

Making Better Mistakes: Leveraging Class Hierarchies with Deep Networks

Luca Bertinetto* Romain Mueller* Konstantinos Tertikas Sina Samangooei Nicholas A. Lord*
 {luca.bertinetto, romain.mueller, konstantinos.tertikas, sina, nick.lord}@five.ai
www.five.ai

Abstract

Deep neural networks have improved image classification dramatically over the past decade, but have done so by focusing on performance measures that treat all classes other than the ground truth as equally wrong. This has led to a situation in which mistakes are less likely to be made than before, but are equally likely to be absurd or catastrophic when they do occur. Past works have recognised and tried to address this issue of mistake severity, often by using graph distances in class hierarchies, but this has largely been neglected since the advent of the current deep learning era in computer vision. In this paper, we aim to renew interest in this problem by reviewing past approaches and proposing two simple modifications of the cross-entropy loss which outperform the prior art under several metrics on two large datasets with complex class hierarchies: *tieredImageNet* and *iNaturalist*'19.

1. Introduction

Image classification networks have improved greatly over recent years, but generalisation remains imperfect, and test-time errors do of course occur. Conventionally, such errors are defined with respect to a single ground-truth class and reported using one or more top- k measures (k typically set to 1 or 5). However, this practice imposes certain notions of what it means to make a mistake, including treating all classes other than the “true” label as equally wrong. This may not actually correspond to our intuitions about desired classifier behaviour, and for some applications this point may prove crucial. Take the example of an autonomous vehicle observing an object on the side of the road: whatever measure of classifier performance we use, we can certainly agree that mistaking a lamppost for a tree is less of a problem than mistaking a person for a tree, as such a mistake would have crucial implications in terms both of prediction and planning. If we want to take such considerations into account, we must incorporate a nontrivial model of the re-

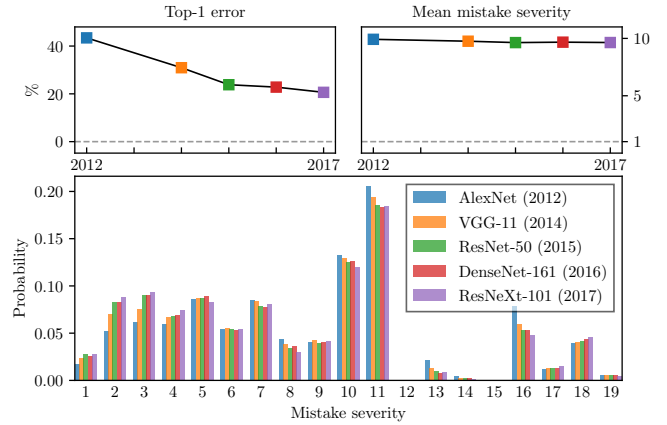


Figure 1: Top-1 error and distribution of mistakes w.r.t. the WordNet hierarchy for well-known deep neural network architectures on ImageNet: see text for definition of mistake severity. The top-1 error has enjoyed a spectacular improvement in the last few years, but even though the number of mistakes has decreased in absolute terms, the severity of the mistakes made has remained fairly unchanged over the same period. The grey dashed lines denote the minimal possible values of the respective measures.

lationships between classes, and accordingly rethink more broadly what it means for a network to “make a mistake”. One natural and convenient way of representing these class relationships is through a taxonomic hierarchy tree.

This idea is not new. In fact, it was once fairly common across various machine learning application domains to consider class hierarchy when designing classifiers, as surveyed in Silla & Freitas [37]. That work assembled and categorised a large collection of hierarchical classification problems and algorithms, and suggested widely applicable measures for quantifying classifier performance in the context of a given class hierarchy. The authors noted that the hierarchy-informed classifiers of the era typically empirically outperformed “flat” (*i.e.* hierarchy-agnostic) classifiers even under standard metrics, with the performance gap increasing further under the suggested hierarchical met-

*Equal contribution.

rics. Furthermore, class hierarchy is at the core of the ImageNet dataset: as detailed in Deng *et al.* [12], it was constructed directly from WordNet [26], itself a hierarchy originally designed solely to represent semantic relationships between words. Shortly after ImageNet’s introduction, works such as Deng *et al.* [11], Zhao *et al.* [49], and Verma *et al.* [42] explicitly noted that the underpinning WordNet hierarchy suggested a way of quantifying the severity of mistakes, and experimented with minimising hierarchical costs accordingly. Likewise, Deng *et al.* [10] presented a straightforward method for using a hierarchy-derived similarity matrix to define a more semantically meaningful compatibility function for image retrieval. Despite this initial surge of interest and the promising results accompanying it, the community shortly decided that hierarchical measures were not communicating substantially different information about classifier performance than top- k measures¹, and effectively discarded them. When the celebrated results in Krizhevsky *et al.* [20] were reported in flat top- k terms only, the precedent was firmly set for the work which followed in the deep learning era of image classification. Interest in optimising hierarchical performance measures waned accordingly.

We argue here that this problem is ripe for revisitation, and we begin by pointing to Fig. 1. Here, a *mistake* is defined as a top-1 prediction which differs from the ground-truth class, and the *severity* of such a mistake is the height of the lowest common ancestor of the predicted and ground-truth classes in the hierarchy. We see that while the flat top-1 accuracies of state-of-the-art classifiers have improved to impressive levels over the years, the distributions of the severities of the errors that *are* made have changed very *little* over this time. We hypothesise that this is due, at least in part, to the scarcity of modern learning methods which attempt to exploit prior information and preferences about class relationships in the interest of “making better mistakes”, whether this information is sourced from an offline taxonomy or otherwise. The few exceptions of which we are aware include Frome *et al.* [14], Wu *et al.* [43], Barz & Denzler [4], and a passing mention in Redmon & Farhadi [32]. In Sec. 2, we suggest a framework for thinking about these pieces of work, their predecessors, and some of their conceptual relatives.

The contributions of this work are as follows:

1. We review relevant literature within an explanatory framework which unifies a fairly disjoint prior art.
2. Building on the perspective gained from the preceding, we propose two methods that are both simple and effective at leveraging class hierarchies.

¹From Russakovsky *et al.* [36]: “[...] we found that all three measures of error (top-5, top-1, and hierarchical) produced the same ordering of results. Thus, since ILSVRC2012 we have been exclusively using the top-5 metric which is the simplest and most suitable to the dataset.”

3. We perform an extensive experimental evaluation to both demonstrate the effectiveness of the said methods compared to prior art and to encourage future work on the topic.

To ensure reproducibility, the PyTorch [29] code of all our experiments will be made available at github.com/fiveai/making-better-mistakes.

2. Framework and related work

We first suggest a simple framework for thinking about methods relevant to the problem of making better mistakes on image classification, beginning with the standard supervised setup. Consider a training set $\mathcal{S} = \{(x_i, C_i)\}_{i=1,\dots,N}$ which pairs N images $x_i \in \mathcal{I}$ with class labels $C_i \in \mathcal{C}$. A network architecture implements the predictor function $\phi(x; \theta)$, whose parameters θ are learned by minimising the empirical risk

$$\frac{1}{N} \sum_{i=1}^N \mathcal{L}(\phi(x_i; \theta), y(C_i)) + \mathcal{R}(\theta), \quad (1)$$

where the loss function \mathcal{L} compares the predictor’s output $\phi(x_i; \theta)$ to an embedded representation $y(C_i)$ of each example’s class, and \mathcal{R} is a regulariser.

Under common choices such as cross-entropy for \mathcal{L} and one-hot embedding for y , it is easy to see that the framework is agnostic of relationships between classes. The question is how such class relationships \mathcal{H} can be incorporated into the loss in Eqn. 1. We identify the following three approaches:

1. Replacing class representation $y(C)$ with an alternate embedding $y^{\mathcal{H}}(C)$. Such “label-embedding” methods, discussed in Sec. 2.1, can draw their embedding both from taxonomic hierarchies and alternative sources.
2. Altering the loss function \mathcal{L} in terms of its arguments to produce $\mathcal{L}^{\mathcal{H}}(\phi(x; \theta), y(C))$, *i.e.* making the penalty assigned to a given output distribution and embedded label dependent on \mathcal{H} . Methods using these “hierarchical losses” are covered in Sec. 2.2.
3. Altering the function $\phi(x; \theta)$ to $\phi^{\mathcal{H}}(x; \theta)$, *i.e.* making hierarchically-informed architectural changes to the network, generally with the hope of introducing a favourable inductive bias. We cover these “hierarchical architectures” in Sec. 2.3.

While a regulariser $\mathcal{R}^{\mathcal{H}}$ is certainly feasible, it is curiously rare in practice: [49] is the only example we know of.

2.1. Label-embedding methods

These methods map class labels to vectors whose relative locations represent semantic relationships, and optimise a loss on these embedded vectors. The DeVISE method

of Frome *et al.* [14] maps target classes onto a unit hypersphere via the method of Mikolov *et al.* [24], assigning terms with similar contexts to similar representations through analysis of unannotated Wikipedia text. The loss function is a ranking loss which penalises the extent to which the output is more cosine-similar to false label embeddings than to the correct one. They learn a linear mapping from a pre-trained visual feature pipeline (based on AlexNet) to the embedded labels, then fine-tune the visual pipeline itself through backpropagation. Romera-Paredes & Torr [34] note that their one-line solution for learning an analogous linear mapping for zero-shot classification should easily extend to accommodating these sorts of embeddings. In Hinton *et al.* [18], the role of the label embedding function is played by a temperature-scaled pre-existing classifier ensemble. This ensemble is “distilled” into a smaller DNN through cross-entropy minimisation against the ensemble’s output². For zero-shot classification, Xian *et al.* [44] experiment with various independent embedding methods, as is also done in Akata *et al.* [2]: annotated attributes, word2vec [25], glove [30], and the WordNet hierarchy. Their ranking loss function is functionally equivalent to that in Frome *et al.* [14], and they learn a choice of linear mappings to these representations from the visual features output by a fixed CNN. Barz & Denzler [4] present an embedding algorithm which iteratively maps examples onto a hypersphere such that all cosine distances represent similarities derived from lowest common ancestor (LCA) height in a given hierarchy tree. They train a full deep model, and alternately present two different losses: (1) a linear loss based on cosine distance to the embedded class vectors, used for retrieval, and (2) the standard cross-entropy loss on the output of a fully connected/softmax layer added after the embedding layer, for classification.

2.2. Hierarchical losses

In these methods, the loss function itself is parametrised by the class hierarchy such that a higher penalty is assigned to the prediction of a more distant relative of the true label. Deng *et al.* [11] simply train kNN- and SVM-based classifiers to minimise the expected WordNet LCA height directly. Zhao *et al.* [49] modify standard multi-class logistic regression by recomputing output class probabilities as the normalised class-similarity-weighted sums of all class probabilities calculated as per the usual regression. They also regularise feature selection using an “overlapping-group-lasso penalty” which encourages the use of similar features for closely related classes, a rare example of a hierarchical regulariser. Verma *et al.* [42] incorporate normalised LCA height into a “context-sensitive loss function” while learning a separate metric at each node in a taxonomy tree for

²This refers to the “distillation” section of [18]. There is a separate hierarchical classification section discussed later.

nearest-neighbour classification. Wu *et al.* [43] consider the specific problem of food classification. They begin with a standard deep network pipeline and add one fully connected/softmax layer to the end of the (shared) pipeline for each level of the hierarchy, so as to explicitly predict each example’s lineage. The loss is the sum of standard log losses over all of the levels, with a single user-specified coefficient for reweighting the loss of all levels other than the bottom one. Because this method is inherently susceptible to producing marginal probabilities that are inconsistent across levels, it incorporates a subsequent label-propagation smoothing step. Alsallakh *et al.* [5] likewise use a standard deep architecture as their starting point, but instead add branches strategically to intermediate pipeline stages. They thereby force the net to classify into offline-determined superclasses at the respective levels, backpropagating error in these intermediate predictions accordingly. On deployment, these additions are simply discarded.

2.3. Hierarchical architectures

These methods attempt to incorporate class hierarchy into the classifier architecture without necessarily changing the loss function otherwise. The core idea is to “divide and conquer” at the structural level, with the classifier assigning inputs to superclasses at earlier layers and making fine-grained distinctions at later ones. In the context of language models, it was noted at least as early as Goodman [15] that classification with respect to an IS-A hierarchy tree could be formulated as a tree of classifiers outputting conditional probabilities, with the product of the conditionals along a given leaf’s ancestry representing its posterior; motivated by efficiency, Morin & Bengio [27] applied this observation to a binary hierarchy derived from WordNet. Redmon & Farhadi [32] propose a modern deep-learning variant of this framework in the design of the YOLOv2 object detection and classification system. Using a version of WordNet pruned into a tree, they effectively train a conditional classifier at every parent node in the tree by using one softmax layer per sibling group and training under the usual cross-entropy loss over leaf posteriors. While their main aim is to enable the integration of the COCO detection dataset with ImageNet, they suggest that graceful degradation on new or unknown object categories might be an incidental benefit. Brust & Denzler [7] propose an extension of conditional classifier chains to the more general case of DAGs.

The above approaches can be seen as a limiting case of hierarchical classification, in which every split in the hierarchy is cast as a separate classification problem. Many hierarchical classifiers fall between this extreme and that of flat classification, working in terms of a coarser-grained conditionality in which a “generalist” makes hard or soft assignments to groupings of the target classes before then distinguishing the group members from one another using

“experts”. Xiao *et al.* [45], the quasi-ensemble section of Hinton *et al.* [18], Yan *et al.* [46], and Ahmed *et al.* [1] all represent modern variations on this theme (which first appears no later than [19]). Additionally, the listed methods all use some form of low-level feature sharing either via architectural constraint or parameter cloning, and all infer the visual hierarchy dynamically through confusion clustering or latent parameter inference. Alsallakh *et al.* [5] make the one proposal of which we are aware which combines hierarchical architectural modifications (at *train* time) with a hierarchical loss, as described in Sec. 2.2. At test time, however, the architecture is that of an unmodified AlexNet, and all superclass “assignment” is purely implicit.

3. Method

We now outline two simple methods that allow us to leverage class hierarchies in order to make better mistakes on image classification. We concentrate on the case where the output of the network is a categorical distribution over classes for each input image and denote the corresponding distribution as $p(C) = \phi_C(x; \theta)$, where subscripts denote vector indices and x and θ are omitted for brevity. In Sec. 3.1, we describe the *hierarchical cross-entropy* (HXE), a straightforward example of the hierarchical losses reviewed in Sec. 2.2. This approach expands each class probability into the chain of conditional probabilities defined by its unique lineage in a given hierarchy tree. It then reweights the corresponding terms in the loss so as to penalise classification mistakes in a way that is informed by the hierarchy. In Sec. 3.2, we suggest an easy choice of embedding function to implement the label-embedding framework covered in Sec. 2.1. The resulting *soft labels* are PMFs over \mathcal{C} whose values decay exponentially w.r.t. an LCA-based distance to the ground truth.

3.1. Hierarchical cross-entropy

When the hierarchy \mathcal{H} is a tree, it corresponds to a unique factorisation of the categorical distribution $p(C)$ over classes in terms of the conditional probabilities along the path connecting each class to the root of the tree. Denoting the path from a leaf node C to the root R as $C^{(0)}, \dots, C^{(h)}$, the probability of class C can be factorised as

$$p(C) = \prod_{l=0}^{h-1} p(C^{(l)}|C^{(l+1)}), \quad (2)$$

where $C^{(0)} = C$, $C^{(h)} = R$, and $h \equiv h(C)$ is the height of the node C . Note that we have omitted the last term $p(C^{(h)}) = p(R) = 1$. Conversely, the conditionals can be written in terms of the class probabilities as

$$p(C^{(l)}|C^{(l+1)}) = \frac{\sum_{A \in \text{Leaves}(C^{(l)})} p(A)}{\sum_{B \in \text{Leaves}(C^{(l+1)})} p(B)}, \quad (3)$$

where $\text{Leaves}(C)$ denotes the set of leaf nodes of the subtree starting at node C .

A direct way to incorporate hierarchical information in the loss is to hierarchically factorise the output of the classifier according to Eqn. 2 and define the total loss as the reweighted sum of the cross-entropies of the conditional probabilities. This leads us to define the *hierarchical cross-entropy* (HXE) as

$$\mathcal{L}_{\text{HXE}}(p, C) = - \sum_{l=0}^{h(C)-1} \lambda(C^{(l)}) \log p(C^{(l)}|C^{(l+1)}), \quad (4)$$

where $\lambda(C^{(l)})$ is the weight associated with the edge node $C^{(l+1)} \rightarrow C^{(l)}$, see Fig. 2a. Even though this loss is expressed in terms of conditional probabilities, it can be easily applied to models that output class probabilities using Eqn. 3. Note that \mathcal{L}_{HXE} reduces to the standard cross-entropy when all weights are equal to 1. This limit case, which was briefly mentioned by Redmon & Farhadi in their YOLO-v2 paper [32], results only in architectural changes but does not incorporate hierarchical information in the loss directly.

Eqn. 4 has an interesting information-theoretical interpretation: since each term $\log p(C^{(l)}|C^{(l+1)})$ corresponds to the information associated with the edge $C^{(l+1)} \rightarrow C^{(l)}$ in the hierarchy, the HXE corresponds to discounting the information associated with each of these edges differently. Note that since the HXE is expressed in terms of conditional probabilities, the reweighting in Eqn. 4 is not equivalent to reweighting the cross-entropy for each possible ground truth class independently (as done, for instance, in [21, 9]).

A sensible choice for the weights is to take

$$\lambda(C) = \exp(-\alpha h(C)), \quad (5)$$

where $h(C)$ is the height of node C and $\alpha > 0$ is a hyperparameter that controls the extent to which information is discounted down the hierarchy. The higher the value of α , the higher the preference for “generic” as opposed to “fine-grained” information, because classification errors related to nodes further away from the root receive a lower loss. While such a definition has the advantage of simplicity, one could think of other meaningful weightings, such as ones depending on the branching factor of the tree or encoding a preference towards specific classes. We concentrate on equation 5 here, as it is both simple and easily interpretable while leaving a more systematic exploration of different weighting strategies for future work.

3.2. Soft labels

Our second approach to incorporating hierarchical information, *soft labels*, is a label-embedding approach as described in Sec. 2.1. These methods use a mapping function $y(C)$ to associate classes with representations which

encode class-relationship information that is absent in the trivial case of the one-hot representation. In the interest of simplicity, we choose a mapping function $y^{\text{soft}}(C)$ which outputs a categorical distribution over the classes. This enables us to simply use the standard cross-entropy loss:

$$\mathcal{L}_{\text{Soft}}(p, C) = - \sum_{A \in \mathcal{C}} y_A^{\text{soft}}(C) \log p(A), \quad (6)$$

where the soft label embedding is given componentwise by

$$y_A^{\text{soft}}(C) = \frac{\exp(-\beta d(A, C))}{\sum_{B \in \mathcal{C}} \exp(-\beta d(B, C))}, \quad (7)$$

for class distance function d and parameter β . This loss is illustrated in Fig. 2b. For the distance function $d(C_i, C_j)$, we use the height of $\text{LCA}(C_i, C_j)$ divided by the height of the tree. To understand the role of the hyperparameter β , note that values of β that are much bigger than the typical inverse distance in the tree result in a label distribution that is nearly one-hot, *i.e.* $y_A(C) \simeq \delta_{AC}$, in which case the cross-entropy reduces to the familiar single-term log-loss expression. Conversely, for very small values of β the label distribution is near-uniform. Between these extremes, greater probability mass is assigned to classes more closely related to the ground truth, with the magnitude of the difference controlled by β .

We offer two complementary interpretations that motivate this representation (besides its ease). For one, the distribution describing each target class can be considered to be a model of the actual uncertainty that a labeller (*e.g.* human) would experience due to visual confusion between closely related classes. It could also be thought of as encoding the extent to which a common response to different classes is *required* of the classifier, *i.e.* the imposition of correlations between outputs, where higher correlations are expected for more closely related classes. This in turn suggests a connection to the superficially different but conceptually related *distillation* method of Hinton *et al.* [18], in which correlations between a large network’s responses to different classes are mimicked by a smaller network to desirable effect. Here, we simply supply these correlations directly, using widely available hierarchies.

Another important connection is the one to the technique of *label smoothing* [38], in which the “peaky” distribution represented by a one-hot label is combined with the uniform distribution. This technique has been widely adopted to regularise the training of large neural networks (*e.g.* [38, 8, 41, 50]), but has only very recently [28] been studied more thoroughly.

4. Evaluation

In the following, we first describe the datasets (Sec. 4.1) and metrics (Sec. 4.2) comprising the setup common to all

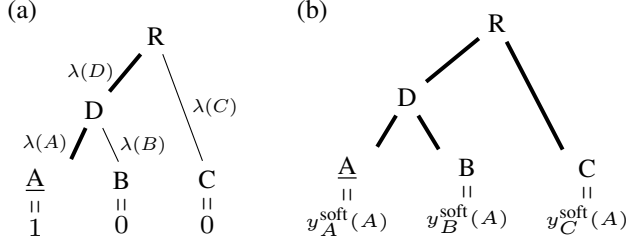


Figure 2: Representations of the *HXE* (Sec. 3.1) and *soft labels* (Sec. 3.2) losses for a simple illustrative hierarchy are drawn in subfigures (a) and (b) respectively. The ground-truth class is underlined, and the edges contributing to the total value of the loss are drawn in bold.

of our experiments. Then, in Sec. 4.3, we empirically evaluate our two simple proposals and compare them to the prior art. Finally, we experiment with random hierarchies to understand when and how information on class relatedness can help classification.

4.1. Datasets

In our experiments, we use *tieredImageNet* [33] (a large subset of ImageNet/ILSVRC’12 [36]) and *iNaturalist’19* [40], two datasets with hierarchies that are *a)* significantly different from one another and *b)* complex enough to cover a large number of visual concepts. ImageNet aims to populate the WordNet [26] hierarchy of nouns, with WordNet itself generated by inspecting IS-A lexical relationships. By contrast, *iNaturalist’19* [40] has a biological taxonomy [35] at its core.

tieredImageNet was originally introduced by Ren *et al.* [33] for the problem of few-shot classification, in which the sets of classes between dataset splits are disjoint. The authors’ motivation in creating the dataset was to use the WordNet hierarchy to generate splits containing significantly different classes, facilitating better assessment of few-shot classifiers by enforcing problem difficulty.

Although our task and motivations are different, we chose this dataset because of the large portion of the WordNet hierarchy spanned by its classes. To make it suitable for the problem of (standard) image classification, we re-sampled the dataset so as to represent all classes across the train, validation, and test splits. Moreover, since the method proposed in Section 3.1 and YOLO-v2 [32] require that the graph representing the hierarchy is a tree, we modified the graph of the spanned WordNet hierarchy slightly to comply with this assumption (more details available in Appendix A). After this procedure, we obtained a tree of height 13 and 606,702 images from 608 different classes, which we randomly assigned to training, validation, and test splits with respective probabilities 0.7, 0.15, and 0.15. We refer to this modified version of *tieredImageNet* as *tieredImageNet*.

H.

iNaturalist is a dataset of images of organisms that has so far mainly been used to evaluate fine-grained visual categorisation methods. The dataset construction protocol differs significantly from the one used for ImageNet in that it relies on passionate volunteers instead of workers paid per task [40]. Importantly, for the 2019 edition of the CVPR Fine-Grained Visual Categorization Workshop, metadata with hierarchical relationships between species have been released. In contrast to WordNet, this taxonomy is an 8-level complete tree that can readily be used in our experiments without modifications. Since the labels for the test set are not public, we randomly re-sampled three splits from the total of 271,273 images from 1010 classes, again with probabilities 0.7, 0.15, and 0.15 for the training, validation, and test set, respectively. We refer to this modified version of iNaturalist’19 as *iNaturalist19-H*.

4.2. Metrics

We consider three measures of performance, covering different interpretations of a classifier’s *mistakes*.

Top- k error. Under this measure, an example is defined as correctly classified if the ground truth is among the top k classes with the highest likelihood. This is the measure normally used to compare classifiers, usually with $k=1$ or $k=5$. Note that this measure considers all mistakes of the classifier equally, irrespective of how “similar” the predicted class is to the ground truth.

Hierarchical measures. We also consider measures that, in contrast to the top- k error, do weight the severity of mistakes. We use the height of the lowest common ancestor (LCA) between the predicted class and the ground truth as a core severity measure, as originally proposed in the papers describing the creation of ImageNet [12, 11]. As remarked in [11], this measure should be thought of in logarithmic terms, as the number of confounded classes is exponential in the height of the ancestor. We also experimented with the Jiang-Conrath distance as suggested by Deselaers & Ferrari [13], but did not observe meaningful differences wrt. the height of the LCA.

We consider two measures that utilise the height of the LCA between nodes in the hierarchy.

- The **hierarchical distance of a mistake** is the height of the LCA between the ground truth and the predicted class *when the input is misclassified, i.e.* when the class with the maximum likelihood is incorrect. Hence, it measures the severity of misclassification when only a single class can be considered as a prediction.
- The **average hierarchical distance of top- k** , instead, takes the mean LCA height between the ground truth and each of the k most likely classes. This measure is important, for example, when multiple hypotheses of

a classifier can be considered for a certain downstream task.

4.3. Experimental results

In the following, we analyse the performance of the two approaches described in Sec. 3.1 and Sec. 3.2, which we denote by *HXE* and *soft labels*, respectively. Besides a vanilla cross-entropy-based flat classifier, we also implemented and compared against the methods proposed by Redmon & Farhadi [32] (YOLO-v2)³, Frome *et al.* [14] (DeViSE), and Barz & Denzler [4]. As mentioned in Sec. 1, these methods represent, to the best of our knowledge, the only modern attempts to deliberately reduce the semantic severity of a classifier’s mistakes that are generally applicable to any modern architecture. Note, though, that we do not run DeViSE on *iNaturalist19-H*, as the class IDs of this dataset are alien to the corpus used by word2vec [24].

Implementation details. Since we are interested in understanding the mechanisms by which the above metrics can be improved, it is essential to use a simple configuration that is common between all of the algorithms taken into account.

We use a ResNet-18 architecture (with weights pre-trained on ImageNet) trained with Adam [31] for 200,000 steps and mini-batches of size 256. We use a learning rate of $1e-5$ unless specified otherwise. To prevent overfitting, we adopt PyTorch’s basic data augmentation routines with default hyperparameters: `RandomHorizontalFlip()` and `RandomResizedCrop()`.

Further implementation details for all of the methods are deferred to Appendix B.

Main results. In Fig. 3 and 4 we show how it is possible to effectively trade off top-1 error to reduce hierarchical error, by simply adjusting the hyperparameters α and β in Eqn. 5 and 7. Specifically, increasing α corresponds to (exponentially) discounting information down the hierarchy, thus more severely penalising mistakes where the predicted class is further away from the ground truth. Similarly, decreasing β in the soft-label method amounts to progressively shifting the label mass away from the ground truth and towards the neighbouring classes. Both methods reduce to the cross-entropy in the respective limits $\alpha \rightarrow 0$ and $\beta \rightarrow \infty$. Moreover, notice that varying β affects the entropy of the distribution representing a soft label, where the two limit cases are $\beta=\infty$ for the standard one-hot case and $\beta=0$ for the uniform distribution. We experiment with $0.1 \leq \alpha \leq 0.6$ and $5 \leq \beta \leq 30$.

To limit noise in the evaluation procedure, for both of our methods and all of the competitors, we fit a 4th-degree polynomial to the validation loss (after having discarded the first 50,000 training steps) and pick the epoch corresponding to

³Note that this refers to the conditional classifier subsystem proposed in Sec. 4 of that work, not the main object detection system.

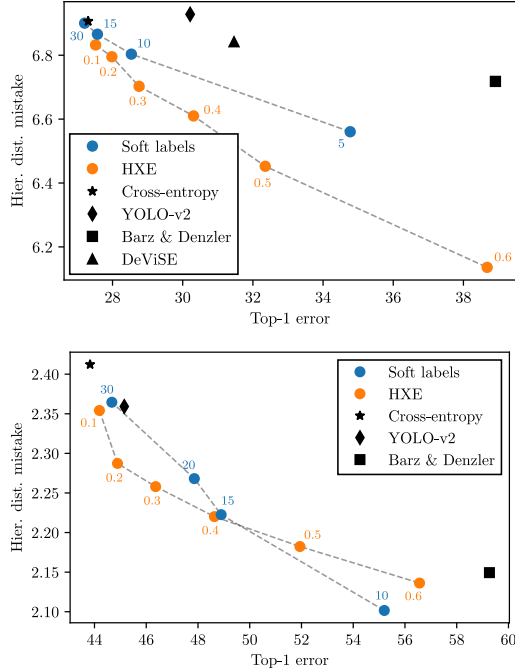


Figure 3: Top-1 error vs. hierarchical distance of mistakes, for *tieredImageNet-H* (top) and *iNaturalist19-H* (bottom). Points closer to the bottom-left corner of the plot are the ones achieving the best tradeoff.

its minimum along with its four neighbours. Then, to produce the points reported in our plots, we average the results obtained from these five epochs on the validation set, while reserving the test set for the experiments of Table 1. Notice how, in Fig. 4, when considering the hierarchical distance with $k=1$, methods are almost perfectly aligned along the plot diagonal, which demonstrates the strong linear correlation between this metric and the top-1 error. This result is consistent with what is observed in [36], which, in 2011, led the organisers of the ILSVRC workshop to discard rankings based on hierarchical distance.

When considering the other metrics described in Sec. 4.2, a different picture emerges. In fact, a tradeoff between top-1 error and hierarchical distance is evident in Fig. 3 and in the plots of Fig. 4 with $k=5$ and $k=20$. Notice how the points on the plots belonging to our methods outline a set of tradeoffs that subsumes the prior art. For example, in Fig. 3, given any desired tradeoff between top-1 error and hierarchical distance of mistakes on *tieredImageNet-H*, it is better to use HXE than any other method. A similar phenomenon is observable when considering the average hierarchical distance of top-5 and top-20 (Fig. 4), although in these cases it is better to use the soft labels. The only exception to this trend is represented by Barz & Denzler [4] on *tieredImageNet-H*, which can achieve slightly lower av-

erage hierarchical distance for $k=5$ or $k=20$ at a significant cost in terms of top-1 error.

Using the results illustrated in Fig. 3 and 4, we pick two reasonable operating points for both of our proposals: one for the high-distance/low-top1-error regime, and one for the low-distance/high-top1-error regime. We then run both of these configurations on the test sets and report our results on Table 1. The means are again obtained from the five best epochs, and we use the standard deviation to compute 95% confidence intervals.

The trends observed on the validation set largely repeat themselves on the test set. When one desires to prioritise top-1 error, then soft labels with high β or HXE with low α are more appropriate, as they outperform the cross-entropy on the hierarchical-distance-based metrics while being practically equivalent in terms of top-1 error. In cases where the hierarchical measures should be prioritised instead, it is preferable to use soft labels with low β or HXE with high α , depending on the particular choice of hierarchical metric. Although the method of Barz & Denzler is competitive in this regime, it also exhibits the worst deterioration in top-1 error with respect to the cross-entropy.

Our experiments generally indicate, over all tested methods, an inherent tension between performance in the top-1 sense and in the hierarchical sense. We speculate that there may be a connection between this tension and observations proceeding from the study of adversarial examples indicating a tradeoff between robustness and (conventional) accuracy, as in *e.g.* [39, 48].

Can hierarchies be arbitrary? Although the lexical WordNet hierarchy and the biological taxonomy of *iNaturalist* are not visual hierarchies per se, they arguably reflect meaningful visual relationships between the objects represented in the underlying datasets. Since deep networks leverage visual features, it is interesting to investigate the extent to which the structure of a particular hierarchy is important. In other words, what would happen with an arbitrary hierarchy, one that does not have any relationship with the visual world?

To answer this question, we randomised the nodes of the hierarchies and repeated our experiments. Results on *iNaturalist19-H* are displayed in Fig. 5 (*tieredImageNet-H* exhibits a similar trend). Again, we report tradeoff plots showing top-1 errors on the x-axis and metrics based on the height of the LCA (on the randomised hierarchy) on the y-axis. It is evident that the hierarchical distance metrics are significantly worse when using the random hierarchy. Although this is not surprising, the extent to which the results deteriorate is remarkable. This suggests that the inherent nature of the structural relationship expressed by a hierarchy is paramount for learning classifiers that, besides achieving competitive top-1 accuracy, are also able to make better mistakes.

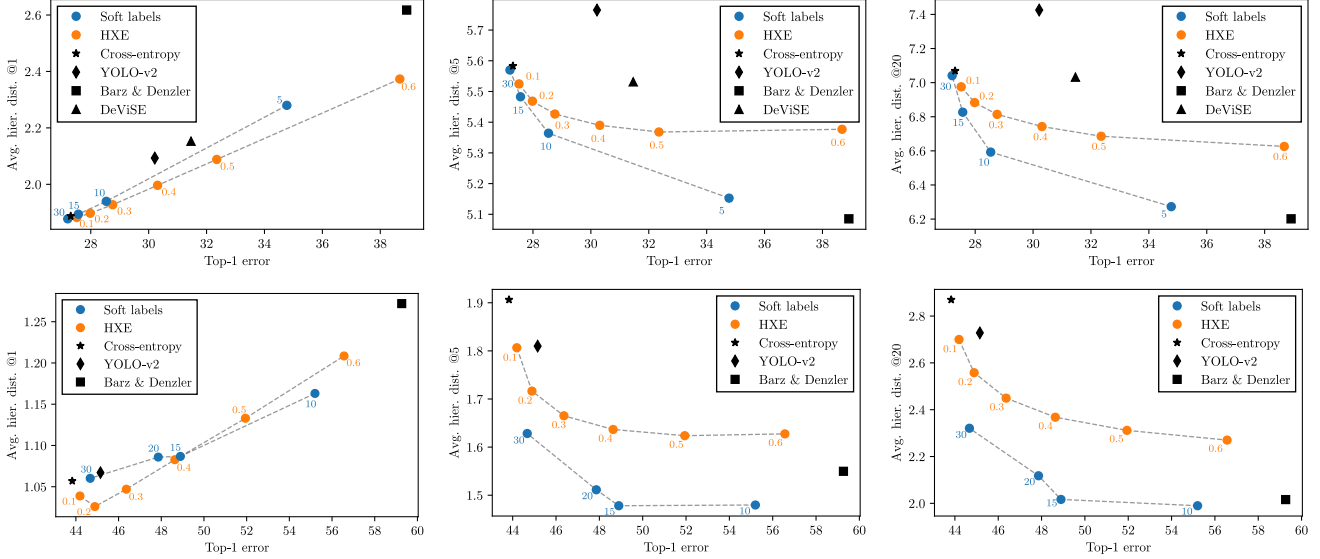


Figure 4: Top-1 error vs. average hierarchical distance of top- k (with $k \in \{1, 5, 20\}$) for *tieredImageNet-H* (top three) and *iNaturalist19-H* (bottom three). Points closer to the bottom-left corner of the plot are the ones achieving the best tradeoff.

Table 1: Results on the test sets of *tieredImageNet-H* (top) and *iNaturalist19-H* (bottom), with 95% confidence intervals. For each column of each dataset, the best entry is highlighted in yellow, while the worst is highlighted in gray.

	Hier. dist. mistake ↓	Avg. hier. dist. @1 ↓	Avg. hier. dist. @5 ↓	Avg. hier. dist. @20 ↓	Top-1 error ↓
CROSS-ENTROPY	6.89 ± 0.004	1.90 ± 0.002	5.59 ± 0.004	7.07 ± 0.007	27.55 ± 0.038
BARZ & DENZLER [4]	6.72 ± 0.017	2.62 ± 0.014	5.09 ± 0.009	6.21 ± 0.007	39.03 ± 0.157
YOLO-v2 [32]	6.91 ± 0.006	2.10 ± 0.002	5.77 ± 0.012	7.42 ± 0.018	30.43 ± 0.030
DEVISE [14]	6.83 ± 0.005	2.17 ± 0.003	5.54 ± 0.003	7.04 ± 0.002	31.69 ± 0.058
HXE $\alpha=0.1$ (ours)	6.83 ± 0.009	1.89 ± 0.003	5.53 ± 0.004	6.98 ± 0.008	27.68 ± 0.066
HXE $\alpha=0.5$ (ours)	6.46 ± 0.026	2.11 ± 0.021	5.37 ± 0.003	6.69 ± 0.008	32.61 ± 0.443
SOFT-LABELS $\beta=15$ (ours)	6.83 ± 0.005	1.90 ± 0.004	5.49 ± 0.002	6.83 ± 0.002	27.78 ± 0.063
SOFT-LABELS $\beta=5$ (ours)	6.56 ± 0.009	2.29 ± 0.008	5.16 ± 0.006	6.28 ± 0.005	35.00 ± 0.096
CROSS-ENTROPY	2.41 ± 0.003	1.05 ± 0.004	1.90 ± 0.004	2.87 ± 0.006	43.77 ± 0.138
BARZ & DENZLER [4]	2.19 ± 0.008	1.27 ± 0.007	1.56 ± 0.006	2.03 ± 0.005	57.83 ± 0.137
YOLO-v2 [32]	2.37 ± 0.006	1.07 ± 0.007	1.81 ± 0.008	2.73 ± 0.009	45.23 ± 0.202
HXE $\alpha=0.1$ (ours)	2.35 ± 0.007	1.04 ± 0.004	1.80 ± 0.004	2.70 ± 0.009	44.28 ± 0.171
HXE $\alpha=0.6$ (ours)	2.13 ± 0.003	1.21 ± 0.004	1.62 ± 0.003	2.68 ± 0.003	56.61 ± 0.241
SOFT-LABELS $\beta=30$ (ours)	2.35 ± 0.002	1.05 ± 0.005	1.62 ± 0.005	2.32 ± 0.004	44.75 ± 0.139
SOFT-LABELS $\beta=10$ (ours)	2.10 ± 0.005	1.16 ± 0.006	1.47 ± 0.004	1.99 ± 0.003	55.16 ± 0.196

Curiously, for the soft labels, the top-1 error of the random hierarchy is consistently *lower* than its “real” hierarchy counterpart. We speculate this might be due to the structural constraints imposed by a hierarchy anchored to the visual world, which can limit a neural network from opportunistically learning correlations that allow it to achieve low top-1 error (at the expense of ever more brittle generalisation). Indeed, the authors of [47] noted that it is more difficult to train a deep network to map real images to random labels than it is to do so with random images. The most likely explanation for this is that common visual features, which are inescapably shared by closely related examples, dictate common responses.

5. Conclusion

Since the advent of deep learning, the computer vision community’s interest in making better classification mistakes seems to have nearly vanished. In this paper, we have shown that this problem is still very much open and ripe for a comeback. We have demonstrated that two simple baselines that modify the cross-entropy loss are able to outperform the few modern methods tackling this problem. Improvements in this task are undoubtedly possible, but it is important to note the delicate balance between standard top-1 accuracy and mistake severity. Our hope is that the results presented in this paper are soon to be surpassed by the new competitors that it has inspired.

- [24] Tomas Mikolov, Kai Chen, Greg Corrado, and Jeffrey Dean. Efficient estimation of word representations in vector space. *arXiv preprint arXiv:1301.3781*, 2013. 2, 6, 11, 12
- [25] Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg S Corrado, and Jeff Dean. Distributed representations of words and phrases and their compositionality. In *Advances in Neural Information Processing Systems*, 2013. 3
- [26] George A Miller. *WordNet: An electronic lexical database*. 1998. 2, 5, 11
- [27] Frederic Morin and Yoshua Bengio. Hierarchical probabilistic neural network language model. In *Aistats*. Citeseer, 2005. 3, 11
- [28] Rafael Müller, Simon Kornblith, and Geoffrey Hinton. When does label smoothing help? In *Advances in Neural Information Processing Systems*, 2019. 5
- [29] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. Pytorch: An imperative style, high-performance deep learning library. In *Advances in Neural Information Processing Systems*, 2019. 2
- [30] Jeffrey Pennington, Richard Socher, and Christopher Manning. Glove: Global vectors for word representation. In *Proceedings of the 2014 conference on empirical methods in natural language processing*, 2014. 3
- [31] Sashank J Reddi, Satyen Kale, and Sanjiv Kumar. On the convergence of adam and beyond. In *International Conference on Learning Representations*, 2019. 6, 11, 12
- [32] Joseph Redmon and Ali Farhadi. Yolo9000: better, faster, stronger. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2017. 2, 3, 4, 5, 6, 8, 11
- [33] Mengye Ren, Eleni Triantafillou, Sachin Ravi, Jake Snell, Kevin Swersky, Joshua B Tenenbaum, Hugo Larochelle, and Richard S Zemel. Meta-learning for semi-supervised few-shot classification. In *International Conference on Learning Representations*, 2018. 5
- [34] Bernardino Romera-Paredes and Philip Torr. An embarrassingly simple approach to zero-shot learning. In *International Conference on Machine Learning*, 2015. 3
- [35] Michael A Ruggiero, Dennis P Gordon, Thomas M Orrell, Nicolas Bailly, Thierry Bourgoïn, Richard C Brusca, Thomas Cavalier-Smith, Michael D Guiry, and Paul M Kirk. A higher level classification of all living organisms. *PloS one*, 2015. 5
- [36] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, et al. Imagenet large scale visual recognition challenge. *International Journal of Computer Vision*, 2015. 2, 5, 7, 11
- [37] Carlos N Silla and Alex A Freitas. A survey of hierarchical classification across different application domains. *Data Mining and Knowledge Discovery*, 2011. 1
- [38] Christian Szegedy, Vincent Vanhoucke, Sergey Ioffe, Jon Shlens, and Zbigniew Wojna. Rethinking the inception architecture for computer vision. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2016. 5
- [39] Dimitris Tsipras, Shibani Santurkar, Logan Engstrom, Alexander Turner, and Aleksander Madry. Robustness may be at odds with accuracy. In *International Conference on Learning Representations*, 2018. 7
- [40] Grant Van Horn, Oisin Mac Aodha, Yang Song, Yin Cui, Chen Sun, Alex Shepard, Hartwig Adam, Pietro Perona, and Serge Belongie. The inaturalist species classification and detection dataset. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2018. 5, 6
- [41] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Advances in neural information processing systems*, pages 5998–6008, 2017. 5
- [42] Nakul Verma, Dhruv Mahajan, Sundararajan Sellamanickam, and Vinod Nair. Learning hierarchical similarity metrics. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2012. 2, 3
- [43] Hui Wu, Michele Merler, Rosario Uceda-Sosa, and John R Smith. Learning to make better mistakes: Semantics-aware visual food recognition. In *Proceedings of the 24th ACM international conference on Multimedia*, 2016. 2, 3
- [44] Yongqin Xian, Zeynep Akata, Gaurav Sharma, Quynh Nguyen, Matthias Hein, and Bernt Schiele. Latent embeddings for zero-shot classification. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2016. 3
- [45] Tianjun Xiao, Jiaxing Zhang, Kuiyuan Yang, Yuxin Peng, and Zheng Zhang. Error-driven incremental learning in deep convolutional neural network for large-scale image classification. In *Proceedings of the 22nd ACM international conference on Multimedia*, 2014. 3
- [46] Zhicheng Yan, Hao Zhang, Robinson Piramuthu, Vignesh Jagadeesh, Dennis DeCoste, Wei Di, and Yizhou Yu. Hd-cnn: hierarchical deep convolutional neural networks for large scale visual recognition. In *IEEE International Conference on Computer Vision*, 2015. 3
- [47] Chiyuan Zhang, Samy Bengio, Moritz Hardt, Benjamin Recht, and Oriol Vinyals. Understanding deep learning requires rethinking generalization. In *International Conference on Learning Representations*, 2016. 8
- [48] Hongyang Zhang, Yaodong Yu, Jiantao Jiao, Eric Xing, Laurent El Ghaoui, and Michael Jordan. Theoretically principled trade-off between robustness and accuracy. In *International Conference on Machine Learning*, 2019. 7
- [49] Bin Zhao, Fei Li, and Eric P Xing. Large-scale category structure aware image categorization. In *Advances in Neural Information Processing Systems*, 2011. 2, 3
- [50] Barret Zoph, Vijay Vasudevan, Jonathon Shlens, and Quoc V Le. Learning transferable architectures for scalable image recognition. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2018. 5

A. Pruning the WordNet hierarchy

The ImageNet dataset [36] was generated by populating the WordNet [26] hierarchy of nouns with images. WordNet is structured as a graph composed of a set of IS-A parent-child relationships. Similarly to the work of Morin & Bengio [27] and Redmon & Farhadi [32], our proposed hierarchical cross entropy loss (HXL, Sec. 3.1) also relies on the assumption that the hierarchy underpinning the data takes the form of a tree. Therefore, we modified the hierarchy to obtain a tree from the WordNet graph.

First, for each class, we found the longest path from the corresponding node to the root. This amounts to selecting the paths with the highest discriminative power with respect to the image classes. When multiple such paths existed, we selected the one with the minimum number of new nodes and added it to the new hierarchy. Second, we removed the few non-leaf nodes with a single child, as they do not possess any discriminative power.

Finally, we observed that the pruned hierarchy’s root is not PHYSICAL ENTITY, as one would expect, but rather the more general ENTITY. This is problematic, since ENTITY contains both physical objects *and* abstract concepts, while *tieredImageNet-H* classes only represent physical objects. Upon inspection, we found that this was caused by the classes BUBBLE, TRAFFIC SIGN, and TRAFFIC LIGHTS being connected to SPHERE and SIGN, which are considered abstract concepts in the WordNet hierarchy. Instead, we connected them to SPHERE, ARTIFACT and SIGNBOARD, respectively, thus connecting them to PHYSICAL ENTITY.

Even though our second proposed method (*soft labels*), as well as the cross-entropy baseline, DeViSE [14] and Barz & Denzler [4], do not make any assumption regarding the structure of the hierarchy, we ran them using this obtained pruned hierarchy for consistency of the experimental setup.

B. More implementation details

In order to perform meaningful comparisons, we adopted a simple configuration (network architecture, optimiser, data augmentation, ...) and used it for all the methods presented in this paper. This configuration is already stated in the implementation details of Sec. 4.3, but we report it again here for convenience.

We used a ResNet-18 architecture (with weights pre-trained on ImageNet) trained with Adam [31] for 200,000 steps and mini-batch size of 256. We used a learning rate of $1e-5$ unless specified otherwise. To prevent overfitting, we adopted PyTorch’s basic data augmentation routines with default hyperparameters: `RandomHorizontalFlip()` and `RandomResizedCrop()`. For both datasets, images have been resized to 224×224 .

Below, we provide further information about the methods we compared against, together with the few minor im-

Listing 1: Network head used for DeViSE.

```
model.fc = torch.nn.Sequential(
    torch.nn.Linear(in_features=512,
                    out_features=512),
    torch.nn.ReLU(),
    torch.nn.BatchNorm1d(512),
    torch.nn.Linear(in_features=512,
                    out_features=300)
)
```

plementation choices we had to make. As mentioned in Sec. 1, these methods represent, to the best of our knowledge, the only modern attempts to deliberately reduce the semantic severity of a classifier’s mistakes that are generally applicable to any modern architecture.

YOLO-v2. In motivating the hierarchical variant of the YOLO-v2 framework, Redmon & Farhadi [32, Sec. 4], mention the need of integrating the smaller COCO detection dataset [22] with the larger ImageNet classification dataset under a unified class hierarchy. Their approach too relies on a heuristic for converting the WordNet graph into a tree, and then effectively training a conditional classifier at every parent node in the tree by using one softmax layer per sibling group and training under the usual softmax loss over leaf posteriors. The authors report only a marginal drop in standard classification accuracy when enforcing this tree-structured prediction, including the additional internal-node concepts. They note that the approach brings benefits, including graceful degradation on new or unknown object categories, as the network is still capable of high confidence in a parent class when unsure as to which of its children is correct.

Since the model outputs conditional probabilities instead of class probabilities, we changed the output dimension of the terminal fully-connected layer, such that it outputs logits for every node in the hierarchy. Proper normalisation of the conditional probabilities is then enforced at every node of the hierarchy using the softmax function. Finally, the loss is computed by summing the individual cross-entropies of the conditional probabilities on the path connecting the ground-truth label to the root of the tree.

DeViSE. Frome *et al.* [14] proposed DeViSE with the aim of both making more semantically reasonable errors and enabling zero-shot prediction. The approach involves modifying a standard deep classification network to instead output vectors representing semantic embeddings of the class labels. The label embeddings are learned through analysis of unannotated text [24] in a separate step, with the classification network modified by replacing the softmax layer with a learned linear mapping to that embedding space. The loss function is a form of ranking loss which penalises the extent of greater cosine similarity to negative examples than

positive ones. Inference comprises finding the nearest class embedding vectors to the output vector, again under cosine similarity.

Since an official implementation of DeVISE is not available to the public, we re-implemented it following the details discussed in the paper [14]. Below the list of changes we found appropriate to make.

- For the generation of the word embeddings, instead of the rather dated method of Mikolov *et al.* [24], we used the high-performing and publicly available⁴ fastText library [6] to obtain word embeddings of length 300 (the maximum made available by the library).
- Instead of a single fully-connected layer mapping the network output to the word embeddings, we used the network “head” described in Listing 1. We empirically verified that this configuration with two fully-connected layers outperforms the one with a single fully-connected layer. Moreover, in this way the number of parameters of DeVISE roughly matches the one of the other experiments, which have architectures with a single fully-connected layer but a higher number of outputs (608, equivalent to the number of classes of *tieredImageNet-H*, as opposed to 300, the word-embedding size).
- Following what described in [14], we performed training in two steps. First, we trained only the fully-connected layers for the first 150,000 steps with a learning rate of $1e-4$. We then trained the entire network for 150,000 extra epochs, using a learning rate of $1e-6$ for the weights of the backbone. Note that [14] did not specify neither how long the two steps of training should last nor the values of the respective learning rates. To decide the above values, we performed a small hyperparameter search.
- [14] says that DeVISE is trained starting from an ImageNet-pretrained architecture. Since we evaluated all methods on *tieredImageNet-H*, we instead initialised DeVISE weights with the ones of an architecture fully trained with the cross-entropy loss on this dataset. We verified that this obtains better results than starting training from ImageNet weights.

Barz&Denzler [4]. This approach involves first mapping class labels into a space in which dot products represent semantic similarity (based on normalised LCA height), then training a deep network to learn matching feature vectors (*before* the fully connected layer) on its inputs. There is a very close relationship to DeVISE [14], with the main difference being that here, the label embedding is derived from a supplied class hierarchy in a straightforward manner instead of via text analysis: iterative arrangement of embedding vectors such that all dot products equal respective

semantic similarities. The authors experiment with two different loss functions: (1) a linear reward for the dot product between the output feature vector and ground-truth class embedding (*i.e.* a penalty on misalignment); and (2) the sum of the preceding and a weighted term of the usual cross-entropy loss on the output of an additional fully connected layer with softmax activation, for classification. We only used (2), since in [4] it attains significantly better results than (1).

We used the code released by the authors⁵ to produce the label embeddings. To ensure consistency with the other experiments, two differences in implementation with respect to the original paper were required.

- We simply used a ResNet-18 instead of the architectures Barz & Denzler experimented with in their paper [4] (*i.e.* ResNet-110w [17], PyramidNet-272-200 [16] and Plain-11 [3]).
- Instead of SGD with warm restarts [23], we used Adam [31] with a learning rate of $1e-4$ (the value performing best on the validation set) for 300,000 steps.

C. Outputting conditional probabilities with HXE

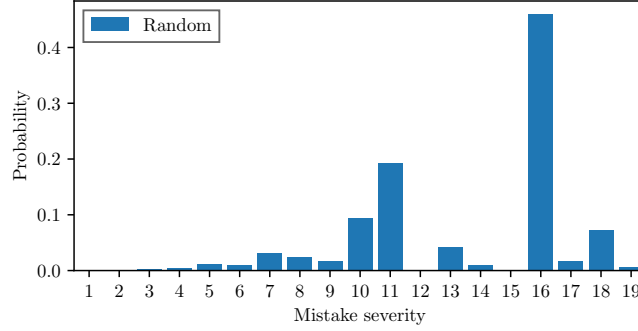
We also investigated whether outputting conditional probabilities instead of class probabilities affects the performance of the classifier represented by our proposed HXE approach (Sec. 3.1). These two options correspond, respectively, to implementing hierarchical information as an architectural change or as modification of the loss only.

Comparing different values of α for otherwise identical training parameters, we observe that outputting the class probabilities consistently results in an improvement of performance across *all* of our metrics, see Suppl. Fig. 2. Moreover, directly considering the class probabilities has also the advantage of not requiring direct knowledge of the hierarchy at test time.

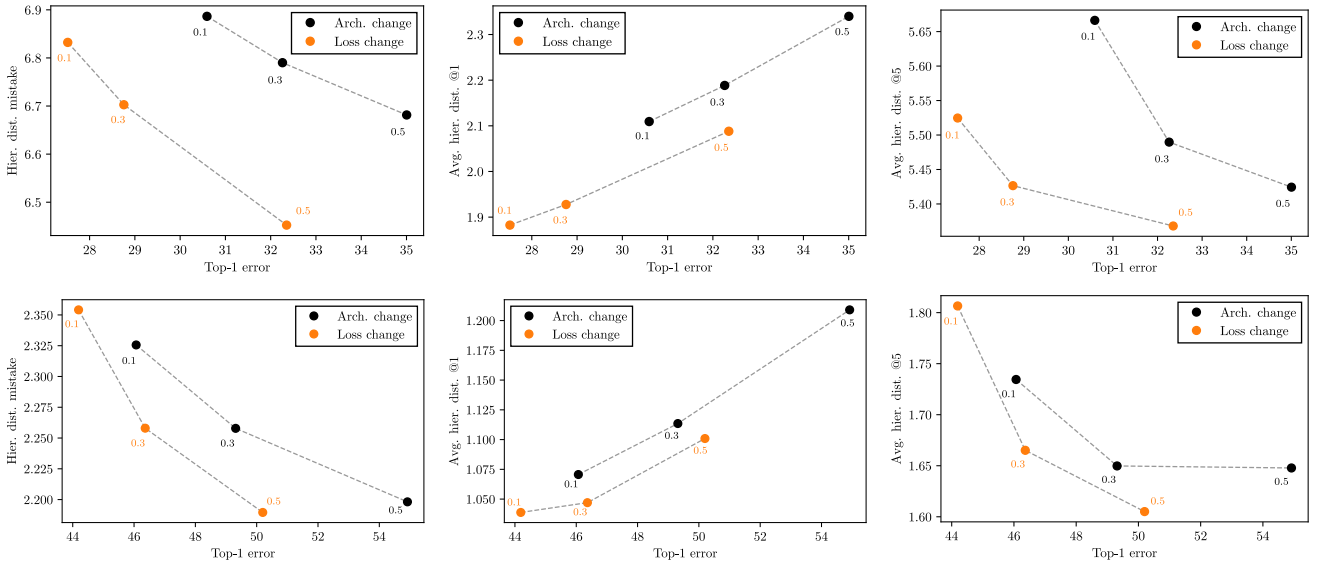
⁴<https://github.com/facebookresearch/fastText>

⁵<https://github.com/cvjena/semantic-embeddings>

D. Supplementary figures



Supplementary Figure 1: Distribution of mistake severity when picking random example pairs in *tieredImageNet-H*. Note that even though this distribution shares some similarities with the ones shown in Fig. 1, it is substantially different. This indicates that the general shape of the distributions of the mistake severities for the various DNN architectures investigated here cannot be explained by properties of the dataset alone.



Supplementary Figure 2: Outputting the conditional probabilities (Arch. change) results in a degradation of performance compared to outputting the class probabilities directly (Loss change) when using the hierarchical cross-entropy loss with exponential weights $\lambda(C) = \exp(-\alpha h(C))$. Results are shown both on *tieredImageNet-H* (top) and *iNaturalist19-H* (bottom). Points closer to the bottom-left corner of the plots are the ones achieving the best tradeoff.