

HIERARCHICAL CLASSIFICATION BY TRAINING TO DIFFUSE ON THE MANIFOLD

Anonymous authors

Paper under double-blind review

ABSTRACT

1 Hierarchical classification, the problem of classifying images according to a hi-
 2 erarchical taxonomy, has practical significance owing to the principle of “making
 3 better mistakes”, i.e., better to predict correct coarse labels than incorrect fine
 4 labels. Nevertheless, the literature does not sufficiently study this problem, pre-
 5 sumably because using top-1 accuracy to benchmark methods tends to yield a
 6 ranking order consistent with those using hierarchical metrics. On the other hand,
 7 for a downstream task of classification, today’s *de facto* practice is to *finetune* a
 8 pretrained deep neural network using the cross-entropy loss on leaf classes, result-
 9 ing in a leaf-class softmax classifier which even rivals sophisticated hierarchical
 10 classifiers atop deep nets. We argue that hierarchical classification should be bet-
 11 ter addressed by regularizing finetuning with explicit consideration of the given
 12 hierarchical taxonomy, because data intuitively lies in hierarchical manifolds in
 13 the raw feature space defined by the pre-trained model. To this end, we propose
 14 a hierarchical cross-modal contrastive loss that computes contrastive losses w.r.t
 15 labels at hierarchical levels in the taxonomy (including both hierarchy and text
 16 concepts). This results into features that can better serve hierarchical classifica-
 17 tion. Moreover, for inference, we re-conceptualize hierarchical classification by
 18 treating the taxonomy as a graph, presenting a diffusion-based methodology that
 19 adjusts posteriors at multiple hierarchical levels altogether. This distinguishes our
 20 method from the existing ones, which are either top-down (using coarse-class pre-
 21 dictions to adjust fine-class predictions) or bottom-up (processing fine-class pre-
 22 dictions towards coarse-label predictions). We evaluate our method by comparing
 23 them against existing ones on two large-scale datasets, iNat18 and iNat21. Ex-
 24 tensive experiments demonstrate that our method resoundingly outperforms prior
 25 arts w.r.t both top-1 accuracy and hierarchical metrics.

26 1 INTRODUCTION

27 Hierarchical classification (Naumoff, 2011; Deng et al., 2012; Zhu & Bain, 2017; Bertinetto et al.,
 28 2020) has long been a pivotal and challenging problem in machine learning. It aims to categorize
 29 images w.r.t a given hierarchical taxonomy, adhering to the principle of “making better mistakes” —
 30 essentially, favouring correct coarse-class predictions over inaccurate fine-class predictions (Deng
 31 et al., 2012; Wu et al., 2020).

32 Methods of hierarchical classification improve either training or inference. Existing inference meth-
 33 ods can be divided into two types: top-down (Redmon & Farhadi, 2017), and bottom-up (Val-
 34 madre, 2022). Top-down methods adjust the posterior for predicting a specific class by using its
 35 parent/ancestor posterior probabilities. They often underperform bottom-up methods Redmon &
 36 Farhadi (2017); Bertinetto et al. (2020), which prioritise predicting the leaf-classes and subsequently
 37 calculate posteriors for the parent/ancestor classes. Valmadre (2022) attributes the underperform-
 38 ance of top-down methods to the high diversity within coarse-level categories, soliciting effec-
 39 tive training methods. Perhaps surprisingly, although these sophisticated hierarchical classification
 40 methods show promising results in certain metrics, they do not consistently rival the simplistic flat-
 41 softmax baseline, which learns a softmax classifier on the leaf classes only. The status quo leads to
 42 a natural question: *Is it still helpful to make predictions for hierarchical classes other than the leaf*
 43 *classes for better hierarchical classification?* That said, it is still an open question how to effectively
 44 exploit hierarchical taxonomy to improve training and inference for hierarchical classification.

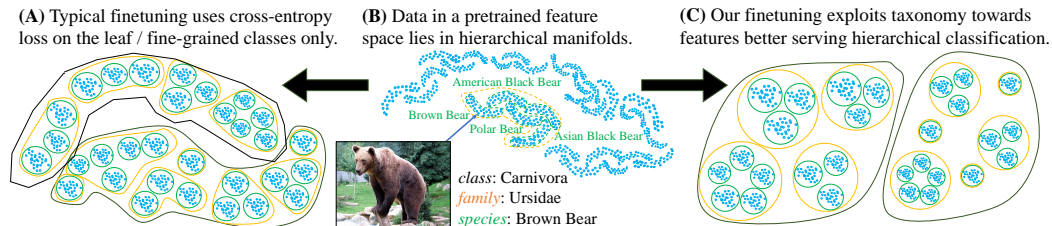


Figure 1: To solve a downstream task of classification, a *de facto* practice is to fine-tune a pretrained model using the cross-entropy loss on leaf classes (e.g., Brown Bear at the species level). (A): This yields features that help leaf-class classification but fail to model their hierarchical relationships w.r.t a taxonomy (e.g., Ursidae at the family level). Hence, it does not necessarily help hierarchical classification. Nevertheless, such features are better than the “raw features” of the pretrained model, which provides a feature space (B) where data hypothetically lie in hierarchical manifolds w.r.t the taxonomy. (C): Differently, we propose to fine-tune the pretrained model by *explicitly* exploiting the hierarchical taxonomy towards features that can better serve the task of hierarchical classification (Fig. 2).

45 We argue that, to better approach hierarchical classification for a downstream task that defines a
 46 hierarchical taxonomy, one should first explicitly exploit this taxonomy to learn features (Fig. 1),
 47 not just finetuning a pretrained model using the cross-entropy loss on leaf classes only (Bertinetto
 48 et al., 2020). Note that a taxonomy contains not only a hierarchy of concepts (e.g., species, genus,
 49 order, family, etc.) but also describable texts or names for the concepts. This motivates us to finetune
 50 a pretrained vision-language model (VLM) (Radford et al., 2021; Wang et al., 2023; Goyal et al.,
 51 2023). For better finetuning, we introduce a hierarchical cross-modal contrastive fine-tuning strategy
 52 (HCCF) (Sec. 3.2). HCCF explicitly exploits hierarchical taxonomy towards learning better features,
 53 which directly mirror the given taxonomy and hence better serve hierarchical classification.

54 Moreover, we argue that one should also collectively adjust posteriors at multiple hierarchical levels
 55 towards the final results of hierarchical classification. To this end, we present a set of diffusion-based
 56 methods for inference (Sec. 3.3), inspired by the literature of information retrieval Page et al. (1998);
 57 Iscen et al. (2017); An et al. (2021) which shows that diffusion is adept at mapping manifolds.
 58 This distinguishes our methods from existing top-down and bottom-up inference approaches that
 59 linearly interpret hierarchical classification. Our methods treat the hierarchical taxonomy as a graph,
 60 enabling probability distribution in the taxonomy. To the best of our knowledge, our work makes
 61 the first attempt to apply diffusion to hierarchical classification. Extensive experiments demonstrate
 62 that our diffusion-based inference methods, along with HCCF, achieve state-of-the-art performance
 63 and resoundingly outperform prior arts (Sec. 4.2).

64 To summarize, we make three major contributions.

- 65 1. We revisit the problem of hierarchical classification from the perspective of manifold learn-
 66 ing, offering new insights in the contemporary deep learning land.
- 67 2. We present the hierarchical cross-modal contrastive finetuning strategy for finetuning a
 68 model to better solve the problem of hierarchical classification.
- 69 3. We introduce a novel diffusion-based inference methodology to exploit posteriors at mul-
 70 tiple levels towards the final prediction.

71 2 RELATED WORKS

72 **Hierarchical classification.** Hierarchical classification holds significance, ensuring broader-level
 73 results even when detailed predictions are elusive. Datasets like ImageNet (Russakovsky et al.,
 74 2015) and WordNet (Miller, 1995) have long emphasized taxonomy, while newer ones like iNat18
 75 (Van Horn et al., 2018) and iNat21 offer finer-grained labels. Research in this domain is robust,
 76 with seminal works like “Hedging Your Bet” (Deng et al., 2012) and contemporary deep learning
 77 approaches employing flat softmax, oftmargin, and descendant softmax training losses (Valmadre,
 78 2022), along with bottom-up (Valmadre, 2022) and top-down (Redmon & Farhadi, 2017) infer-
 79 ences. Its practical applications are evident in areas like long-tailed 3D detection for autonomous
 80 driving (Peri et al., 2023), emphasizing specific metrics, methods, and joint training. Despite ex-

81 tensive research, recent findings suggest that advanced training and inference methods don't always
 82 surpass the flat softmax baseline (Valmadre, 2022). This paper presents innovative techniques that
 83 harness hierarchical data more efficiently during both the training and inference stages.

84 **Long-tailed recognition (LTR).** Long-tail categorization is an active research topic, as the long-
 85 tail feature is prevalent across coarse-level, fine-grained, and instance-level categorizations. Cur-
 86 rent strategies often employ data rebalancing (Mahajan et al., 2018; Chawla et al., 2002) or class-
 87 balanced loss functions (Cao et al., 2019) to improve the classification accuracy of infrequent classes.
 88 Despite these advancements, the exploration of the long-tail attribute within hierarchical categoriza-
 89 tion remains less investigated, indicating a need for further research in this area.

90 **Fine-grained visual categorization (FGVC).** Fine-grained categorization, a task bridging coarse-
 91 level classification and instance-level classification, presents both significant value and substantial
 92 challenges (Akata et al., 2015; Yang et al., 2018). In cases where predicting the fine-grained level
 93 tag proves difficult, users often still prefer an accurate coarse-level result, highlighting the impor-
 94 tance of hierarchical research within the fine-grained classification (Deng et al., 2012). This paper
 95 contributes to this aspect, pushing forward the understanding and application of hierarchical fine-
 96 grained categorization in the context of long-tail distributions.

97 **Diffusion.** Diffusion is an advanced methodology adept at faithfully delineating the manifold within
 98 a data distribution by leveraging the interconnectedness inherent in a Markov chain (Zhou et al.,
 99 2003a;b). A renowned variation of this method, PageRank (Page et al., 1998), has achieved con-
 100 siderable success in various business endeavors. Moreover, it has been extensively employed in
 101 the realm of image retrieval (Isken et al., 2017; An et al., 2021), an application of instance-level
 102 classification. However, its potential in broader classifications, such as fine-grained and hierarchi-
 103 cal categorizations, has not been extensively researched. In this paper, we pioneer the exploration
 104 of its utility in understanding and utilizing the relationships within these broader, fine-grained, and
 105 hierarchical classifications.

106 3 METHODS

107 **Hierarchical classification and notations.** This paper delves into the intricacies of Single-Path
 108 Labels (SPL) and Non-Mandatory Leaf-Node Prediction (NMLNP) in hierarchical classification.
 109 In SPL, a sample is restricted from belonging to multiple distinct classes unless there exists a
 110 superclass-subclass relationship. On the other hand, NMLNP allows the classifier to predict any
 111 class within the hierarchy, not being confined to just the leaf nodes. In this study, we let Y denote
 112 the entirety of categories within the taxonomy tree. For a given node $y \in Y$, $C(y)$ signifies its child
 113 nodes, while $A(y)$ stands for its ancestor nodes. The set of leaf nodes is represented by L .

114 3.1 HIERARCHICAL MANIFOLD

115 We introduce a hierarchical manifold model in the embedding space to elucidate the intricacies of
 116 hierarchical classification. Although data manifolds are prevalent in high-dimensional spaces, what
 117 sets hierarchical classification apart is its distinct manifold structure. As depicted in Fig 1, before
 118 optimization, each category in the embedding space can be visualized as a separate manifold. Draw-
 119 ing an analogy to the parent-child node relationship, **parent manifolds envelop child manifolds.**
 120 An optimally refined embedding space should discern manifolds across all hierarchical levels.

121 The hierarchical manifold assumption holds merit. Given that manifolds are frequently observed
 122 in diverse real-world datasets, it's plausible that the embedding space houses these hierarchical
 123 manifolds prior to achieving an optimal training solution. This sheds light on the limitations of
 124 current techniques in addressing the hierarchical classification challenge. As illustrated in Fig. 1,
 125 existing methods, **failing to grasp the nuances of higher-level manifolds**, might misclassify an
 126 image under the family level, even if they correctly identify it at the species level.

127 While there are extant hierarchical loss functions aimed at this problem, they predominantly predict
 128 only the leaf node categories. Consequently, the hierarchical loss equation ultimately converges to
 129 supervision solely at the leaf level. For instance, when employing bottom-up inference for interior

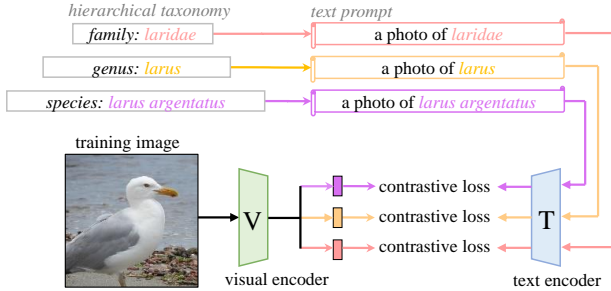


Figure 2: The proposed Hierarchical Cross-modal Contrastive Finetuning (HCCF) exploits hierarchical taxonomy to adapt a pretrained visual encoder to the downstream task of hierarchical classification. It sums contrastive losses between a training image and its taxonomic names at multiple levels. To the best of our knowledge, we make the first attempt to fine-tune a vision-language model using a predefined taxonomy for hierarchical classification.

130 node prediction results as:

$$q_y(\theta) = \begin{cases} [\text{softmax}_L(\theta)]_y & \text{if } y \in L \\ \sum_{v \in C(y)} q_v(\theta) & \text{if } y \notin L \end{cases} \quad (1)$$

131 The negative log-likelihood concerning the interior nodes **is reduced to the leaf nodes** as $\ell(y, \theta) =$
 132 $-\log q_y(\theta) = -\log \left(\sum_{u \in L(y)} \exp \theta_u \right) + \log \left(\sum_{u \in L(y)} \exp \theta_u \right)$. Advanced losses, such as soft-
 133 margin and descendant softmax (Valmadre, 2022), also focus on the leaf level, neglecting the sepa-
 134 ration of higher-level manifolds. This results in suboptimal outcomes for hierarchical classification.

135 The hierarchical manifold model inspires novel strategies for both training and inference. For the
 136 training phase, the model suggests that we should: 1) **Effectively leverage the multiple labels** asso-
 137 ciated with each training image, and 2) **Enhance the separation** between sample distributions from
 138 different categories across various levels in the embedding space, thereby reducing misclassification
 139 risks. During inference, the model motivates us to use **diffusion—a technique renowned for its**
 140 **efficacy with manifolds**—to refine the scores predicted by the neural network.

141 3.2 HIERARCHICAL CROSS-MODAL CONTRASTIVE FINE-TUNING

142 To more **effectively map the taxonomy relations in the embedding space**, we initially employ
 143 the Vision-language pretrained model, CLIP (Radford et al., 2021), as our primary visual encoder.
 144 Using textual descriptions for each image provides a more comprehensive supervisory signal, cap-
 145 turing both leaf and interior node relationships in the taxonomy tree. While CLIP’s superiority over
 146 ImageNet as a pretrained model is somewhat recognized, its efficacy in hierarchical classification
 147 remains untested. Our experiments on the renowned iNat18 dataset (Van Horn et al., 2018) indicate
 148 significant improvements (Table 1).

149 Our advancements extend beyond the utilization of the CLIP pre-trained model. We propose a hier-
 150 archical cross-modal contrastive loss, aiming to **extend the distance** between sample distributions
 151 across varied categories and levels (shown in Fig. 2). This strategy is anchored in two core tenets
 152 of our hierarchical manifold model. Firstly, we harness the full potential of textual descriptions for
 153 each training image. By employing the CLIP text encoder, we encode the hierarchical labels of
 154 these images. Distinct from prevailing hierarchical losses, our interior node prediction isn’t merely
 155 inferred from leaf nodes. Instead, it’s directly guided by the embedding vectors of text labels across
 156 different levels, enabling a more nuanced understanding of category relationships and better cap-
 157 turing of higher-level manifolds. Secondly, our methodology employs contrastive loss, ensuring
 158 maximal separation between samples from diverse categories, thereby mitigating the complexities
 159 introduced by hierarchical manifolds. Our hierarchical cross-modal contrastive fine-tuning loss is
 160 defined as:

$$L(f, g) := \sum_{l=1}^L \left(\sum_{i=1}^N -\log \frac{\exp(\bar{f}^l(I_i) \cdot \bar{g}(T_i^l))}{\sum_{j=1}^N \exp(\bar{f}^l(I_i) \cdot \bar{g}(T_j^l))} + \sum_{i=1}^N -\log \frac{\exp(\bar{f}^l(I_i) \cdot \bar{g}(T_i^l))}{\sum_{j=1}^N \exp(\bar{f}^l(I_j) \cdot \bar{g}(T_i^l))} \right), \quad (2)$$

161 where $\bar{f}^l(I_i)$ is normalized embedding of the i -th image I_i from the visual encoder f^l , which con-
 162 sists of visual backbone and level-specific head. $\bar{g}(T_j^l)$ is the normalized text embedding of the
 163 text T_j^l , that is the j -th sample of level l extracted from text encoder g . Assuming there are N
 164 image-texts pairs in one batch, I_i is input image and T_i^l denotes the ground truth label at level l . All

165 taxonomy tree text prompts utilize a shared text encoder, mitigating overfitting risks and conserving
 166 training and inference resources. The visual encoder comprises a shallow feature extractor and a
 167 level-specific extractor head for every level, ensuring encoding aligns with the hierarchical taxon-
 168 omy level. Both visual and text encoders are updated during training, and text encoding of every
 169 taxonomy level serve as linear classifier weights during inference.

170 3.3 DIFFUSION-BASED INFERENCE

171 Through our new training strategy, we generate prediction scores for all taxonomy categories. The
 172 ensuing challenge is to utilize these scores for inference and robust decision-making effectively.

173 **Existing inference techniques**, namely the top-down (Redmon & Farhadi, 2017) and down-
 174 top (Valmadre, 2022) approaches, can be further improved. The top-down method computes the
 175 conditional likelihood of each child node based on its parent nodes. While theoretically appeal-
 176 ing, it is empirically outperformed by the down-top approach (Redmon & Farhadi, 2017; Bertinetto
 177 et al., 2020; Valmadre, 2022). Valmadre (Valmadre, 2022) attributes this underperformance to the
 178 high diversity within coarse-level categories and advocates using fine-grained scores to infer hierar-
 179 chical outcomes. We align with Valmadre’s observations, yet we assert that predictions for mid-level
 180 categories have inherent value when utilizing our innovative diffusion-based inference.

181 **Motivation.** When a category receives an anomalously high or low score from the neural network,
 182 we can recalibrate this score based on the scores of its neighboring categories within the taxonomy
 183 tree. Essentially, sub-categories under the same parent category should exhibit consistent scoring
 184 patterns, either high or low. By diffusing the scores across the taxonomy’s structural connections to
 185 achieve equilibrium, we can enhance the initial predictions made by the neural network Remarkably,
 186 experiment results show that our method enhances both the leaf-level top-1 accuracy and the overall
 187 hierarchical performance, outperforming existing techniques (Sec. 4.4).

188 **Notation.** Given a total of n categories (including intermediate categories) in the taxonomy graph,
 189 we define a connection matrix $W \in R^{n \times n}$ to describe the interrelationships among categories within
 190 the graph. Let $f^0 \in R^n$ be the prediction output of the neural network. Our target is to refine f^0
 191 based on W to get the final f^* , which gives both better leaf-level and hierarchical performance.

192 **Connection matrix.** We first use the expert-designed taxonomy given by each dataset to define the
 193 connection matrix W . That is, $w_{ij} = 1$ if category i and j have the parent-children relation in the
 194 taxonomy tree. Otherwise $w_{ij} = 0$. Here, we assume the graph is undirected, and the connection
 195 matrix is symmetric ($W = W^T$). The self-similarity is set as zero ($\text{diag}(W) = \mathbf{0}$). We will explore
 196 more weight options within this matrix in subsequent sections.

197 Normalization for the connection matrix is an essential step for diffusion in information retrieval.
 198 We find it is also necessary in the hierarchical classification. In this paper, we use the symmetrically
 199 normalization as follows:

$$W_n = D^{-1/2} W D^{-1/2}, \quad D = \text{diag}(W \mathbf{1}_n). \quad (3)$$

200 **Iteration.** Our diffusion mechanism iteratively updates the category scores according to the follow-
 201 ing:

$$f^{t+1} = \alpha W_n f^t + (1 - \alpha) f^0, \quad (4)$$

202 where α is set among $(0, 1)$. This is a “random walk” algorithm in the taxonomy graph. Intu-
 203 itively, for each iteration, each category spreads its prediction score to its neighbor categories with
 204 probability α , and follows the initial neural network prediction with probability $1 - \alpha$.

205 **Convergence:** The iterative process is assured to converge towards a stationary distribution (Zhou
 206 et al., 2003b). We provide a straightforward proof here. By recursively integrating $f^1 = \alpha W_n f^0 +$
 207 $(1 - \alpha) f^0$ into subsequent iterations f^2, f^3 , and so on, we derive:

$$f^t = (\alpha W_n)^t f^0 + (1 - \alpha) \sum_{i=0}^{t-1} (\alpha W_n)^i f^0. \quad (5)$$

208 As t approaches infinity, the term $(\alpha W_n)^t$ approaches zero, and the summation term converges
 209 to $(I - \alpha W_n)^{-1}$, where I denotes the identity matrix of size n . Thus, the eventual stationary
 210 distribution is expressed as:

$$f^* = (1 - \alpha)(I - \alpha W_n)^{-1} f^0. \quad (6)$$

211 **Relation to Spectral Clustering:** It’s pertinent to elucidate the connection between our hierarchi-
 212 cal classification diffusion and spectral clustering, given that both methodologies emphasize node
 213 grouping within a graph. Notably, the term $(I - \alpha W_n)$ in Equation 6 can be interpreted as a variant
 214 of the symmetrically normalized Laplacian $(I - W_n)$ for the taxonomy graph. This Laplacian is
 215 instrumental in spectral clustering, enabling the capture of the data’s intrinsic topological character-
 216 istics. In the spectral clustering paradigm, each node is characterized by a k -dimensional spectral
 217 space vector, derived from the k eigenvectors satisfying $v = (I - W_n)^{-1} \lambda v$. Conversely, our dif-
 218 fusion process assigns each node a singular scalar score, as dictated by Equation 6. Conceptually,
 219 our diffusion approach can be perceived as a tailored spectral clustering for the neural network’s
 220 predicted vector f^0 , pinpointing a category subset with peak scores in the spectral domain.

221 **Differentiable diffusion:** As demonstrated in Eq. 6, the diffusion process converges to a closed
 222 form. Intriguingly, this represents a linear transformation from the initial scores f^0 to the final state
 223 f^* . Currently, the connection matrix W is constructed based on the provided taxonomy tree struc-
 224 ture, comprising binary values that might not accurately capture the genuine relationships between
 225 category pairs. Given a substantial sample size from the training set, we investigate the potential of
 226 training a linear mapping directly to supplant the closed form. This differentiable method could offer
 227 a more nuanced understanding of the relationships between categories. We call this new approach
 228 differentiable diffusion.

229 Our main contribution lies in introducing **an advanced diffusion method**, specially designed for
 230 using the taxonomy graph’s structure. To the best of our knowledge, this is the first work to apply
 231 diffusion techniques to hierarchical classification problems. While existing literature has extensively
 232 explored the diffusion of instance space (like web and image) with considerable success (Page et al.,
 233 1998; An et al., 2021), the impact of diffusion on the category space (how to group the instances)
 234 remains largely uncharted territory. This diffusion approach offers several distinct advantages over
 235 existing top-down and down-top inference:

- 236 1. **Comprehensive graph utilization:** Unlike traditional methods that focus solely on di-
 237 rect parent-child relationships, our diffusion technique leverages the entire graph structure,
 238 including sibling relationships.
- 239 2. **Iterative information blending:** While existing methods transfer information once
 240 through the graph edge, our diffusion process iteratively blends information at each node
 241 until a stable state is achieved, thereby maximizing the utility of all predicted category
 242 nodes.
- 243 3. **Manifold problem resolution:** Our method addresses the manifold problem by utilizing
 244 inter-category relationships, on which we elaborate subsequently.

245 4 EXPERIMENTS

246 4.1 IMPLEMENTATIONS

247 To assess the efficacy of our novel training and inference approach for hierarchical classification,
 248 we employ the metrics and dataset from the recent study by Valmadre (Valmadre, 2022). This study
 249 presents state-of-the-art (SOTA) methods, comprehensive experiments on existing techniques, and
 250 a suite of robust metrics tailored for hierarchical classification. Similar to Valmadre (2022), all
 251 the experiments use ResNet 50 (He et al., 2016) as the backbone. Valmadre’s benchmark dataset
 252 is the balanced iNaturalist 21-mini (iNat21). In our evaluation, we extend the datasets to include
 253 iNaturalist 18 (iNat18), showcasing the versatility of our method and its performance under long-
 254 tailed distributions. In line with Valmadre’s approach (Valmadre, 2022), our metrics are derived
 255 from operating curves, encompassing Average Precision (AP), Average Correct (AC), Recall at X%
 256 Correct (R@XC), and a specificity measure. We also incorporate single prediction metrics such as
 257 Majority F1, Leaf F1, and Leaf Top1 Accuracy. Notably, while Leaf Top1 Accuracy gauges leaf-
 258 level accuracy, the other metrics focus on hierarchical classification performance. Our methods are
 259 benchmarked against various SOTA hierarchical classification techniques, including flat softmax
 260 (Bertinetto et al., 2020), Multilabel focal (Lin et al., 2017), Cond softmax (Redmon & Farhadi,
 261 2017), Cond sigmoid (Brust & Denzler, 2019), DeepRTC (Wu et al., 2020), PS softmax (Wu et al.,
 262 2020), Softmargin and (Valmadre, 2022) descendent softmax (Valmadre, 2022).

Table 1: Benchmarking results on the iNat18 dataset. We report numbers w.r.t both hierarchical metrics (Valmadre, 2022) and the standard top-1 accuracy on leaf classes (dubbed Leaf Top1 in the last column). Our HCCF, which contrastively finetunes a pretrained model using all the taxonomic levels, significantly outperforms prior arts. Additionally applying diffusion improves performance notably further.

Model	AP	AC	R@90C	R@95C	Majority F1	Leaf F1	Leaf Top1
Flat softmax (Bertinetto et al., 2020)	61.18	58.94	45.44	37.58	64.27	64.57	47.33
Multilabel focal (Lin et al., 2017)	46.70	43.97	34.05	28.05	50.69	49.91	14.85
Cond softmax (Redmon & Farhadi, 2017)	54.13	51.12	36.68	30.07	58.74	58.60	36.94
Cond sigmoid (Brust & Denzler, 2019)	52.04	49.29	35.23	29.31	55.46	58.29	36.36
Deep RTC (Wu et al., 2020)	60.07	54.25	23.69	14.33	66.72	66.72	47.13
PS softmax (Wu et al., 2020)	64.15	62.02	49.54	42.02	67.50	67.44	49.21
Softmargin (Valmadre, 2022)	58.53	55.86	40.28	33.71	58.73	63.70	45.10
Descendant softmax (Valmadre, 2022)	61.88	59.65	46.79	38.49	65.48	65.32	48.71
HCCF	<u>72.75</u>	<u>70.60</u>	<u>59.56</u>	<u>52.60</u>	<u>72.73</u>	<u>75.16</u>	<u>55.78</u>
HCCF + diffusion	73.48	71.88	62.48	55.53	75.94	75.71	56.33

Table 2: Benchmarking results on the iNat21 dataset. We report numbers w.r.t both hierarchical metrics (Valmadre, 2022) and the standard top-1 accuracy on leaf classes (dubbed Leaf Top1). Our HCCF finetuning, which contrastively finetunes a pretrained model using all the taxonomic levels, significantly outperforms prior arts. Additionally applying diffusion to inference improves performance notably further.

Model	AP	AC	R@90C	R@95C	Majority F1	Leaf F1	Leaf Top1
Flat softmax (Bertinetto et al., 2020)	66.17	64.32	53.85	47.02	68.87	68.69	50.89
Multilabel focal (Lin et al., 2017)	54.58	50.35	36.16	30.45	50.62	60.27	31.05
Cond softmax (Redmon & Farhadi, 2017)	58.88	56.26	42.95	36.23	62.85	62.80	41.64
Cond sigmoid (Brust & Denzler, 2019)	59.24	56.74	42.84	35.61	61.41	65.11	44.64
Deep RTC (Wu et al., 2020)	63.92	58.07	25.36	14.10	70.17	70.22	51.43
PS softmax (Wu et al., 2020)	68.22	66.49	56.20	49.85	71.07	70.80	52.76
descendant softmax (Valmadre, 2022)	64.95	62.71	48.84	42.59	64.64	69.03	50.55
softmargin (Valmadre, 2022)	66.53	64.72	54.41	47.91	69.39	69.09	52.22
HCCF (Ours)	<u>72.46</u>	<u>70.52</u>	<u>60.49</u>	<u>53.66</u>	<u>73.35</u>	<u>74.72</u>	<u>55.11</u>
HCCF + diffusion (Ours)	73.16	71.62	62.81	55.97	75.31	75.32	55.86

263 4.2 COMPARE WITH STATE-OF-THE-ART METHODS

264 We performed a comparative analysis of our innovative training and inference methods against estab-
 265 lished state-of-the-art (SOTA) hierarchical techniques using the iNat18 and iNat21 datasets. Table 1
 266 and Table 2 demonstrate the enhanced performance of our approach over existing SOTA methods
 267 across both datasets. Unless otherwise indicated, all methods utilized the same CLIP pretrained
 268 model.

269 In our implementation of the SOTA methods, we strictly adhered to the code provided by Val-
 270 madre (Valmadre, 2022). Our results for the iNat21 are consistent with those presented by Valmadre.
 271 Although (Valmadre, 2022) did not provide outcomes for iNat18, we included results for this dataset
 272 to illustrate our model’s capability in handling long-tailed distributions, noting that iNat18 is long-
 273 tailed while iNat21 is balanced.

274 While a direct comparison with (Valmadre, 2022) for iNat18 is not available, we ensured the reli-
 275 ability of our results by using the reproduction code and settings from Valmadre’s open-source re-
 276 sources. These results emphasize the advantages of our method over existing SOTA methodologies,
 277 proving effective for both balanced and long-tailed datasets. Our hierarchical contrastive training
 278 approach sets new standards in the field, outperforming existing SOTA methods for both the iNat18
 279 and iNat21 datasets.

280 4.3 ABLATION STUDY ABOUT HIERARCHICAL CROSS-MODAL CONTRASTIVE FINE-TUNING

281 We conducted an ablation study to assess the impact of each component in hierarchical cross-modal
 282 contrastive fine-tuning (HCCF), as detailed in Table 3. In contrast to the traditional training using
 283 cross-entropy loss (flat softmax (Bertinetto et al., 2020) combined with negative log-likelihood), our
 284 HCCF incorporates several enhancements:

Table 3: Ablation Study of Hierarchical Cross-Modal Fine-Tuning (HCCF) on iNat18. This study highlights three key modifications from the Cross-Entropy (CE) loss baseline to our HCCF: using CLIP pre-trained text encoder (text embedding), hierarchical training (L67 and L123456), and contrastive loss (CL). The adoption of the CLIP pre-trained text encoder markedly boosts model performance, with hierarchical training and contrastive loss providing additional enhancements. For a comprehensive explanation, refer to Sec. 4.3.

Models	AP	AC	R@90C	R@95C	Majority F1	Leaf F1	Leaf Top1
CE loss baseline (Bertinetto et al., 2020)	61.18	58.94	45.44	37.58	64.27	64.57	47.33
CE loss + text embedding	66.25	64.09	51.72	43.66	69.42	69.31	53.10
CE loss + text embedding + L67	67.81	65.7	54.13	46.09	70.81	70.66	54.07
CE loss + text embedding + L1234567	69.18	67.07	56.32	48.28	71.99	71.81	53.68
CL + text encoder + L1234567 (HCCF)	72.75	70.60	59.56	52.60	72.73	75.16	55.78

Table 4: An ablation study of Hierarchical Cross-Modal Fine-Tuning (HCCF) over different training levels on iNat18 reveals intriguing insights. While training across more levels consistently enhances all metrics under CE loss, as illustrated in Table 3, the same doesn’t hold true for contrastive loss. Training at the leaf level (denoted as L7) yields the highest leaf Top1 accuracy but falls short in hierarchical metrics compared to multi-level encoder head training. For metrics like AP, AC, and Leaf F1, comprehensive training across all levels (denoted as L123467) outperforms other configurations. Training on levels 6 and 7 alone achieves the peak for R@90C and R@95C. Broadening the training levels benefits hierarchical metrics, with the coarsest (level 1) and sub-finest (level 6) levels proving most advantageous. It’s noteworthy that these findings diverge from the prevailing belief that top-1 accuracy benchmarks align with hierarchical metric rankings (Russakovsky et al., 2015), underscoring the importance of studying hierarchical metrics.

Model	AP	AC	R@90C	R@95C	Majority F1	Leaf F1	Leaf Top1
HCCF L7	72.40	70.33	59.36	52.42	72.33	74.72	56.69
HCCF L67	72.64	70.65	60.53	53.22	72.85	74.88	56.10
HCCF L567	72.62	70.51	59.69	52.92	72.72	74.97	55.80
HCCF L4567	72.50	70.34	59.26	52.29	72.58	74.89	55.43
HCCF L34567	72.52	70.36	59.46	52.27	72.65	74.87	55.29
HCCF L234567	72.55	70.37	59.38	51.63	72.45	74.98	55.72
HCCF L1234567	72.75	70.60	59.56	52.60	72.73	75.16	55.78

285 **Use of CLIP pre-trained text encoder:** To assess the benefits of the CLIP pre-trained text encoder,
 286 we modified the initial weights of the final fully connected layer in CE loss training by incorporat-
 287 ing the CLIP pre-trained text embeddings for each category. This strategy harnesses the knowledge
 288 from the cross-modal pre-training set, creating a more optimized initial embedding space for the
 289 categories. This straightforward adjustment leads to a marked improvement in the CE baseline
 290 performance. While the effectiveness of leveraging the CLIP pre-trained encoder has been previ-
 291 ously noted in contexts like few-shot classification (Xiao et al., 2022) and object detection (Jin et al.,
 292 2021), our work stands out as the first to apply this technique to hierarchical classification, achieving
 293 notable gains.

294 **Hierarchical training:** Unlike the flat softmax which aggregates the probabilities of child nodes
 295 to determine the mid-level node probability, our hierarchical training instructs the model to directly
 296 estimate the probability for each mid-level node. This strategy aims to better delineate the mid-
 297 level manifolds, as depicted in Fig. 1. This method further enhances performance, particularly in
 298 hierarchical metrics.

299 **Incorporation of contrastive loss:** As discussed in Sec. 3.2, the addition of the contrastive loss
 300 further augments the model’s performance.

301 In summary, our HCCF approach, with its multiple enhancements, demonstrates superior perfor-
 302 mance compared to traditional training methods. We additionally performed hierarchical cross-
 303 modal fine-tuning at various levels, beginning exclusively with the leaf level and culminating with
 304 all levels. As indicated in Table 4, the utilization of all levels yielded the optimal hierarchical perfor-
 305 mance. However, it adversely affected the leaf-level performance. Harnessing the bottom two levels
 306 proved to be the most cost-efficient strategy. Intriguingly, incorporating additional levels, such as
 307 levels 5, 6, and 7, did not improve performance compared to just using levels 6 and 7. It’s notewor-
 308 thy that these findings diverge from the prevailing belief that top-1 accuracy benchmarks align with
 309 hierarchical metric rankings (Russakovsky et al., 2015), underscoring the importance of studying
 310 hierarchical metrics.

Table 5: Evaluation of our cutting-edge diffusion-based inference against established state-of-the-art (SOTA) methods on iNat18. Despite all inference techniques utilizing the same trained model, our diffusion and differentiable diffusion approaches surpass all the SOTA methods. Notably, this enhancement is achieved without any modifications to the trained model.

Model	AP	AC	R@90C	R@95C	Majority F1	Leaf F1	Leaf Top1
Top-down (Redmon & Farhadi, 2017)	64.36	61.72	46.10	34.97	68.54	68.36	46.62
Advanced-top-down (Jain et al., 2023)	72.11	69.98	58.09	46.96	76.23	75.96	55.71
Bottom-up (Valmadre, 2022)	72.75	70.60	59.56	52.60	72.73	75.16	55.78
Diffusion (Ours)	73.48	71.88	62.48	55.53	75.94	75.71	56.33
Differentiable diffusion (Ours)	73.82	71.91	61.99	53.36	76.01	76.09	59.70

Table 6: Our diffusion-based inference method is model-agnostic, enhancing classifier performance across all metrics. This improvement is consistent whether the model is trained comprehensively across all levels (HCCF L123456), on level 6 and level 7 (HCCF L67), or solely at the leaf level (Flat softmax).

Model	AP	AC	R@90C	R@95C	Maj F1	Leaf F1	Leaf Top1
HCCF L1234567 bottom-up	72.75	70.60	59.56	52.60	72.73	75.16	55.78
HCCF L1234567 diffusion	73.60	71.85	62.06	54.97	74.79	75.82	56.50
HCCF L1234567 differentiable diffusion	73.82	71.91	61.99	53.36	76.01	76.09	59.70
HCCF L67 bottom-up	72.64	70.65	60.53	53.22	72.85	74.88	56.10
HCCF L67 diffusion	73.35	71.63	62.26	55.25	74.57	75.51	56.84
HCCF L67 differentiable diffusion	73.23	71.37	61.38	53.30	75.48	75.44	59.51
Flat softmax bottom-up	69.18	67.07	56.32	48.28	71.99	71.81	53.68
Flat softmax diffusion	69.45	67.56	56.47	48.61	72.57	72.31	54.14
Flat softmax differentiable diffusion	69.20	67.12	56.40	48.75	71.96	71.81	53.84

311 4.4 COMPARE DIFFUSION WITH OTHER INFERENCE METHODS

312 In addition to training, inference plays a pivotal role in hierarchical classification for final decision-
 313 making. We evaluated our innovative diffusion-based techniques, including both general and dif-
 314 ferentiable diffusion, against traditional top-down and bottom-up inference methods. The results,
 315 presented in Table 5, reveal that our methods notably surpass existing ones. Intriguingly, diffusion
 316 not only enhances hierarchical metrics but also boosts the leaf-level top 1 accuracy. The fact that
 317 our general diffusion doesn’t necessitate extra training makes this discovery particularly noteworthy.
 318 When trained using our differentiable diffusion, the performance escalates even further.

319 Differentiable diffusion excels in numerous metrics over general diffusion except in R@90C and
 320 R@95C. The advantage of general diffusion is its simplicity and the absence of a training require-
 321 ment. Further experiments, as seen in Table 6, confirm the consistency of these findings across
 322 various models. This underscores the novelty and success of our diffusion-centric approach to clas-
 323 sification.

324 4.5 SOCIAL IMPACT AND LIMITATIONS

325 Our research introduces innovative training methodologies and novel diffusion mechanisms for hi-
 326 erarchical classification. Extensive experiments show that our proposed methods deliver more ac-
 327 curate and impactful hierarchical classification results. These advancements have potential impli-
 328 cations for various applications, from object detection to the realm of autonomous driving. While
 329 our techniques represent a significant leap forward, they have limitations. Our empirical evaluations
 330 have been primarily anchored to the well-structured iNat18 and iNat21 datasets. As a next step, it
 331 would be pivotal to assess the versatility of our method in diverse real-world contexts, including its
 332 potential role in autonomous driving systems.

333 5 CONCLUSIONS

334 This paper introduces a fresh perspective on the hierarchical classification problem by viewing it
 335 through the lens of manifold learning. Leveraging this approach, we present innovative strategies
 336 for training and inference. Our proposed hierarchical cross-modal contrastive loss and graph-based
 337 diffusion methods for hierarchical predictions offer a nuanced balance between coarse and fine-
 338 class predictions. Evaluations on iNat18 and iNat21 datasets demonstrate the superior performance
 339 of our methods in terms of both top-1 accuracy and various hierarchical metrics, marking a notable
 340 advancement in the field of hierarchical classification.

341 REFERENCES

- 342 Zeynep Akata, Scott Reed, Daniel Walter, Honglak Lee, and Bernt Schiele. Evaluation of output
343 embeddings for fine-grained image classification. In *Proceedings of the IEEE conference on*
344 *computer vision and pattern recognition*, pp. 2927–2936, 2015.
- 345 Guoyuan An, Yuchi Huo, and Sung-Eui Yoon. Hypergraph propagation and community selection
346 for objects retrieval. *Advances in Neural Information Processing Systems*, 34:3596–3608, 2021.
- 347 Luca Bertinetto, Romain Mueller, Konstantinos Tertikas, Sina Samangooei, and Nicholas A Lord.
348 Making better mistakes: Leveraging class hierarchies with deep networks. In *Proceedings of the*
349 *IEEE/CVF conference on computer vision and pattern recognition*, pp. 12506–12515, 2020.
- 350 Clemens-Alexander Brust and Joachim Denzler. Integrating domain knowledge: using hierarchies
351 to improve deep classifiers. In *Asian conference on pattern recognition*, pp. 3–16. Springer, 2019.
- 352 Kaidi Cao, Colin Wei, Adrien Gaidon, Nikos Arechiga, and Tengyu Ma. Learning imbalanced
353 datasets with label-distribution-aware margin loss. *Advances in neural information processing*
354 *systems*, 32, 2019.
- 355 Nitesh V Chawla, Kevin W Bowyer, Lawrence O Hall, and W Philip Kegelmeyer. Smote: synthetic
356 minority over-sampling technique. *Journal of artificial intelligence research*, 16:321–357, 2002.
- 357 Jia Deng, Jonathan Krause, Alexander C Berg, and Li Fei-Fei. Hedging your bets: Optimizing
358 accuracy-specificity trade-offs in large scale visual recognition. In *2012 IEEE Conference on*
359 *Computer Vision and Pattern Recognition*, pp. 3450–3457. IEEE, 2012.
- 360 Sachin Goyal, Ananya Kumar, Sankalp Garg, Zico Kolter, and Aditi Raghunathan. Finetune like
361 you pretrain: Improved finetuning of zero-shot vision models. In *Proceedings of the IEEE/CVF*
362 *Conference on Computer Vision and Pattern Recognition*, pp. 19338–19347, 2023.
- 363 Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recog-
364 nition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp.
365 770–778, 2016.
- 366 Ahmet Iscen, Giorgos Tolias, Yannis Avrithis, Teddy Furon, and Ondrej Chum. Efficient diffusion
367 on region manifolds: Recovering small objects with compact cnn representations. In *Proceedings*
368 *of the IEEE conference on computer vision and pattern recognition*, pp. 2077–2086, 2017.
- 369 Kanishk Jain, Shyamgopal Karthik, and Vineet Gandhi. Test-time amendment with a coarse classi-
370 fier for fine-grained classification. *arXiv preprint arXiv:2302.00368*, 2023.
- 371 Ying Jin, Yinpeng Chen, Lijuan Wang, Jianfeng Wang, Pei Yu, Zicheng Liu, and Jenq-Neng
372 Hwang. Is object detection necessary for human-object interaction recognition? *arXiv preprint*
373 *arXiv:2107.13083*, 2021.
- 374 Tsung-Yi Lin, Priya Goyal, Ross Girshick, Kaiming He, and Piotr Dollár. Focal loss for dense
375 object detection. In *Proceedings of the IEEE international conference on computer vision*, pp.
376 2980–2988, 2017.
- 377 Dhruv Mahajan, Ross Girshick, Vignesh Ramanathan, Kaiming He, Manohar Paluri, Yixuan Li,
378 Ashwin Bharambe, and Laurens Van Der Maaten. Exploring the limits of weakly supervised
379 pretraining. In *Proceedings of the European conference on computer vision (ECCV)*, pp. 181–
380 196, 2018.
- 381 George A Miller. Wordnet: a lexical database for english. *Communications of the ACM*, 38(11):
382 39–41, 1995.
- 383 DG Naumoff. Hierarchical classification of glycoside hydrolases. *Biochemistry (Moscow)*, 76:
384 622–635, 2011.
- 385 Lawrence Page, Sergey Brin, Rajeev Motwani, and Terry Winograd. The pagerank citation ranking:
386 Bring order to the web. Technical report, Technical report, stanford University, 1998.

- 387 Neehar Peri, Achal Dave, Deva Ramanan, and Shu Kong. Towards long-tailed 3d detection. In
388 *Conference on Robot Learning*, pp. 1904–1915. PMLR, 2023.
- 389 Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal,
390 Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual
391 models from natural language supervision. In *International conference on machine learning*, pp.
392 8748–8763. PMLR, 2021.
- 393 Joseph Redmon and Ali Farhadi. Yolo9000: better, faster, stronger. In *Proceedings of the IEEE*
394 *conference on computer vision and pattern recognition*, pp. 7263–7271, 2017.
- 395 Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng
396 Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, et al. Imagenet large scale visual
397 recognition challenge. *International journal of computer vision*, 115:211–252, 2015.
- 398 Jack Valmadre. Hierarchical classification at multiple operating points. *Advances in Neural Infor-*
399 *mation Processing Systems*, 35:18034–18045, 2022.
- 400 Grant Van Horn, Oisin Mac Aodha, Yang Song, Yin Cui, Chen Sun, Alex Shepard, Hartwig Adam,
401 Pietro Perona, and Serge Belongie. The inaturalist species classification and detection dataset. In
402 *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 8769–8778,
403 2018.
- 404 Yidong Wang, Zhuohao Yu, Jindong Wang, Qiang Heng, Hao Chen, Wei Ye, Rui Xie, Xing Xie,
405 and Shikun Zhang. Exploring vision-language models for imbalanced learning. *arXiv preprint*
406 *arXiv:2304.01457*, 2023.
- 407 Tz-Ying Wu, Pedro Morgado, Pei Wang, Chih-Hui Ho, and Nuno Vasconcelos. Solving long-tailed
408 recognition with deep realistic taxonomic classifier. In *European Conference on Computer Vision*
409 *(ECCV)*, 2020.
- 410 Taihong Xiao, Zirui Wang, Liangliang Cao, Jiahui Yu, Shengyang Dai, and Ming-Hsuan Yang. Ex-
411 ploiting category names for few-shot classification with vision-language models. *arXiv preprint*
412 *arXiv:2211.16594*, 2022.
- 413 Ze Yang, Tiange Luo, Dong Wang, Zhiqiang Hu, Jun Gao, and Liwei Wang. Learning to navigate
414 for fine-grained classification. In *Proceedings of the European conference on computer vision*
415 *(ECCV)*, pp. 420–435, 2018.
- 416 Dengyong Zhou, Olivier Bousquet, Thomas Lal, Jason Weston, and Bernhard Schölkopf. Learning
417 with local and global consistency. *Advances in neural information processing systems*, 16, 2003a.
- 418 Dengyong Zhou, Jason Weston, Arthur Gretton, Olivier Bousquet, and Bernhard Schölkopf. Rank-
419 ing on data manifolds. *Advances in neural information processing systems*, 16, 2003b.
- 420 Xinqi Zhu and Michael Bain. B-cnn: branch convolutional neural network for hierarchical classifi-
421 cation. *arXiv preprint arXiv:1709.09890*, 2017.