5, 2.5, 2.5	0, 0, 0, 0	0, 0, 0, 0
LDAM [7]*	90.7, 90.7, 90.5, 90.5	6.6, 14.4, 26.6, 41.6	12.4, 24.8, 41.1, 57.0
REVISE [43]	-	-	56.1, 60.3, 64.1, 67.8
CA-VAE [38]	73.4, 77.7, 81.0, 81.0	56.8, 66.0, 72.8, 77.1	64.0, 71.3, 76.6, 79.0
DA-VAE [38]	74.0, 74.6, 73.5, 73.9	63.0, 71.4, 77.7, 79.8	68.0, 73.0, 75.6, 76.8
CADA-VAE [38]	76.6, 79.4, 81.9, 82.6	63.8, 68.7, 74.8, 78.0	<b>69.6, 73.7, 78.2, 80.2</b>
CE Loss* (VE)	90.7, 90.9, 89.7, 87.9	5.9, 11.2, 32.6, 57.9	
LAGO [4]* (SE)	82.6, 81.9, 81.7, 81.5	11.5, 20.6, 46.2, 59.4	20.2, 33.0, 59.0, 68.7
<b>DRAGON (ours)</b>	74.5, 76.7, 79.2, 81.7	61.1, 62.9, 74.3, 82.1	67.1, 69.1, 76.7, <b>81.9</b>

(c) Two-Level AWA

Table S 15: Comparing DRAGON with SoTA GFSL models and baselines with increasing number of few-shot training samples on the CUB, SUN and AWA datasets. We report per-class  $Acc_{ms}$ ,  $Acc_{fs}$  and  $Acc_H$ . Each cell represents 1-shot, 2-shot, 5-shot and 10-shot accuracies

same model variant on the test set.

Tables S12 and S13 show the results of the ablation study on the test set. It shows the same behavior as the ablation study on the validation set that was reported in the main paper. Table S14 compares three different training protocols: (1) *All-Train*: Training the DRAGON fusion-module naively without a hold-out set. (2) *End-To-End*: Training all the architecture (both experts and fusion-module) end to end in an early fusion manner. (3) *Three-Stage-Training*: Training

our models as explained in section D.