

Received 28 February 2023, accepted 22 March 2023, date of publication 31 March 2023,
date of current version 5 April 2023.

Digital Object Identifier 10.1109/ACCESS.2023.3263472

RESEARCH ARTICLE

An Attention-Based Architecture for Hierarchical Classification With CNNs

IVÁN PIZARRO¹, RICARDO ÑANCULEF¹, AND CARLOS VALLE²

¹Department of Informatics, Universidad Técnica Federico Santa María, Valparaíso 2390123, Chile

²Department of Data Science and Informatics, University of Playa Ancha, Valparaíso 2360072, Chile

Corresponding author: Iván Pizarro (ivan.pizarro@usm.cl)

The work of Carlos Valle was supported by the National Agency for Research and Development (ANID) Chile through the National Fund for Scientific and Technological Development (FONDECYT) Iniciación Project under Grant 11230351.

ABSTRACT Branch Convolutional Neural Nets have become a popular approach for hierarchical classification in computer vision and other areas. Unfortunately, these models often led to hierarchical inconsistency: predictions for the different hierarchy levels do not necessarily respect the class-subclass constraints imposed by the hierarchy. Several architectures to connect the branches have arisen to overcome this limitation. In this paper, we propose a more straightforward and flexible method: let the neural net decide how these branches must be connected. We achieve this by formulating an attention mechanism that dynamically determines how branches influence each other during training and inference. Experiments on image classification benchmarks show that the proposed method can outperform state-of-the-art models in terms of hierarchical performance metrics and consistency. Furthermore, although sometimes we found a slightly lower performance at the deeper level of the hierarchy, the model predicts much more accurately the ground-truth path between a concept and its ancestors in the hierarchy. This result suggests that the model does learn not only local class memberships but also hierarchical dependencies between concepts.

INDEX TERMS Attention mechanisms, deep learning, hierarchical classification.

I. INTRODUCTION

In many pattern classification applications, class labels can be organized into a hierarchical taxonomy [1], [2]. For example, to help shoppers find products easily, e-commerce platforms such as Amazon maintain a detailed product taxonomy in which categories (e.g., Electronics Bags & Cases) branch into more specific sub-categories (e.g., Laptop Computer Briefcases). Other applications in which hierarchies of this type arise include sentiment analysis [3], web content categorization [4], disease detection [5], and gene function prediction [6], [7]. Hierarchical classification methods are devised to organize data into the hierarchy while respecting and exploiting the class relationships encoded by the taxonomy.

Many methods exist that adapt traditional machine learning algorithms for hierarchical classification. They include top-down approaches, which train different classifiers per node or level of the hierarchy [8], and global approaches, which

train a single classifier for the entire hierarchy [9], [10]. Unfortunately, only a few approaches address this task using deep learning methods such as convolutional neural nets (CNNs) [11], [12]. As these models become state-of-the-art in more and more classification problems [13], methods to embed class hierarchies into deep architectures may become increasingly useful and necessary.

The popular approach of augmenting CNNs to include *branches* that support hierarchical classification was first introduced in [14] and then extended in [15]. In this model, referred to as Branch Convolutional Neural Network (B-CNN), a *branch* is a fully-connected subnet that receives a feature representation from a main convolutional block and computes predictions for a specific hierarchy level. The main advantage of this approach is the implementation's simplicity and generality: it can accommodate any standard architecture as the central building block. The main disadvantage is that predictions per level are often inconsistent with the hierarchy: the class predicted for a given level may not be an ancestor of the classes predicted for lower levels. This

The associate editor coordinating the review of this manuscript and approving it for publication was Sudipta Roy¹.

limitation has motivated many refinements of the original model. For instance, [16] proposed an architecture in which skip connections allow hidden features from a branch to propagate directly to the next branch. This way, the prediction for a given level is explicitly conditioned on the previous one. More recently, [17] proposed to connect a branch with all the previous branches in the network. Finally, in [18], the authors used a similar connectivity pattern, but instead of propagating hidden representations, they propagate probability estimates from each branch. Although all these methods slightly improve hierarchical classification metrics, we empirically show that their improvements often came from the more specific level of the hierarchy only. Therefore, these models can still produce many predictions inconsistent with the hierarchy at inference time. This result shows that current CNNs do not accurately learn hierarchical class relationships.

This work highlights two limitations of current CNNs for hierarchical classification. First, existing approaches are limited to top-down connectivity patterns: features propagate unidirectionally from the top to deeper branches. This way, predictions at detailed levels of the class hierarchy do not explicitly give feedback to the upper levels. In addition, the connectivity pattern is static; it cannot change during training or inference. We hypothesize that a more flexible communication pattern can improve the ability of CNNs to learn hierarchical relationships and, thus, hierarchical consistency.

Inspired by recent advances in deep learning, this work proposes BA-CNN, a method that extends the original B-CNN model by connecting its branches through an attention mechanism. The main motivation for our method is that Attention [19] currently allows models to selectively and dynamically aggregate information from different parts of the computational graph with a minimal computational burden. Thus, by connecting B-CNN's branches through an attention mechanism, we allow a branch to condition its predictions on context vectors of the whole class hierarchy without cumbersome hand-crafted connections. This mechanism breaks the static top-down architecture of current models. In particular, fused vectors can include features from lower and deeper branches and, therefore, from coarser and deeper levels of the class hierarchy. Moreover, by adjusting the attention weights, attention allows the model to change how branches influence each other during learning and prediction.

Experiments on image classification benchmarks show that the proposed method can outperform state-of-the-art models in terms of hierarchical performance metrics, sometimes at the cost of slightly lower performance on the most specific level. Furthermore, the model's predictions reconstruct more accurately the ground-truth path between a concept and its ancestors in the hierarchy, which suggests the model does learn not only class memberships but also their dependencies. In addition, we show that BA-CNN has a memory footprint and training time close to that of a flat CNN and lower than that of more explicit and dense connectivity patterns for hierarchical classification.

The contribution of this paper is threefold:

- First, we propose a novel extension of the B-CNN model [15] for deep hierarchical classification. The model uses an attention module that allows feature maps to flow in different directions of the hierarchy, top-down and bottom-up. Moreover, the attention mechanism determines how branches influence each other in a dynamic and data-driven way.
- Second, we extend the actual comparison in the literature by including four current branched CNN architectures for hierarchical classification, three benchmark datasets, and four different metrics.
- Third, we propose two custom metrics named hierarchical accuracy and hierarchical consistency. Those metrics complement the chosen state-of-the-art metrics and help to measure the ability of the model to learn hierarchical dependencies between concepts more directly.

To facilitate future research, we provide code that applies the proposed method and other four CNN architectures in hierarchical classification datasets https://github.com/IvanPizarroQ/BA_CNN.

The remainder of this paper is organized as follows. Section II provides a brief background on hierarchical classification, reviewing fundamental concepts and traditional approaches to the problem. Section III reviews other methods based on convolutional neural nets that leverage and exploit class hierarchies. The technical background of the attention mechanisms related to our formulation is presented in Section IV. Our architecture is motivated and described in Section V. Section VI explains the experimental setting and results. In this section, we also compare the proposed method with four other deep hierarchical classification architectures. Finally, we summarize our conclusions and discuss future research directions in Section VII.

II. HIERARCHICAL CLASSIFICATION

In pattern recognition, a classifier is a mapping $f : X \rightarrow Y$ between data domain X and a finite set of class labels C , representing different concepts in an application domain. Classification algorithms often assume classes are unrelated. Hierarchical classification instead copes with problems in which the classes Y organize into a *taxonomy* [1] that agglomerates classes to create more abstract concepts.

Wu et al. [20] defined such a taxonomy as a tree-structured traditional concept hierarchy defined over a partial order set $(C; <)$, where $<$ represents the “is-a” relationship. Silla and Freitas [1] defined the “is-a” relationship as asymmetric, antireflexive, and transitive:

- 1) The only one greatest element is the root of the tree R , i.e., $\exists! R \in C : \forall c \in C, c \neq R \Rightarrow c < R$.
- 2) $\forall c_i, c_j \in C$ if $c_i < c_j$ then $c_j \not< c_i$.
- 3) $\forall c_i \in C, c_i \not< c_i$.
- 4) $\forall c_i, c_j, c_k \in C, c_i < c_j$ and $c_j < c_k$ imply $c_i < c_k$.

On the one hand, the existence of a taxonomy for the classes of interest represents an opportunity to improve

classification performance. As concepts at the lower levels of the hierarchy are a specialization of concepts at the upper levels, it should be easier to classify an instance by following the taxonomy. The classifier can first discriminate between simpler/coarser concepts, discarding a large set of possible outcomes, and then focus on the specific features that make the instance part of a sub-category. For example, knowing that *birds* cannot be *felines* can help the classifier discard feline-specific features when the *bird* category is much more likely than the *mammal* category. On the other hand, the existence of a taxonomy presents a learning challenge. If an instance belongs to a given class, it automatically belongs to all its superclasses, and the classifier should respect this *hierarchy constraint*. For example, if the model classifies an instance as a *feline*, it should classify the instance as a *mammal* and not as a *bird*.

Methods to leverage and respect concept hierarchy differ from three main criteria [21], [22]. The first is the hierarchical structure used, either a tree or DAG, as depicted in Fig. 1. The main difference between using a DAG instead of a tree is that a node can have more than one parent node in the DAG (a category may belong to multiple superclasses). The second criterion is the completeness of the decisions made by the classifier, i.e., whether the method can stop the classification at any node of the class hierarchy (non-mandatory leaf-node prediction) or is constrained to stop at a leaf (mandatory leaf-node prediction). The third criterion is related to how the classifier explores the hierarchical structure. The most straightforward approach ignores the class hierarchy, typically predicting only classes at leaf nodes. Methods in this category are called *flat classifiers*. Another way to achieve this is to learn the whole class hierarchy using a single classifier. Methods in this category are called *global* (or big-bang) classifiers. A third possible scenario involves learning the taxonomy class using a set of local classifiers. These methods are also known as top-down classifiers [1].

A. FLAT CLASSIFICATION

This method is also known as the direct approach [23] or bottom-up approach [24]. It completely ignores the class hierarchy, typically predicting only the classes at leaf nodes. Therefore, this approach acts as a traditional classification algorithm during training and inference (testing). Despite neglecting the hierarchy, this approach can indirectly lead to hierarchical classifications in some cases. Indeed, assuming the “is-a” relationship, we can assign to an instance all the ancestor classes corresponding to the leaf class predicted by the model. It is worth noting that the latter applies only to a tree hierarchy and if the taxonomy is perfectly known. To illustrate this approach, Fig. 2 shows the use of a flat multi-class classification algorithm.

The main disadvantage of flat classifiers is that they do not exploit possible correlations between the target concepts, which help the learner to improve generalization such as multi-task learning and transfer learning benefit from

learning inter-related tasks jointly [25]. Furthermore, theoretical studies have revealed that neglecting the hierarchical structure is especially unfavorable in large-scale multi-class scenarios, where a flat classifier must simultaneously discriminate between many unbalanced categories [26].

B. GLOBAL CLASSIFIERS

As Fig. 3 shows, algorithms in this category train a single (usually complex) classification model to learn the class hierarchy. For example, Labrou and Finin [27] proposed a text-mining classifier to address *Yahoo!*'s hierarchical categories. It learns the hierarchy using a set of topic prototypes and a classification method resembling Rocchio's document categorization approach [28]. For inference, this method first computes the similarity of a test document with each topic. The method then classifies the document in the corresponding topic if the similarity is beyond a certain threshold. Unfortunately, the authors do not suggest a method to determine this hyperparameter's value automatically. Kiritchenko et al. [9], [10] considered the hierarchical class problem a multi-label classification problem. During the training process, the method expands the label set of all the training examples with their corresponding ancestor labels and treats them as different possible outcomes. This approach is naturally prone to hierarchical inconsistencies, so the authors considered a post-processing stage that considers all outputs to ensure that the hierarchical constraints are respected. It is worth noting that the output of a global classifier might be easier to interpret than that of a local classifier because the complexity of the decision procedure implemented by the former is often lower. For example, the experiments reported in [29] showed that the number of rules generated by the global approach was much smaller than the number of rules generated by the local approach. In addition, the global classifier approach does not suffer from the significant drawback of the local classifier approach, namely, the fact that a misclassification at a given class level is propagated to the lower levels of the class hierarchy.

C. LOCAL CLASSIFIERS

These approaches handle the hierarchy using a local top-down approach. For a new example in the test set, the algorithm first predicts the coarsest level and then uses this prediction to narrow the choices at the following (finer) level. This procedure is applied recursively until the classifier reaches a leaf node. The major drawback of this approach is that misclassification at a given class level propagates to the lower levels of the class hierarchy. According to Silla and Freitas [1], there are three main approaches to implementing a classifier of this type:

- 1) Local classifier per node.
- 2) Local classifier per parent node.
- 3) Local classifier per level.

The local classifier per node approach is the most used in the literature, and it consists of training a binary classifier for each node of the class hierarchy except the root. For

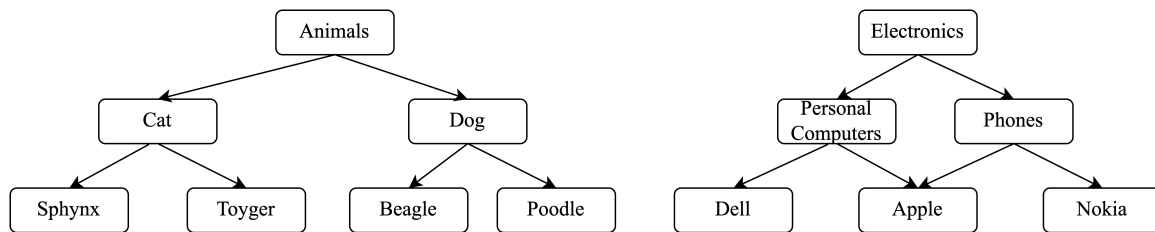


FIGURE 1. Hierarchical classification structures: Tree (left) and DAG (right).

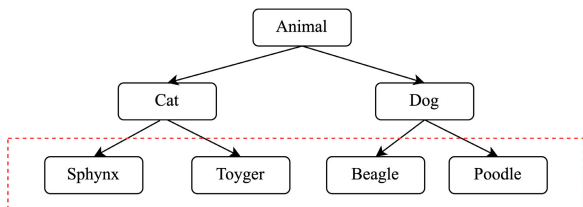


FIGURE 2. Local classifier approach using a multi-class classification algorithm.

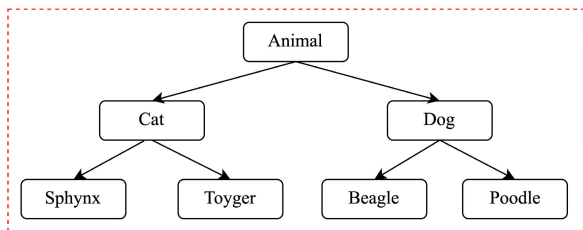


FIGURE 3. Example of the global classifier approach.

training, there are several ways to define each binary classifier's positive and negative examples. Then, for inference, each classifier predicts whether or not an example belongs to the class corresponding to the actual node. An advantage of this approach is that it allows assigning more than one label per level, as is required for multi-label hierarchical classification problems. However, a significant disadvantage of this method is that prediction errors in the nodes can generate inconsistent predictions, in other words, predictions that do not respect the class taxonomy. Therefore, due to the above, methods with this approach must have some way of correcting inconsistencies. The local classifier per parent node approach implements a multi-class classifier per parent node, trained to distinguish between its child nodes. It is worth mentioning that this method is also prone to inconsistencies if a post-processing method to correct the predictions is not applied. Finally, the local classifier per level approach consists of implementing a multi-class classifier per level of the class taxonomy.

D. HIERARCHICAL CONSISTENCY

We expect that a hierarchical classification algorithm produces predictions that respect class hierarchy. This property is referred to as *hierarchical consistency*. Wu et al. [20] defines

it as: "A label set C_i assigned to an instance d_i is called consistent with a given hierarchy if C_i forms a connected proper subgraph of the hierarchy graph rooted in the top node." For example, in the hierarchy in Fig. 4, the correct label set for the instance $d_1 = Nokia$ is {"Electronics", "Phones", "Nokia"}. The label set {"Electronics", "Phones", "Dell"} is called a hierarchical inconsistency or a class-membership inconsistency [1]. As the hierarchical structures are usually Tree or DAG, we can exclude the root node from any ancestor set since it does not provide additional information.

E. HIERARCHICAL METRICS

Most current related works evaluate hierarchical classification models using flat measures such as accuracy, precision, and recall. Thus, to measure the level-wise performance, we used the flat metric *accuracy* per level of the hierarchy. However, these measures do not consider the relations between different levels of the hierarchy. Many hierarchical performance metrics have been proposed in the literature [30]. From these, we used Kiritchenko's metrics [9] because these metrics are suitable for more different hierarchical classification problems, and Silla and Freitas recommended his usage in [1]. In addition, we propose two custom metrics: hierarchical consistency to evaluate the consistency of the predictions and hierarchical accuracy.

1) HIERARCHICAL PRECISION, RECALL, AND F-SCORE

To measure the performance of hierarchical classification, Kiritchenko et al. [9] proposed a hierarchical version of the flat metrics precision, recall, and F-score. The authors called the metrics hP (Hierarchical precision), hR (Hierarchical recall), and hF -score (Hierarchical F-score). Formally, let C_i be the true label set for an instance d_i , i.e., the set of elements composed of the ground-truth label for the finest hierarchy level and all its ancestors. On the other hand, let \hat{C}_i be the set of labels predicted for an instance d_i at each level of the class hierarchy. Therefore, hP and hR are computed as follows:

$$hP = \frac{\sum_i |C_i \cap \hat{C}_i|}{\sum_i |\hat{C}_i|}, \quad hR = \frac{\sum_i |C_i \cap \hat{C}_i|}{\sum_i |C_i|}. \quad (1)$$

Then we can combine the two values into one hF -measure:

$$hF_\beta = \frac{(\beta^2 + 1) \cdot hP \cdot hR}{(\beta^2 \cdot hP + hR)}, \quad \beta \in [0, +\infty]. \quad (2)$$

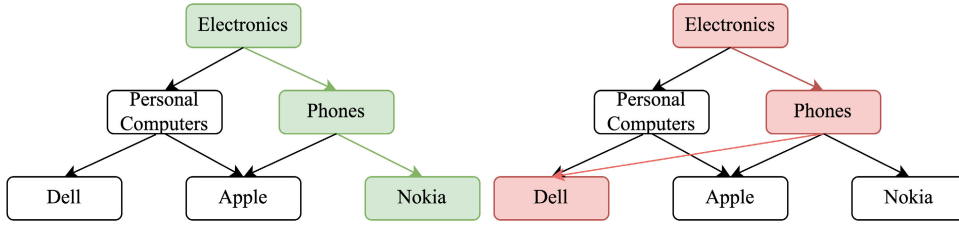


FIGURE 4. Example of a hierarchical consistent label path (left) and a hierarchical inconsistency (right).

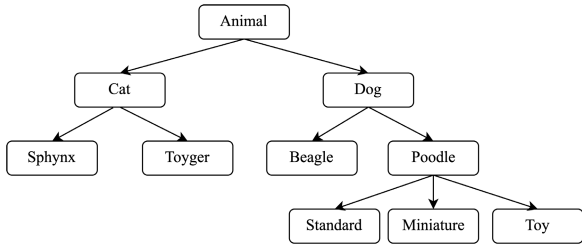


FIGURE 5. Example of a taxonomy to exemplify the use of hierarchical metrics.

In particular, authors recommend $\beta = 1$, giving precision and recall equal weights.

$$hF_1 = \frac{2 \cdot hP \cdot hR}{hP + hR}. \tag{3}$$

These metrics fulfilling the following requirements formulated by the authors:

- 1) The measure gives credit to partially correct classification.
- 2) For non-mandatory leaf node prediction problems:
 - a) The measure gives a higher evaluation for correctly classifying one level down compared to staying at the parent node.
 - b) The measure gives a lower evaluation for incorrectly classifying one level down compared to staying at the parent node.
- 3) The measure punishes errors at higher levels of a hierarchy more heavily.

To exemplify the use of this metric, Table 1 shows five instances associated with the taxonomy illustrated in Fig. 5. It is worth noticing that if all the classes have labels for every level in the taxonomy (full labeling) and the task is mandatory leaf node prediction (the model always assigns leaf classes), the cardinality of the true label sets will be equal to the cardinality of the predicted labels, so hR , hP , and hF_1 -score will have the same numerical value.

Then Kiritchenko’s hierarchical metrics are computed as follows:

$$hP = \frac{\sum_i |C_i \cap \hat{C}_i|}{\sum_i |\hat{C}_i|} = \frac{2 + 1 + 1 + 1 + 2}{2 + 2 + 2 + 2 + 3} = \frac{7}{11} = 0.64, \tag{4}$$

$$hR = \frac{\sum_i |C_i \cap \hat{C}_i|}{\sum_i |C_i|} = \frac{2 + 1 + 1 + 1 + 2}{2 + 2 + 2 + 3 + 3} = \frac{7}{12} = 0.58, \tag{5}$$

$$hF_1 = \frac{2 \cdot hP \cdot hR}{hP + hR} = \frac{2 \cdot 0.64 \cdot 0.58}{0.64 + 0.58} = 0.61. \tag{6}$$

2) HIERARCHICAL ACCURACY

In addition to Kiritchenko’s metrics, we propose to adapt the standard accuracy metric to hierarchical classification. For this modification, we only consider an object as correctly classified if the entire prediction set for that object is correct. For example, in Table 1, we can notice that for five instances, only one is correctly classified (d_1), so hierarchical accuracy is computed as follows:

$$\text{Hierarchical accuracy} = \frac{1}{5} = 0.2. \tag{7}$$

3) HIERARCHICAL CONSISTENCY

As argued in [9] and [1], some models are prone to inconsistent class predictions across different levels. This problem is called hierarchical inconsistency or class-membership inconsistency.

Let G be the set of all correct paths from the root to a leaf in the class hierarchy (if the hierarchy is a tree, $|G|$ = the number of class leaves). Then, to verify the hierarchical consistency, we propose to verify that the predicted label set \hat{C}_i for every instance d_i is in G . This metric does not give credit to partially consistent classifications and only considers a prediction as consistent if \hat{C}_i exists in the taxonomy. It is worth noting that hierarchical accuracy is always lower than or equal to hierarchical consistency because a prediction can respect the hierarchy but be incorrect for a particular input case (d_3 is an example).

For the example presented in Table 1, $G = [\{\text{“Cat”}, \text{“Sphynx”}\}, \{\text{“Cat”}, \text{“Toyger”}\}, \{\text{“Dog”}, \text{“Beagle”}\}, \{\text{“Dog”}, \text{“Poodle”}\}, \{\text{“Dog”}, \text{“Poodle”}, \text{“Standard”}\}, \{\text{“Dog”}, \text{“Poodle”}, \text{“Miniature”}\}, \{\text{“Dog”}, \text{“Poodle”}, \text{“Toy”}\}]$. Then the proposed hierarchical consistency metric (hC) is computed as follows:

$$hC = \frac{\text{No. Consistent Predictions}}{\text{No. Instances}} = \frac{2}{5} = 0.4. \tag{8}$$

The previous examples show that while Kiritchenko’s metrics present values close to 0.6, the custom hierarchical consistency and hierarchical accuracy metrics get values of 0.2 and 0.4, respectively. These new metrics will complement the Kiritchenko metric, providing a more severe penalty for errors since they do not penalize entirely correct predictions for the case of hierarchical accuracy and completely correct relationships for the case of hierarchical consistency.

TABLE 1. Example of the hierarchical metrics calculation.

| Instance d_i | True Labels C_i | Predicted Labels \hat{C}_i | $ C_i \cap \hat{C}_i $ | $ C_i $ | $ \hat{C}_i $ | $\hat{C}_i \in G$ |
|-------------------|--------------------------|---------------------------------|------------------------|---------|---------------|-------------------|
| d_1 | {“Cat”, “Toyger”} | {“Cat”, “Toyger”} | 2 | 2 | 2 | 1 |
| d_2 | {“Dog”, “Beagle”} | {“Dog”, “Sphynx”} | 1 | 2 | 2 | 0 |
| d_3 | {“Cat”, “Sphynx”} | {“Cat”, “Toyger”} | 1 | 2 | 2 | 1 |
| d_4 | {“Dog”, “Poodle”, “Toy”} | {“Dog”, “Toyger”} | 1 | 3 | 2 | 0 |
| d_5 | {“Dog”, “Poodle”, “Toy”} | {“Dog”, “Beagle”, “Toy”} | 2 | 3 | 3 | 0 |

III. RELATED WORK

Yan et al. [14] were the first to introduce the idea of extracting information from the hierarchical structure of a Convolutional Neural Network to enrich the prediction quality. The authors proposed the Hierarchical Deep Convolutional Neural Network (HD-CNN) for image recognition. This approach uses a specialized CNN for image recognition as the base (building block). Next, it uses a classifier for the coarse level that contributes to the finest level. One disadvantage of this method is that its architecture only accepts a hierarchy of two levels. Subsequently, Zhu and Bain [15] presented the Branch Convolutional Neural Network (B-CNN). They were inspired by the idea that the first layers of the CNN obtain information from higher-level layers. As shown in Fig. 6, this method uses an existent convolutional net (e.g., VGG16) as a central feature extractor (common backbone) and generates many branch neural networks as levels have the class hierarchy. Each branch neural network is a feedforward neural network that predicts a given level of the hierarchy. Next, a final loss function computes the weighted sum of the loss functions of each branch. They also presented a training method for this type of structure called *Branch Training Strategy* (BT-Strategy). It consists of modifying the weights of the branch losses while training the network. Thus, to improve performance, a weight update process is performed from the lower-level parameters to the higher-level parameters.

Various authors have implemented the B-CNN to solve different problems. For example, Seo and Shin in [31] applied the B-CNN model in the context of retail with the Fashion MNIST dataset, grouping the ten classes of the dataset into six superclasses and then, in turn, grouping these superclasses into two classes. To implement this, the authors used VGG16 and VGG19 as the central feature extractor. In addition, Sali et al. [32] employed the B-CNN model to classify gastrointestinal disorders on histopathological images using a two-level hierarchy.

To that moment, branches were treated independently, but authors such as Inoue et al. [16] and Zhang et al. [17] demonstrated the benefit of interconnecting branches. Indeed, they showed that branches can complement information from different hierarchy levels during training. Inoue et al. proposed the so-called Concat-net and Add-net, which consists of connecting the last dense layer of the branches in a top-down manner. For example, the second branch would receive the last activation tensor from branch 1, and branch three would receive that from branches 1 and 2. This interconnection

concatenates these layers for the Concat-net model and adds them in the Add-net model, demonstrating advantages over the approach of using the branches independently (B-CNN). Zhang et al., similar to Inoue et al., also propose to connect the branches in a model called Hierarchical Bilinear Convolutional Neural Network (HB-CNN). As shown in Fig. 7, this was performed from the first dense layer of the branch. This connection is called Connectivity Pattern (CP), and the authors showed that it is beneficial not only to extract one branch per level but it can be more than one; in their experiments, they extracted three and five branches. Furthermore, it uses the Bilinear CNN (B-CNN) proposed by Lin et al. in [33] to provide further enhancement to the fine level. The authors presented HB-CNN without the Bilinear CNN as Hierarchical Convolutional Neural Network (H-CNN). All the aforementioned authors applied the BT-Strategy for training, which entails an exhaustive process of adjusting the weights of the different loss functions and estimating the number of times to apply this change.

All works mentioned above present two main limitations. First, features between branches propagate unidirectionally from coarse to fine levels, causing the fine level to benefit most with this connection and not consider bottom-up feedback. The second limitation is that the connectivity patterns between branches can not change during training or inference, so it has to be carefully tuned and assuming that said interconnection is beneficial for all classes, denying the possibility that for some, the most beneficial connection is different.

Other Deep Learning approaches presented in the literature do not exploit the hierarchical structure of the features extracted by the CNN, i.e., they do not make correspondences between levels of the CNN and levels of the hierarchy to enrich the quality of predictions. For instance, Kolisnik et al. [18] trained a modified VGG16 architecture with teacher forcing, using the true labels of a higher level to train lower levels. They validated their results on the Kaggle Fashion Product Images data set [34]. La Grassa et al. [35] proposed an architectural extension that could be adapted to generic neural networks. They chose a base model and added a set of neural layers equal to the number of levels in the hierarchy. Each new layer is associated with a cross-entropy loss function and then computed a center loss function. In the medical field, Kowsari et al. [36] proposed the Hierarchical Medical Image classification (HMIC) model that uses one CNN architecture for each parent node in the hierarchy. They

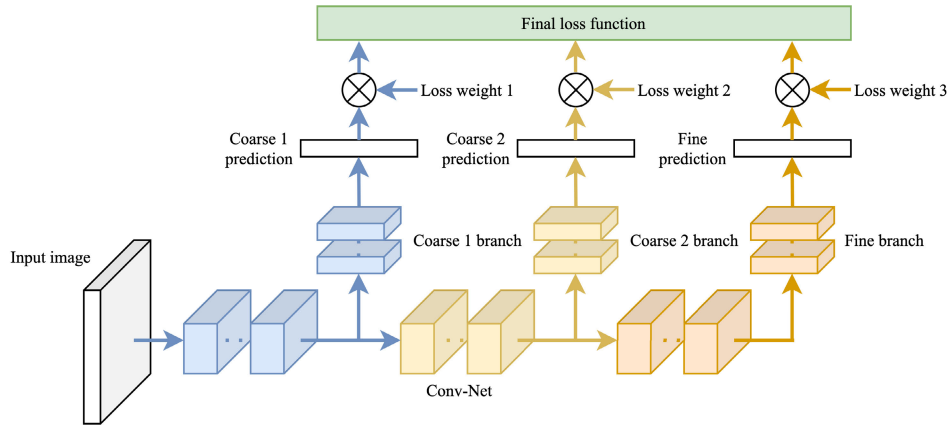


FIGURE 6. B-CNN architecture.

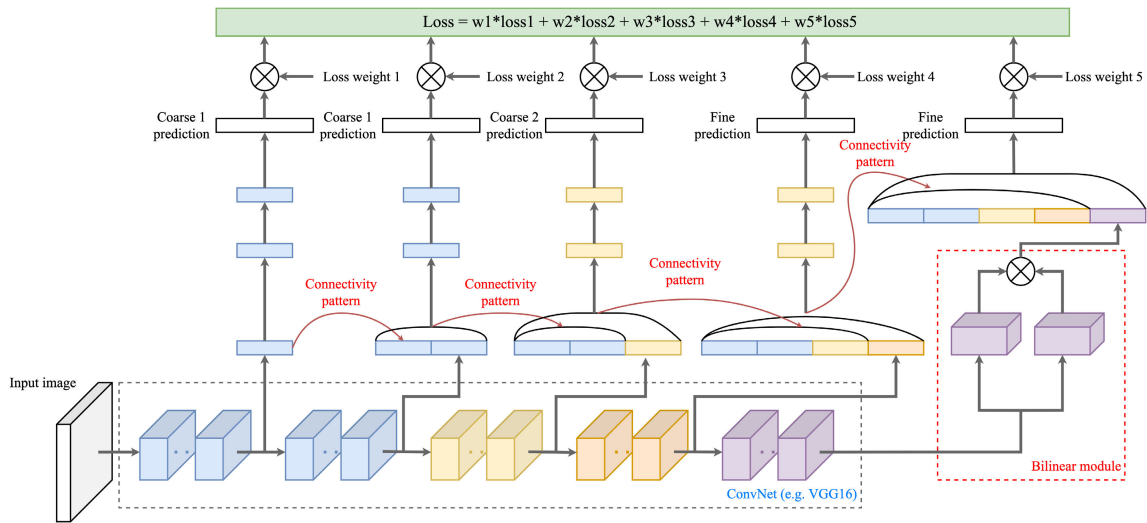


FIGURE 7. HB-CNN architecture.

validated this idea using a medical dataset in which the first level has three classes (Normal, Environmental Enteropathy, and Celiac Disease), and the child level of Celiac Disease is based on severity (I, IIIa, IIIb, and IIIc). Gao [37] proposed a deep hierarchical classification framework, tested on text and images, composed of three parts: a Feedforward Neural Network, a Hierarchical Embedding Network, and a Hierarchical Loss Network. The first is to obtain a root representation; the second is composed of a representation per level, where the representation of a level is computed by concatenating the representations of all previous levels; and the third consists of a dependence loss to punish if the model does not predict the classes according to the hierarchy, and computes a total loss by adding the losses per level.

IV. ATTENTION MECHANISMS

Attention is a recent technique for improving deep learning models that have shown promising results in Natural Language Processing [19], Speech Recognition [38], and

Computer Vision [39]. Inspired by the ability of humans to selectively concentrate on parts of the information when processing large amounts of data, attention mechanisms allow a model to focus on salient features of the input data or its internal representations to solve a task.

Formally, given a sequence of vectors (or more generally tensors) $V = \{v_1, v_2, \dots, v_M\}$ produced at a given level of the neural net, an attention mechanism computes a new set of feature representations $A = \{a_1, a_2, \dots, a_N\}$ by recombining the elements of V as follows:

$$a_n = \sum_{m=1}^M \alpha(q_n, k_m) v_m, \quad (9)$$

where k_m is a vector referred to as a key, q_n is a vector referred to as a query, and $\alpha(q_n, k_m) \in [0, 1]$ is the attention weight assigned by the attention mechanism to the value v_m . In this abstraction, the query encodes an information need or task, while the key k_m is a descriptor of the information contained in v_m . The attention weight $\alpha(q_n, k_m)$ determines

the relevance of the information source v_m for the task q_n and is computed by applying the Softmax function to the scores vector ξ_{nm} . For instance, the values v_m may correspond to the latent vectors computed by an encoder recurrent neural network (RNN) fed with input text, and the queries q_n may correspond to the latent vectors computed by a decoder RNN attempting to translate the text into another language. The keys may correspond to the values themselves or their projection into a space in which the languages can be more easily aligned. In this example, α_{nm} quantifies the attention that the model needs to give the m -th input word to predict the n -th output word.

Different attention models differ in the implementation of the scoring function that computes ξ_{nm} . Bahdanau et al. [40] proposed an additive attention mechanism in which the scores vector was first computed as

$$\xi_{nm} = \tanh \left(W^{(1)}q_n + W^{(2)}k_m \right), \quad (10)$$

were $W^{(1)}$ and $W^{(2)}$ are learnable matrices. Then, a Softmax function is applied to obtain positive attention weights that sum to 1. Later, Luong et al. [41] proposed a multiplicative attention mechanism in which the scores vector is computed as

$$\xi_{nm} = q_n W k_m, \quad (11)$$

where W is a learnable matrix. It is noteworthy that, in these attention models, the keys coincide with the values. Key/value/query abstraction became popular after the introduction of Transformers [19], which came with the concept of multi-head dot-product attention. This mechanism extends the multiplicative model of (11) with two significant improvements. First, the keys, values, and queries are projected onto a subspace using the learnable matrices $W^{(k)}$, $W^{(v)}$, and $W^{(q)}$. Second, multiple attention mechanisms or *heads* can operate in parallel. Finally, the results of these K *heads* are concatenated and projected onto the desired dimensionality using a learnable matrix $W^{(o)}$.

Applying the attention mechanism to internal network representations has led to the concept of *self-attention*, a method that allows some models to learn from sequences without recurrent connections [19]. Recent research has also shown that attention is a powerful method for improving the interpretability of deep-learning models [42]. Indeed, if the mechanism is strategically placed to attend to the input or internal features with clear semantics, visualizing the attention weights can help humans to interpret the model's predictions [43].

V. PROPOSED METHOD

This section presents our method to train a CNN for hierarchical classification. According to Silla's categorization framework [1], the proposed model corresponds to:

- SPP (Single Path Prediction): the model can assign at most one path of the predicted labels to each data instance.

- MLNP (Mandatory Leaf Node Prediction): the model always assigns leaf classes.
- GC (Global Classifier): one single model assigns labels for all levels.

In addition, the model does not assume that the hierarchy is a tree or a DAG and works in both cases. Indeed, the method only assumes that the classes form a top-down taxonomy composed of B different levels and does not require knowledge of the precise relationships between concepts of different levels in advance: the method learns these relationships from data.

Building on B-CNN, our model uses a multi-branch architecture that includes a central feature extractor and a series of branches or sub-models devoted to classifying the input data at each level of the hierarchy. As depicted in Fig. 8, our main contribution is an extension of B-CNN based on an attention mechanism that allows feature maps to flow in different directions of the hierarchy. In addition, the attention mechanism determines how the branches influence each other in a dynamic and data-driven way.

A. BRANCHED ARCHITECTURE

Formalizing recent approximations, if \mathbf{x} denotes a possible input to the system, the central block implements a transformation of the form $\mathbf{z}^{(L)} = F(\mathbf{x})$ constructed as a composition of L simpler transformations or "layers"

$$\begin{aligned} \mathbf{z}^{(\ell)} &= f_\ell(\mathbf{z}^{(\ell-1)}) \quad \forall \ell \in [L], \\ \mathbf{z}^{(0)} &= \mathbf{x}. \end{aligned} \quad (12)$$

Similarly, each branch implements a layered transformation $\mathbf{y}_b^{(K_b)} = G_b(\mathbf{x}_b)$ defined recursively as

$$\begin{aligned} \mathbf{y}_b^{(\ell)} &= g_\ell(\mathbf{y}_b^{(\ell-1)}) \quad \forall \ell \in [K_b], \\ \mathbf{y}_b^{(0)} &= \mathbf{x}_b \quad \forall b \in [B], \end{aligned} \quad (13)$$

where \mathbf{x}_b is a tensor specifically prepared for level b and g_ℓ is the ℓ -th transformation or layer.¹ Predictions for each level of the hierarchy can be obtained from a simple fully connected layer fed with $\mathbf{y}_b^{(K_b)}$,

$$\mathbf{y}_b = \sigma(W_b \mathbf{y}_b^{(K_b)} + \theta_b), \quad (14)$$

where σ denotes a Softmax activation function. W_b and θ_b are the parameters of the fully connected layer for level b .

A key difference between existing models is how a branch is allowed to condition the prediction of another branch. For example, in [15], the branches were not explicitly connected and depended on each other only through the shared feature extractor. The input of branch b is a feature map extracted from the central block, that is $\mathbf{x}_b = \mathbf{z}^{(\ell_b)}$ for some $\ell_b \in [L]$ such that $\ell_b > \ell_{b'}, \forall b > b'$. The output of branch b is not independent of branch $b' < b$ because x_b depends on $x_{b'}$, which is a more primitive representation of the input data.

¹These layers g_ℓ are usually fully connected (FC) layers, i.e., an affine transformation followed by an element-wise non-linear activation.

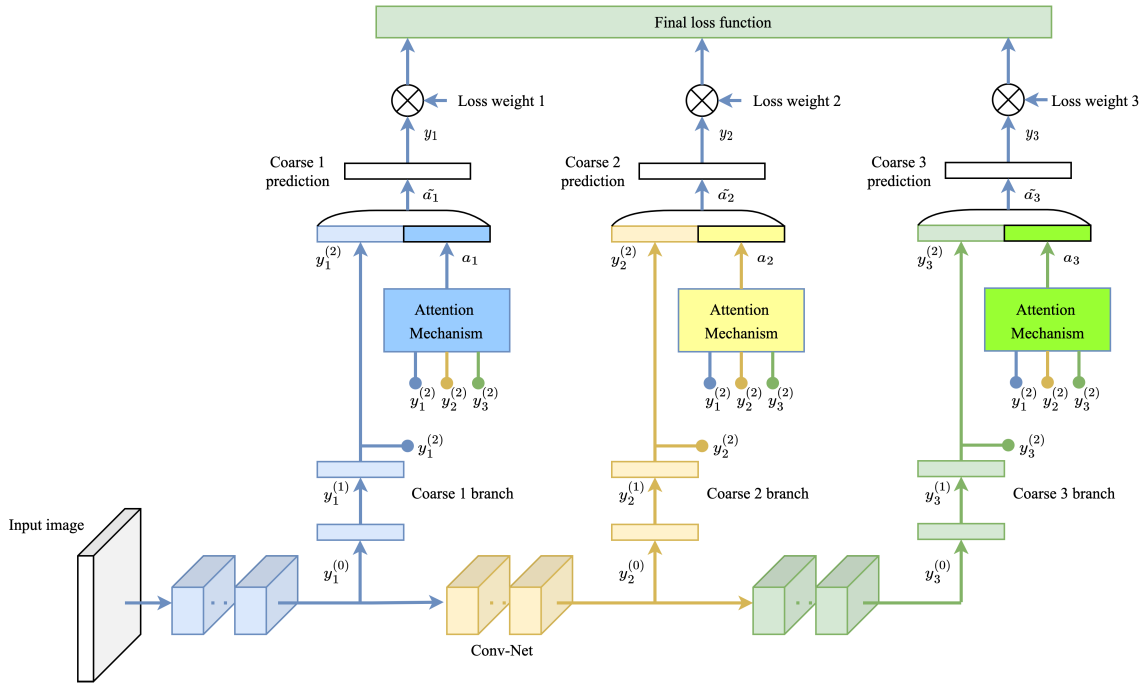


FIGURE 8. Proposed architecture BA-CNN for a 3-level hierarchy.

In [18], two contiguous branches get explicitly connected by defining x_b as the concatenation of the previous branch's output with the feature representation drawn from the central block, i.e., $x_b = z^{(\ell_b)} \oplus y_{b-1}^{(K_{b-1})}$. This architecture allows branch b to explicitly condition the prediction of branch $b + 1$ and resembles the Jordan recurrent connections used in sequence modelling. In [31], the previous branch's output is substituted by pre-output activations, i.e., $x_b = z^{(\ell_b)} \oplus y_{b-1}^{(k_{b-1})}$ for some $k_b \in [K_b]$. In [17], each branch gets explicitly connected to all previous branches by defining x_b as

$$x_b = z^{(\ell_b)} \oplus \left(\bigoplus_{\tilde{b}=1}^{b-1} z^{(\ell_{\tilde{b}})} \right). \quad (15)$$

In [16], branches are connected similarly to [17], but instead of connecting the representation extracted from the central block, they connect the pre-softmax layer $y_b^{(K_b)}$. This connection generates a new pre-softmax layer $\tilde{y}_b^{(K_b)}$, defined for Concat-net as

$$\tilde{y}_b^{(K_b)} = y_b^{(K_b)} \oplus \left(\bigoplus_{\tilde{b}=1}^{b-1} y_{\tilde{b}}^{(K_{\tilde{b}})} \right). \quad (16)$$

Add-net uses the same definition, but instead of performing the concatenation, they use the summation.

B. ATTENTION MECHANISM

Generalizing the above ideas and inspired by recent advances in sequence modelling, we propose replacing these (increasingly complex) connectivity patterns with an attention module. Our aim is to feed branch b using a context vector a_b constructed from representations $\mathcal{B} = \{y_1^{(k)}, y_2^{(k)}, \dots, y_B^{(k)}\}$

drawn from all the other branches in the network, coarser and finer ones. To this end, we propose an attention mechanism that recombines the elements of \mathcal{B} as follows:

$$a_b = \sum_{m=1}^M \alpha_{bm} y_m^{(k)}, \quad (17)$$

where α_{bm} is the attention weight assigned to branch m to compute the context vector of branch b . This mechanism is shown in Fig. 9.

As discussed in the previous section, there are many different ways to compute the weights, α_{bm} . We propose computing α_{bm} using a simple one-hidden-layer neural net with the linear activation g , that is:

$$\alpha_{bm} = \frac{\exp(\xi_{bm})}{\sum_{m'} \exp(\xi_{bm'})}, \quad (18)$$

$$\xi_{bm} = g(W_{bm} y_m^{(k)} + \beta_{bm}). \quad (19)$$

Note that each branch b attends to the representations in \mathcal{B} using the branch-specific parameters W_b and β_b . Note also that these parameters are not shared across the elements in \mathcal{B} , i.e., the model learns different parameters to obtain the logits ξ_{b1} , ξ_{b2} , and ξ_{bB} , which determine how branch b attends to branch m . It would be possible to use the same parameters for each m , but this reduces the flexibility of the attention module.² An interpretation of this approach is as

²If the attention parameters were shared among branches, i.e., $W_{bm} = W_{b'm'} \forall b, b', m \in \mathcal{B}$, each branch b would be constrained to attend the representations in \mathcal{B} with the same attention weights α_{bm} used by the other branches. Therefore, as the attention module's input is the same for every branch b , the context vectors a_b would no longer be branch-specific but global context vectors about the model's predictions for the hierarchy.

follows: to be linearly combined, the tensors in \mathcal{B} must have a consistent dimension, i.e., we need to fix the dimensionality of the feature maps flowing from the branches to the attention module. One way to compensate for this condition is to introduce a map that projects $\mathbf{y}_m^{(k)} \in \mathcal{B}$ into a feature space that disentangles the information needed by branch b about branch m . Our attention module implements that embedding using a fully connected layer with the parameters W_{bm} and β_{bm} .

For further flexibility, we employ residual connections around the attention module; that is, we concatenate \mathbf{a}_b with the original latent representation flowing from branch b to the attention module as follows:

$$\tilde{\mathbf{a}}_b = \mathbf{a}_b \oplus \mathbf{y}_b^{(k)}. \quad (20)$$

Once the attention weights have been computed, the predictions for each level of the hierarchy are obtained using a softmax-activated output layer conditioned on $\tilde{\mathbf{a}}_b$, i.e.:

$$\mathbf{y}_b = \sigma(W_b \tilde{\mathbf{a}}_b + \theta_b). \quad (21)$$

Note that Equations (12) and (13), which define the branched approach to hierarchical classification, are still valid. The fundamental difference between our method and the previous techniques is that we place an attention module between the features of (13) and the output in (14).

C. LEARNING

Provided that all the subnets in the model are differentiable, we can train the system end-to-end using backpropagation. We used the cross-entropy loss to guide the learning of each branch. That is if \mathbf{y}_b^* denotes the desired probability distribution for branch b and $\mathbf{y}_b = \mathbf{y}_b^{(k)}$ is the predicted distribution of that branch, the loss corresponding to level b in the taxonomy is computed as:

$$L_b(\mathbf{y}^*, \mathbf{y}) = -\mathbb{E}_{\mathbf{y}^*} \ln(\mathbf{y}). \quad (22)$$

In practice, the expected value in (22) is estimated using data $S = \{(\mathbf{x}^{(n)}, \mathbf{y}_b^{(n*)})\}$, which have been annotated for that level of the class hierarchy. Previous works (see, for example, [15] and [17]) often assumed that data were classified at a finer level in the hierarchy. In this case, level-wise annotations can be obtained by tracing the hierarchy back. However, note that the current formulation supports partial annotations.

The loss corresponding to the entire taxonomy is defined as the weighted sum of these losses.

$$L(\mathbf{y}_{1:b}^*, \mathbf{y}_{1:b}) = \sum_b \omega_b L_b(\mathbf{y}_b^*, \mathbf{y}_b). \quad (23)$$

VI. EXPERIMENTS

In this section, we describe the experiments conducted to evaluate the proposed method on three image classification datasets: CIFAR-10 [44], CIFAR-100 [44], and Fashion MNIST [45]. These datasets have been widely used in recent studies to assess the performances of deep hierarchical classifiers. We compared the performance of our model with

four baselines: B-CNN [15], H-CNN [17], Add-net [16] and Concat-net [16]. We also compared the performance of the hierarchical CNN models with a traditional CNN trained to predict the last level of the taxonomy. We chose B-CNN as a baseline because it is the base architecture on which most recent works are based. On the other hand, H-CNN is the most recent method proposed for this task. In addition, Add-net and Concat-net are recent architectures slightly simpler than HCNN and, thus, are worth considering.

A. DATASETS

The CIFAR-10 dataset comprises 60,000 32×32 color images organized into ten natural classes, with 6,000 images per class. We employed the three-level class taxonomy proposed by Zhu and Bain [15] and then used by [17], which included the classes *animal* and *transport* at the coarsest level. For the next level, the class *animal* separates into four sub-categories (*medium*, *pet*, *reptile*, and *bird*) and the class *transport* into three (*road*, *water*, and *sky*). Finally, the fine level accommodates the ten original classes. CIFAR-100 consists of 60,000 32×32 RGB images divided into 100 natural classes with 600 images per class. The class hierarchy used in [15] and [17] organizes the classes into three levels: the first level contains eight coarse categories (not specific names used), the second and the third level corresponds to the 20 coarse and 100 fine-grained categories originally included in the dataset [44]. Fashion MNIST comprises 70,000 28×28 gray images extracted from Zalando, an online fashion platform. In [31], the authors built a three-level taxonomy for this dataset which included two classes at the coarsest level: *clothes* and *goods*. The next level separates the class *clothes* into four sub-categories (*tops*, *bottoms*, *dresses*, and *outers*) and the class *goods* into two sub-categories (*accessories* and *shoes*), for a total of six new classes. Finally, the last level accommodates the ten original classes. We adapted our model and the baselines to learn this taxonomy.

B. EXPERIMENTAL SETUP

To implement all the methods, we adopted an architecture first introduced for hierarchical image classification in [15]. The model, inspired by the VGG model's architecture and then adopted by many studies to facilitate comparisons, includes a central feature extractor referred to as Base-C. Base-C bifurcates into different sub-models or branches devoted to predicting class assignments for the different levels of the hierarchy. In [15], authors showed that Base-C outperforms an alternative architecture named Base-B in all metrics, but in Section VI-D, we assess the effects of this choice.

Base-C includes five layer sequences referred to as *convolutional blocks*. The first two convolutional blocks consist of two convolutional layers with 3×3 filters, followed by a 2×2 max pooling layer. The last three convolutional blocks increase the number of convolutions to three, and the fifth convolutional block does not have the max pooling layer. The

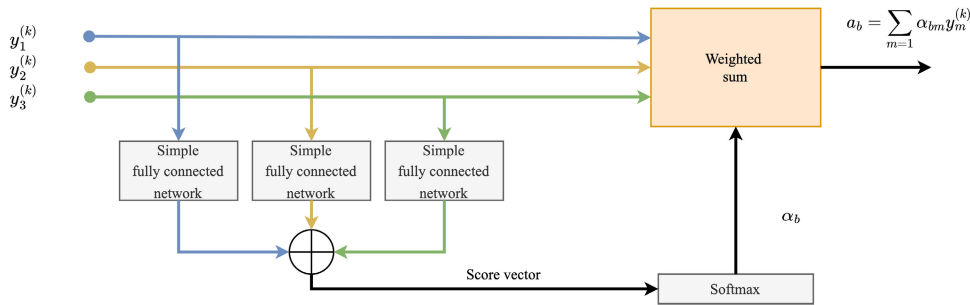


FIGURE 9. Attention mechanism for branch b in a 3-level hierarchy.

convolutional filters are 64, 128, 256, and 512 in the first, second, third, and last two blocks, respectively. The activation function is ReLU for all convolutional layers.

As the taxonomies previously used to validate the baselines have three levels, all the models added three branches to the central backbone. Each branch consisted of two fully-connected layers followed by a Softmax output layer. The number of neurons for each layer was 512, 1024, and 4096 for the first, second, and third branches, respectively. Besides level-specific branches, H-CNN, Add-net, and Concat-net add skip-connections that connect each branch with the subsequent branches. Following the author's recommendations, we implemented these skip-connections by propagating each branch's input vector to the following branches. For the comparisons, we chose VGG16 as a flat CNN representative because it is equivalent to B-CNN without the first two branches.

As discussed before, all assessed methods use a step-by-step training method named BT-Strategy, which gradually modifies the weights of the multiple losses in the objective function. This strategy encourages the model to focus on coarser hierarchy levels during the first epochs of training and gives more importance to fine-grained predictions at the end of the training. As finding the proper schedule for new datasets is difficult and time-consuming, we evaluated the impact of the BT-Strategy in all experiments.

For our attention modules, we implemented fully-connected nets with one hidden layer of 64 neurons each. As the method linearly combines the feature vectors extracted from each branch, these vectors must have a consistent dimension. Therefore, we fixed the dimensionality of the hidden layers in all three branches in our model to 256. This decision was a parsimonious choice that reduced the total number of parameters in the model. Subsequently, we studied the effect of this decision. Since Add-net also linearly combines feature vectors from each branch, we adopted the same standard for this model. In addition to the parsimonious version of our proposal, we included a version using the best hyperparameters from an exhaustive grid search performed on the model. To tune our model, we considered three hyperparameters: The number of neurons of the hidden layers in all branches in the range $[2^4, 2^{12}]$, $[2^6, 2^{12}]$, and

$[2^3, 2^{12}]$ for CIFAR-10, CIFAR-100 and Fashion MNIST respectively, the number of neurons of the hidden layers in the attention mechanism in the range $[2^5, 2^{11}]$, and the usage of BT-Strategy. For all combinations, we calculated all the performance metrics using 5-fold cross-validation. Table 2 shows the optimal settings.

We conducted two statistical tests to evaluate the significance of the experimental results. First, we employed Friedman's test to assess the (null) hypothesis that the methods we compared were statistically equivalent under a given metric. In this design, the method serves as the group variable, and the level of supervision serves as the blocking variable. Second, when rejecting the null hypothesis of Friedman's test, we compared the proposed method with the other algorithms using the Wilcoxon test with Bonferroni correction to check for pairwise differences. Note that the Bonferroni correction yields much more conservative p -values than those obtained by assuming independence of the pairwise differences. For all the tables, the highlighted values correspond to the maximum per category as long as it has a statistically significant difference at 5% level from the value that follows it. If there is more than one highlighted value, it is because the statistical tests show them as equivalent, showing both as the best value.

All the models were implemented using Python with the functional Deep Learning API Keras, running on a desktop computer, composed of an AMD(R) Ryzen 7(R) @3.80GHz processor and an Nvidia Geforce RTX 3060ti GPU. To facilitate future research, all the code required to reproduce our model and the baselines have been made available at https://github.com/IvanPizarroQ/BA_CNN.

The specific experimental settings for each dataset are the following:

1) CIFAR-10

We trained all the models using the settings of [15] and [17]. Stochastic Gradient Descent (SGD) was iterated for 60 epochs using a batch size of 128. The learning rate schedule started with a value of 0.003, decreased to 0.0005 after 42 epochs, and to 0.0001 after 52 epochs. For the BT-Strategy, we start with weights of 0.98, 0.01, and 0.01 for the first, second, and third hierarchy levels, respectively. These weights were modified to 0.1, 0.8, and 0.1 after 10 epochs; to 0.1,

TABLE 2. Best hyperparameters after the exhaustive grid search for CIFAR-10, CIFAR-100 and Fashion MNIST.

| Dataset | Branch No. Neurons | Mechanism No. Neurons | BT-Strategy |
|---------------|--------------------|-----------------------|-------------|
| CIFAR-10 | 32 | 2048 | True |
| CIFAR-100 | 256 | 2048 | False |
| Fashion MNIST | 16 | 2048 | True |

0.2, and 0.7 after 20 epochs; and 0, 0, and 1 after 30 epochs. For the experiments without the BT-Strategy, we set the weights to 0.33, 0.33, and 0.34 during the entire training process.

2) CIFAR-100

The branches used to predict the different levels of the hierarchy were organized following [15]. We trained all the models for 80 epochs using SGD and batch size of 128. The learning rate starts with a value of 0.001. At epoch 55, it decreases to 0.0002. Finally, at epoch 70, it is updated to 0.00005. The schedule of weights for the BT-Strategy was the same as for CIFAR-10.

3) FASHION MNIST

For all the experiments, we adopted the training settings described in [31]. We trained all the models using vanilla SGD for 60 epochs and a batch size of 128. The learning rate starts with a value of 0.001. It decreases to 0.0002 after 42 epochs and to 0.00005 after 52 epochs. When the BT-Strategy was active, the weights for the first, second, and third hierarchy levels were 0.98, 0.01, and 0.01, respectively. These weights were modified to 0.1, 0.8, and 0.1 after 15 epochs; to 0.1, 0.2, and 0.7 after 25 epochs; and 0, 0, and 1 after 35 epochs.

C. RESULTS

Table 3, 4, and 5 present the hierarchical and level-wise performance metrics for CIFAR-10, CIFAR-100, and Fashion MNIST, respectively. The best results within the level of statistical significance are in bold. The parsimonious version of our proposal is presented as BA-CNN and the tuned version as BA-CNN*. For details about the p -values of the statistical tests performed, please refer to Tables 8, 9, 10, 11, 12, 13, 14, and 15 from the Appendix.

In all datasets, the proposed model outperformed all the baselines in terms of hierarchical accuracy and hierarchical consistency with statistically significant differences. This result confirms that our attention mechanism can equalize the information flow among the branches, enabling the model to make more accurate predictions respecting the labels' hierarchy. In particular, for CIFAR-100, hierarchical accuracy was at least ten points higher than baseline values, and hierarchical consistency was higher with margins of at least 25 points. The above demonstrates more clearly the principal advantages of the proposed model in hierarchical terms. In this dataset, the difference between the tuned and

parsimonious models is not statistically significant, so we presented the latter as a convenient choice that reduces the number of parameters. Furthermore, we obtained this result despite having only 40% of the number of parameters used by B-CNN and 18% of the number of parameters required by H-CNN to implement its highly dense connectivity pattern. Regarding hierarchical F1, the proposed model outperformed the baselines on CIFAR-10 and CIFAR-100. For Fashion MNIST, the differences were not statistically significant.

Regarding level-wise performance, on CIFAR-10 and CIFAR-100, the proposed method outperformed the baselines with statistically significant differences at the first two levels of the class hierarchy and competitive results at the last level. Furthermore, in some cases, this advantage was noticeably greater than the slight decrease observed at the third level. For instance, on CIFAR-100 using BT-Strategy, we improve the baselines' accuracy at the coarser level by more than 10 points, whereas the decrease at the third level is 4 points in the worst case. For Fashion MNIST, all the methods achieve similar results with remarkably narrow numerical differences that the statistical tests do not find significant in most cases. We must note that in this dataset, the B-CNN model already obtains an excellent performance at the coarsest level of the hierarchy (about 0.998 of accuracy), significantly reducing the margin of improvement that its variants can achieve, including the ours, which usually improves results at that level.

The above results suggest that the proposed method is a reliable extension of the referenced model: it can improve B-CNN in complex tasks, keep its performance in more straightforward tasks, and especially improve hierarchical accuracy and consistency.

D. ABLATIONS

In this section, we evaluate the robustness of the proposed model to the choice of the central backbone architecture. We also present ablation studies on the effect of two hyperparameters: (i) the dimensionality of the feature vectors propagated from the branches and (ii) the number of neurons of the attention mechanism. These hyperparameters were analyzed as a combination since they are strictly linked.

To investigate the impact of the central backbone architecture on the model's performance, we adopt an alternative model proposed by [15] and referred to as Base-B. Base-B simplifies Base-C by reducing the number of convolutional blocks to four and the number of convolutional layers in the last two blocks to two. Following [15], all models based on this architecture also simplify the branches by reducing the number of neurons in each layer to 256, 512, and 1024 for the first, second, and third branch, respectively. Following the author's recommendations, skip-connections included in H-CNN, Add-net, and Concat-net still propagates each branch's input vector to the following branches.

TABLE 3. Results on CIFAR-10.

| Model w/BT-Strategy | Accuracy by Level | | | Hierarchical Metrics | | | Param. (MM) |
|-----------------------|--------------------------|--------------------------|--------------------------|--------------------------|--------------------------|--------------------------|-------------|
| | Level 1 | Level 2 | Level 3 | Accuracy | Consistency | F1 | |
| Flat CNN [46] | - | - | 0.8823 (8.53e-04) | - | - | - | 39.98 |
| B-CNN [15] | 0.9592 (1.34e-03) | 0.9049 (1.51e-03) | 0.8862 (2.06e-03) | 0.8255 (1.52e-03) | 0.8997 (1.48e-03) | 0.9168 (8.70e-04) | 49.67 |
| BA-CNN | 0.9823 (6.70e-04) | 0.9191 (2.88e-03) | 0.8837 (2.48e-03) | 0.8751 (3.33e-03) | 0.9803 (2.31e-03) | 0.9284 (1.77e-03) | 18.76 |
| BA-CNN* | 0.9836 (6.93e-04) | 0.9230 (1.16e-03) | 0.8880 (1.65e-03) | 0.8829 (1.52e-03) | 0.9877 (1.59e-03) | 0.9315 (9.43e-04) | 15.81 |
| H-CNN [17] | 0.9604 (1.46e-03) | 0.9019 (1.47e-03) | 0.8809 (1.17e-03) | 0.8249 (1.77e-03) | 0.9070 (2.80e-03) | 0.9144 (7.57e-04) | 108.39 |
| Add-net [16] | 0.9589 (1.08e-03) | 0.8994 (1.63e-03) | 0.8831 (1.46e-03) | 0.8332 (1.23e-03) | 0.9212 (1.53e-03) | 0.9138 (7.84e-04) | 18.60 |
| Concat-net [16] | 0.9578 (1.39e-03) | 0.8994 (1.43e-03) | 0.8810 (1.69e-03) | 0.8312 (2.26e-03) | 0.9212 (2.85e-03) | 0.9127 (8.38e-04) | 49.69 |
| Model w/o BT-Strategy | | | | | | | |
| B-CNN [15] | 0.9622 (8.99e-04) | 0.9094 (1.07e-03) | 0.8742 (1.58e-03) | 0.8217 (2.17e-03) | 0.9023 (2.23e-03) | 0.9152 (7.66e-04) | 49.67 |
| BA-CNN | 0.9824 (3.61e-04) | 0.9161 (1.72e-03) | 0.8775 (1.84e-03) | 0.8724 (1.56e-03) | 0.9859 (1.17e-03) | 0.9253 (1.11e-03) | 18.75 |
| H-CNN [17] | 0.9617 (1.19e-03) | 0.9035 (1.12e-03) | 0.8743 (1.67e-03) | 0.8215 (1.85e-03) | 0.9049 (1.64e-03) | 0.9131 (5.75e-04) | 108.39 |
| Add-net [16] | 0.9643 (1.07e-03) | 0.9067 (1.30e-03) | 0.8770 (8.79e-04) | 0.8305 (1.10e-03) | 0.9174 (2.10e-03) | 0.9160 (4.93e-04) | 18.60 |
| Concat-net [16] | 0.9621 (8.16e-04) | 0.9079 (1.57e-03) | 0.8748 (1.08e-03) | 0.8263 (1.79e-03) | 0.9133 (2.02e-03) | 0.9150 (7.85e-04) | 49.69 |

TABLE 4. Results on CIFAR-100.

| Model w/BT-Strategy | Accuracy by Level | | | Hierarchical Metrics | | | Param. (MM) |
|-----------------------|--------------------------|--------------------------|--------------------------|--------------------------|--------------------------|--------------------------|-------------|
| | Level 1 | Level 2 | Level 3 | Accuracy | Consistency | F1 | |
| Flat CNN [15] | - | - | 0.6256 (2.48e-03) | - | - | - | 40.35 |
| B-CNN [15] | 0.7230 (1.79e-03) | 0.7021 (2.06e-03) | 0.6380 (3.36e-03) | 0.4478 (2.22e-03) | 0.5842 (3.03e-03) | 0.6877 (1.36e-03) | 50.06 |
| BA-CNN | 0.8333 (2.06e-03) | 0.7380 (3.25e-03) | 0.6147 (3.16e-03) | 0.5721 (3.59e-03) | 0.8540 (3.47e-03) | 0.7286 (2.27e-03) | 18.81 |
| H-CNN [17] | 0.7279 (2.12e-03) | 0.7014 (2.21e-03) | 0.6504 (2.76e-03) | 0.4639 (2.81e-03) | 0.6092 (3.10e-03) | 0.6932 (1.18e-03) | 108.78 |
| Add-net [16] | 0.6961 (4.32e-03) | 0.6480 (2.88e-03) | 0.6345 (2.62e-03) | 0.4285 (3.79e-03) | 0.5722 (3.74e-03) | 0.6596 (1.95e-03) | 18.63 |
| Concat-net [16] | 0.7102 (3.10e-03) | 0.6793 (2.36e-03) | 0.6358 (2.30e-03) | 0.4510 (2.66e-03) | 0.5971 (4.70e-03) | 0.6751 (1.27e-03) | 50.22 |
| Model w/o BT-Strategy | | | | | | | |
| B-CNN [15] | 0.7413 (2.62e-03) | 0.7320 (2.28e-03) | 0.6194 (2.85e-03) | 0.4638 (2.76e-03) | 0.6050 (3.91e-03) | 0.6976 (2.00e-03) | 50.06 |
| BA-CNN | 0.8415 (4.00e-03) | 0.7427 (3.63e-03) | 0.6080 (3.68e-03) | 0.5854 (3.94e-03) | 0.9028 (2.44e-03) | 0.7307 (3.43e-03) | 18.81 |
| BA-CNN* | 0.8409 (1.49e-03) | 0.7447 (2.66e-03) | 0.6124 (3.63e-03) | 0.5887 (3.26e-03) | 0.9025 (4.10e-03) | 0.7327 (2.21e-03) | 23.42 |
| H-CNN [17] | 0.7409 (1.99e-03) | 0.7175 (2.69e-03) | 0.6432 (2.67e-03) | 0.4773 (2.18e-03) | 0.6245 (3.31e-03) | 0.7005 (6.20e-04) | 108.78 |
| Add-net [16] | 0.7430 (2.93e-03) | 0.7096 (3.96e-03) | 0.6208 (2.64e-03) | 0.4771 (3.09e-03) | 0.6510 (3.45e-03) | 0.6911 (1.12e-03) | 18.63 |
| Concat-net [16] | 0.7391 (2.69e-03) | 0.7157 (2.07e-03) | 0.6213 (1.77e-03) | 0.4723 (3.11e-03) | 0.6272 (4.92e-03) | 0.6920 (1.21e-03) | 50.22 |

TABLE 5. Results on fashion MNIST.

| Model w/BT-Strategy | Accuracy by Level | | | Hierarchical Metrics | | | Param. (MM) |
|-----------------------|-------------------|--------------------------|--------------------------|--------------------------|--------------------------|-------------------|-------------|
| | Level 1 | Level 2 | Level 3 | Accuracy | Consistency | F1 | |
| Flat CNN [15] | - | - | 0.9299 (1.50e-03) | - | - | - | 33.69 |
| B-CNN [15] | 0.9978 (3.25e-04) | 0.9645 (1.25e-03) | 0.9309 (1.17e-03) | 0.9206 (1.36e-03) | 0.9807 (1.03e-03) | 0.9644 (6.39e-04) | 40.56 |
| BA-CNN | 0.9977 (2.97e-04) | 0.9662 (1.30e-03) | 0.9312 (2.25e-03) | 0.9262 (2.22e-03) | 0.9903 (1.00e-03) | 0.9650 (1.05e-03) | 17.41 |
| BA-CNN* | 0.9978 (3.63e-04) | 0.9663 (1.72e-03) | 0.9313 (2.16e-03) | 0.9280 (2.21e-03) | 0.9932 (7.91e-04) | 0.9652 (1.25e-03) | 15.20 |
| H-CNN [17] | 0.9981 (2.79e-04) | 0.9657 (1.06e-03) | 0.9339 (1.52e-03) | 0.9232 (1.49e-03) | 0.9794 (1.25e-03) | 0.9659 (5.85e-04) | 82.11 |
| Add-net [16] | 0.9980 (2.61e-04) | 0.9645 (1.05e-03) | 0.9318 (1.42e-03) | 0.9218 (1.78e-03) | 0.9814 (1.26e-03) | 0.9648 (7.14e-04) | 17.25 |
| Concat-net [16] | 0.9979 (2.05e-04) | 0.9638 (1.14e-03) | 0.9314 (1.43e-03) | 0.9209 (1.68e-03) | 0.9809 (1.34e-03) | 0.9644 (6.91e-04) | 40.58 |
| Model w/o BT-Strategy | | | | | | | |
| B-CNN [15] | 0.9979 (1.78e-04) | 0.9648 (1.00e-03) | 0.9260 (2.07e-03) | 0.9177 (2.30e-03) | 0.9823 (1.42e-03) | 0.9629 (9.34e-04) | 40.56 |
| BA-CNN | 0.9980 (2.16e-04) | 0.9660 (1.14e-03) | 0.9286 (1.22e-03) | 0.9262 (1.33e-03) | 0.9942 (1.04e-03) | 0.9642 (6.78e-04) | 17.41 |
| H-CNN [17] | 0.9978 (2.54e-04) | 0.9666 (9.61e-04) | 0.9288 (1.86e-03) | 0.9203 (1.84e-03) | 0.9810 (1.24e-03) | 0.9644 (8.45e-04) | 82.11 |
| Add-net [16] | 0.9978 (2.47e-04) | 0.9645 (1.25e-03) | 0.9273 (1.44e-03) | 0.9189 (1.30e-03) | 0.9829 (9.16e-04) | 0.9632 (6.20e-04) | 17.25 |
| Concat-net [16] | 0.9978 (3.67e-04) | 0.9647 (7.66e-04) | 0.9271 (1.36e-03) | 0.9192 (1.14e-03) | 0.9832 (6.37e-04) | 0.9632 (5.52e-04) | 40.58 |

Table 6 presents the hierarchical and level-wise performance metrics achieved by the proposed and baseline models when we change the central backbone architecture from Base-C to Base-B. We can observe that all models exhibit a decrease in performance. These results confirm previous studies in [15]. However, despite these differences, the relative advantage of our model is preserved. Indeed, in most cases, the proposed model outperforms the baselines in accuracy at the coarser levels of the hierarchy, and its performance is similar to the baselines at the finer level. Moreover, our model obtains a considerable and statistically significant advantage in hierarchical accuracy and consistency,

showing that the proposed approach is robust to changes in the central backbone. Details about the p -values of the statistical tests can be checked in Tables 16, 17, and 18 of the Appendix.

To evaluate the effect of varying the number of neurons in the attention modules and the dimensionality of the feature vectors propagated from the branches, we measured hierarchical accuracy for different combinations of these two hyperparameters. The dimensionality of propagated feature vectors depends on the dimensionality of the branches, so we varied that value in the range $[2^4, 2^{12}]$. For the number of neurons in our attention mechanisms, we considered values

TABLE 6. Main results on CIFAR-10 for Base-B models.

| Model w/BT-Strategy | Accuracy by Level | | | Hierarchical Metrics | | | Param. (MM) |
|-----------------------|--------------------------|--------------------------|--------------------------|--------------------------|--------------------------|--------------------------|-------------|
| | Level 1 | Level 2 | Level 3 | Accuracy | Consistency | F1 | |
| Flat CNN [15] | - | - | 0.8267 (2.64e-03) | - | - | - | 7.86 |
| B-CNN [15] | 0.9589 (8.75e-04) | 0.8687 (2.33e-03) | 0.8409 (3.02e-03) | 0.7871 (2.57e-03) | 0.9006 (1.99e-03) | 0.8895 (1.54e-03) | 12.38 |
| BA-CNN | 0.9683 (1.19e-03) | 0.8796 (2.21e-03) | 0.8395 (2.64e-03) | 0.8148 (3.43e-03) | 0.9495 (1.84e-03) | 0.8958 (1.81e-03) | 8.72 |
| H-CNN [17] | 0.9585 (1.37e-03) | 0.8702 (1.45e-03) | 0.8390 (1.62e-03) | 0.7876 (1.91e-03) | 0.9030 (3.18e-03) | 0.8892 (1.12e-03) | 29.16 |
| Add-net [16] | 0.9577 (1.87e-03) | 0.8658 (1.53e-03) | 0.8342 (2.08e-03) | 0.7946 (2.07e-03) | 0.9309 (2.62e-03) | 0.8859 (1.15e-03) | 8.57 |
| Concat-net [16] | 0.9570 (1.49e-03) | 0.8662 (1.70e-03) | 0.8360 (2.52e-03) | 0.7926 (2.49e-03) | 0.9226 (1.73e-03) | 0.8864 (1.19e-03) | 12.39 |
| Model w/o BT-Strategy | | | | | | | |
| B-CNN [15] | 0.9640 (1.34e-03) | 0.8665 (3.76e-03) | 0.8222 (3.86e-03) | 0.7774 (4.09e-03) | 0.9051 (2.36e-03) | 0.8842 (2.33e-03) | 12.38 |
| BA-CNN | 0.9724 (8.54e-04) | 0.8717 (2.45e-03) | 0.8233 (2.86e-03) | 0.8103 (2.88e-03) | 0.9677 (2.46e-03) | 0.8891 (1.79e-03) | 8.72 |
| H-CNN [17] | 0.9653 (1.45e-03) | 0.8713 (2.01e-03) | 0.8260 (2.27e-03) | 0.7867 (2.01e-03) | 0.9138 (2.08e-03) | 0.8875 (9.03e-04) | 29.16 |
| Add-net [16] | 0.9654 (9.38e-04) | 0.8675 (1.48e-03) | 0.8236 (3.18e-03) | 0.7928 (2.76e-03) | 0.9376 (1.68e-03) | 0.8855 (1.61e-03) | 8.57 |
| Concat-net [16] | 0.9645 (1.42e-03) | 0.8667 (2.30e-03) | 0.8238 (2.68e-03) | 0.7881 (2.62e-03) | 0.9280 (1.93e-03) | 0.8850 (1.81e-03) | 12.39 |

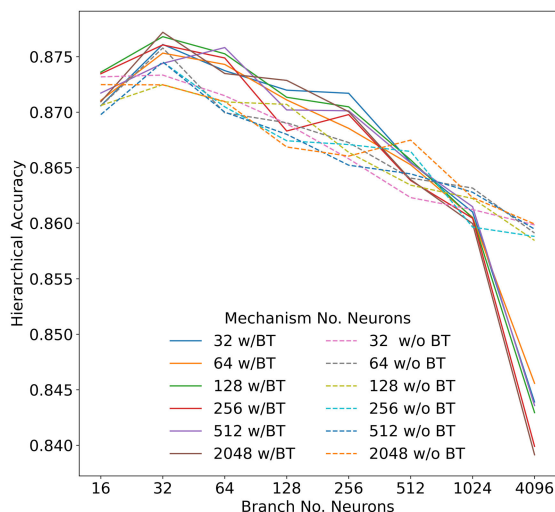


FIGURE 10. Hierarchical accuracy on CIFAR-10 as we increase the number of neurons in the branches.

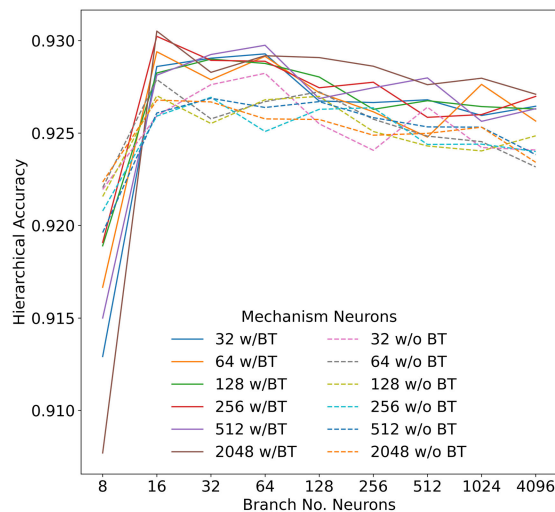


FIGURE 12. Hierarchical accuracy on Fashion MNIST as we increase the number of neurons in the branches.

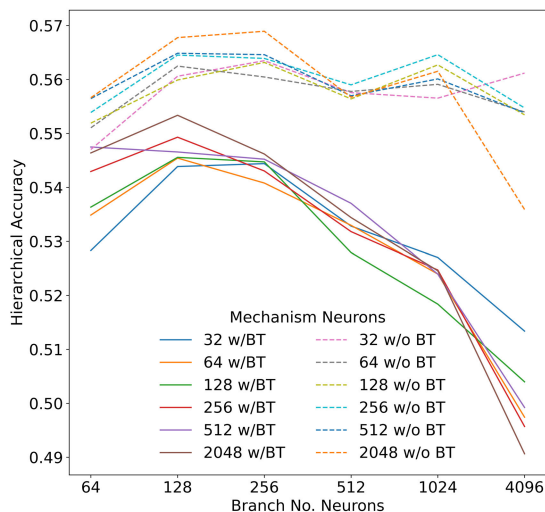


FIGURE 11. Hierarchical accuracy on CIFAR-100 as we increase the number of neurons in the branches.

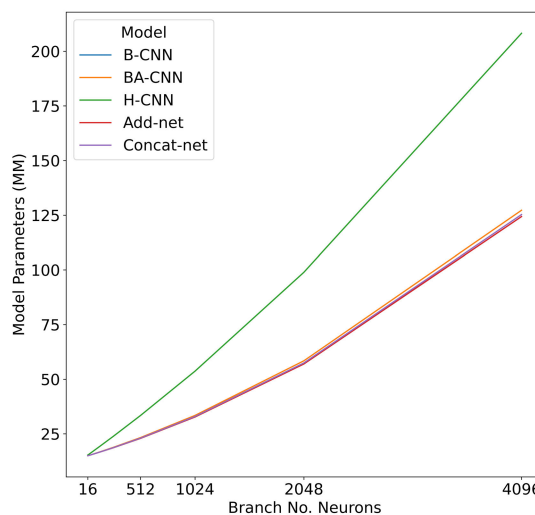


FIGURE 13. Number of model parameters on CIFAR-100 as we increase the number of neurons in the branches.

in the range $[2^5, 2^{11}]$. We adopted a logarithmic grid in both cases. Figs. 10, 11, and 12 summarize the results in the CIFAR-10, CIFAR-100, and Fashion MNIST datasets respectively. In each figure, we represent the branch number

of neurons on the horizontal axis and hierarchical accuracy on the vertical axis. Then we plot a different curve for each

TABLE 7. Time and Space Complexity of the models. As a Flat CNN representative we chose VGG-16 which is equivalent to B-CNN without the first two branches.

| Model | Number of FLOPs (MFLOPs) | GPU Memory Requirement (MB) | Model Parameters (MM) | Memory Required by Model Weights (MB) | Training Time (min) |
|----------------------|--------------------------|-----------------------------|-----------------------|---------------------------------------|---------------------|
| CIFAR-10 | | | | | |
| Flat CNN | 6.26 | 301.63 | 14.76 | 56.32 | 14.77 (0.01) |
| B-CNN | 6.27 | 308.28 | 15.19 | 57.96 | 15.98 (0.09) |
| BA-CNN | 6.29 | 318.19 | 15.82 | 60.36 | 16.79 (0.01) |
| H-CNN | 6.29 | 321.91 | 15.85 | 60.46 | 15.80 (0.30) |
| Add-net | 6.27 | 308.31 | 15.20 | 57.96 | 15.38 (0.03) |
| Concat-net | 6.27 | 308.45 | 15.20 | 57.97 | 15.37 (0.02) |
| CIFAR-100 | | | | | |
| Flat CNN | 6.28 | 302.94 | 15.35 | 58.55 | 19.93 (0.18) |
| B-CNN | 6.34 | 313.59 | 18.64 | 71.10 | 21.43 (0.03) |
| BA-CNN | 6.44 | 329.44 | 23.43 | 89.37 | 22.67 (0.03) |
| H-CNN | 6.45 | 331.59 | 23.88 | 91.10 | 22.09 (0.03) |
| Add-net | 6.34 | 313.84 | 18.64 | 71.10 | 21.52 (0.04) |
| Concat-net | 6.34 | 315.02 | 18.70 | 71.33 | 21.41 (0.25) |
| Fashion MNIST | | | | | |
| Flat CNN | 4.10 | 224.05 | 14.74 | 56.23 | 14.99 (0.02) |
| B-CNN | 4.11 | 228.47 | 14.88 | 56.75 | 15.73 (0.01) |
| BA-CNN | 4.11 | 237.96 | 15.21 | 58.02 | 16.37 (0.26) |
| H-CNN | 4.11 | 237.32 | 15.12 | 57.66 | 15.22 (0.01) |
| Add-net | 4.11 | 228.49 | 14.87 | 56.75 | 15.09 (0.02) |
| Concat-net | 4.11 | 228.55 | 14.88 | 56.76 | 15.13 (0.01) |

value of the attention mechanism's number of neurons. For completeness, we report results with and without using the BT-Strategy for training. We can observe that the model is very robust to the choice of the attention mechanism's number of neurons with a general trend to (slightly) increase as we increase the dimensionality. The branches' dimensionality is the main factor to determine the model's proneness to overfitting. Generally, small values of this parameter obtains the best results. Furthermore, the breaking point at the left (underfitting) is significantly more difficult to reach than the breaking point at the right (overfitting). The above suggests that the proposed attention module does not need a large number of trainable parameters to work effectively. We attribute this result to the module's flexibility. The model can dynamically select the pieces of data that are useful to make a prediction, and thus it does not need to explicitly keep track of many individual patterns.

As regards the use of the BT-Strategy, we observe mixed results. In two datasets (CIFAR-10 and Fashion-MNIST), hierarchical accuracy tends to be better with this training methodology. However, in CIFAR-100, the BT-Strategy causes a significant decrease in performance. Therefore, we can conclude that a model selection phase for tuning the number of neurons in the attention mechanism and predicting the BT-Strategy's efficacy on validation data is worth the effort.

E. TIME AND SPACE COMPLEXITY

In this section, we discuss the proposed model's time and space complexity and the baselines. First, the scalability of the models in terms of the number of parameters depending on the branch's dimension. Then different measures were

performed on the models to evaluate the impact of our proposal in comparison with the state-of-the-art models presented in this work.

Regarding scalability, Fig. 13 shows the effect of increasing the number of neurons in all branches. We observe that all the models increase their space complexity similarly, except H-CNN, which quickly surpasses the rest of the models. It is worth mentioning that this behavior is for the proposed model using 64 neurons in the attention mechanism. If that value increases, the space complexity of the model increases, and, at some point (2048 neurons in this setting), it can reach H-CNN's complexity.

The following measurements were computed to models and presented in Table 7:

- Number of floating point operations per second (FLOPs).
- GPU memory requirement.
- Number of model parameters.
- Memory required by model weights.
- Training time.

The above measurements were computed using the best hyperparameters. We also adjusted the baselines to that hyperparameters to make a fair comparison and evaluate the complexity of the models in the same scenario.

In terms of training time, considering the fastest model (Flat CNN) as the baseline, the proposed model presented, for all datasets, a less than 1.13x slowdown. Between models, the differences are less than 5%. It is worth noticing that the differences in training time between datasets are primarily due to the number of epochs: 80 for CIFAR-100 and 60 for CIFAR-10 and Fashion MNIST.

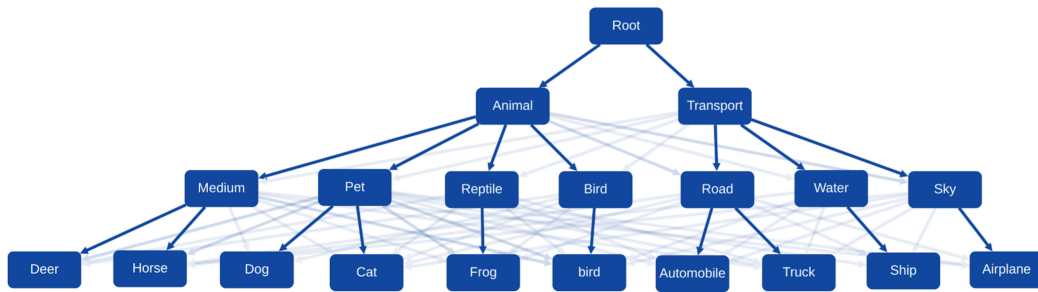


FIGURE 14. Hierarchy tree reconstructed by BA-CNN with BT-Strategy (score = 0.9803).

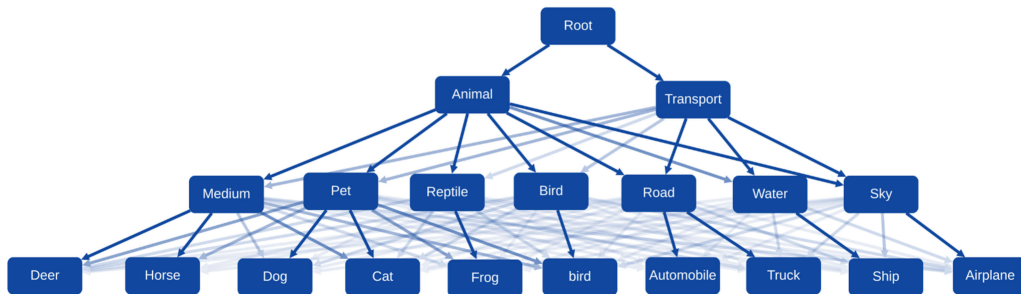


FIGURE 15. Hierarchy tree reconstructed by B-CNN with BT-Strategy (score = 0.8997).

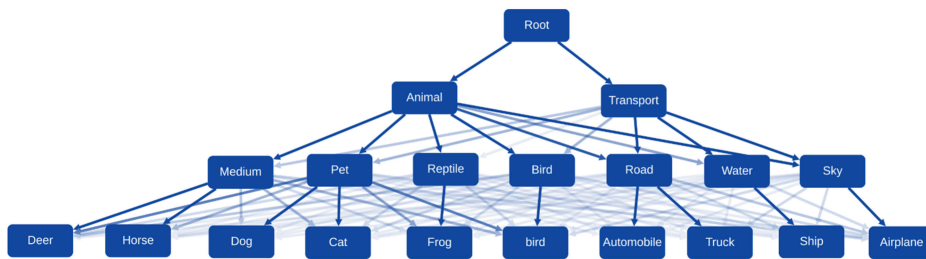


FIGURE 16. Hierarchy tree reconstructed by H-CNN with BT-Strategy (score = 0.9070).

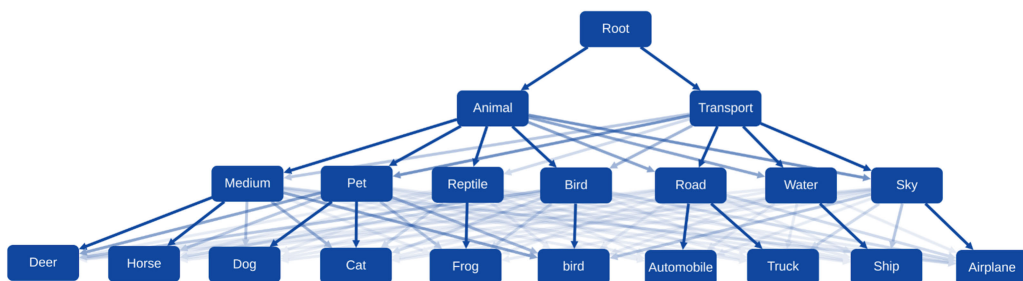


FIGURE 17. Hierarchy tree reconstructed by Add-net with BT-Strategy (score = 0.9212).

The FLOPs measurements confirm that by adopting the same branches' hyper-parameters, the time complexities of the different models are comparable. The main reason for the differences in FLOPs between the datasets is the input size of the network. CIFAR-100 and CIFAR-10 use 32×32 images, and Fashion MNIST 28×28 .

In terms of model parameters and the corresponding memory required, due to the concatenation of branches from the first dense layer, the H-CNN model presents more parameters

than the proposed model, excluding Fashion MNIST, always with a tight margin.

In summary, adding the attention mechanisms to a branched architecture does not negatively impact time and space complexity.

F. HIERARCHICAL CONSISTENCY

To visualize the relationships learned by the model, we reconstructed a hierarchy tree from the model predictions.

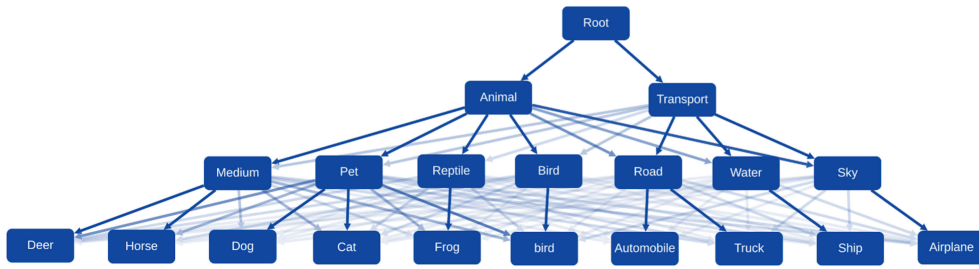


FIGURE 18. Hierarchy tree reconstructed by Concat-net with BT-Strategy (score = 0.9212).

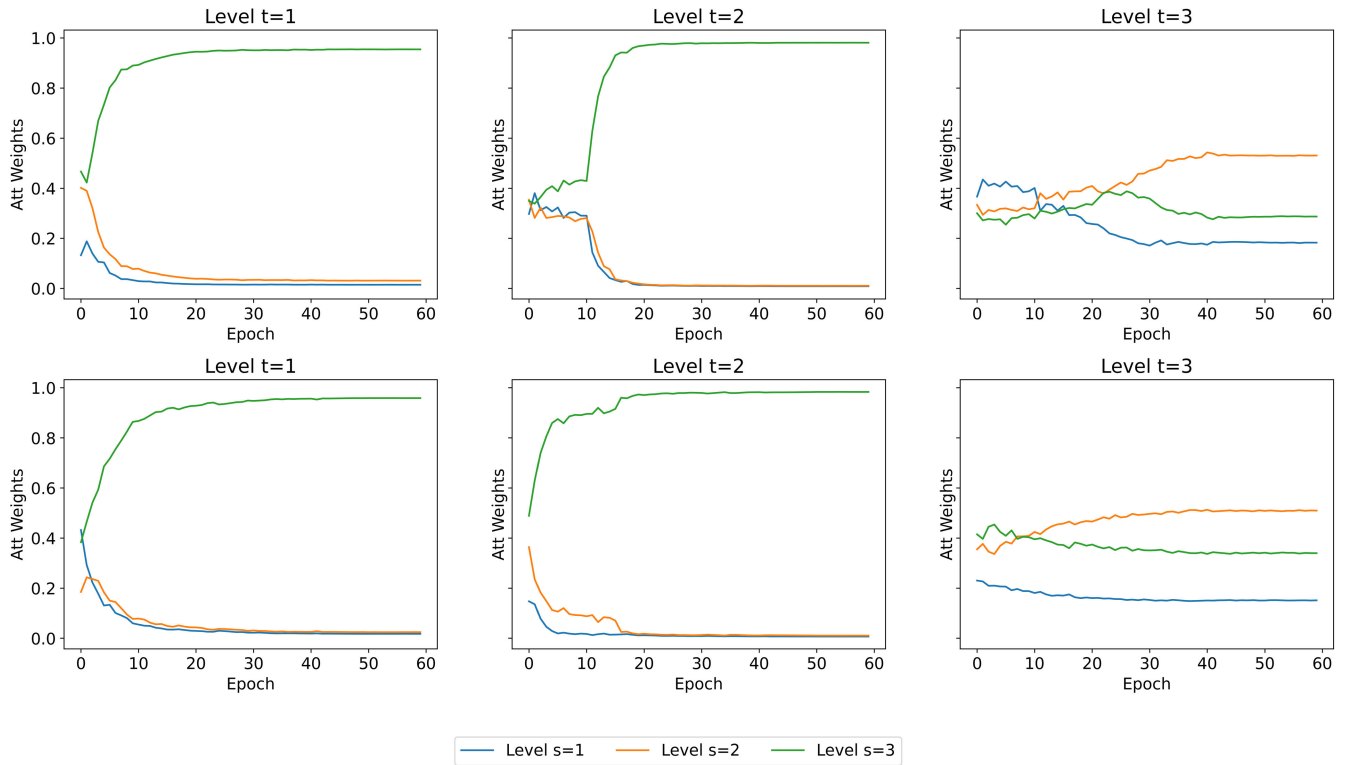


FIGURE 19. Attention weights per level on CIFAR-10. On the first row, using BT-Strategy, and in the second row without it.

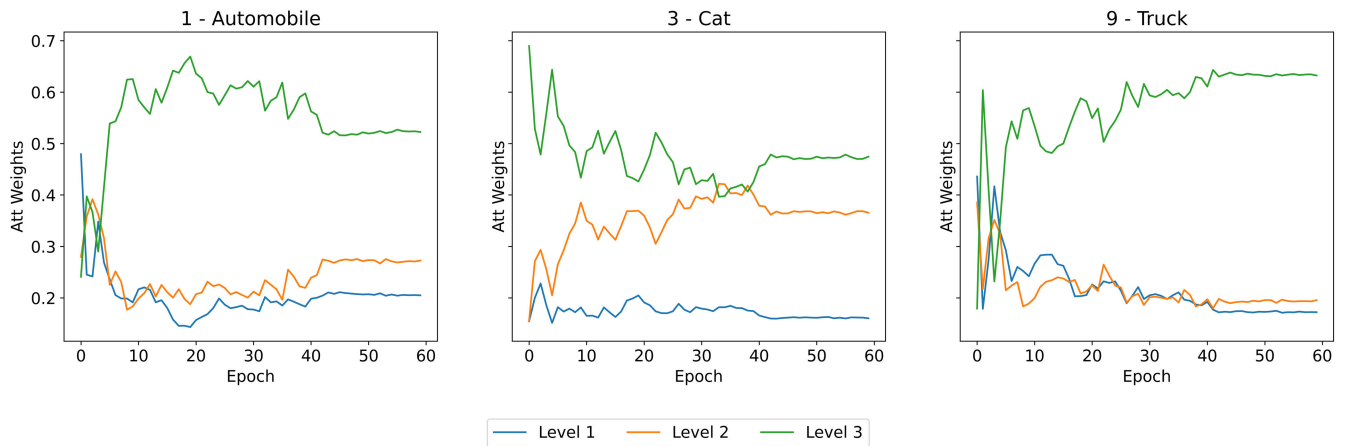


FIGURE 20. Attention weights from three classes on CIFAR-10 for t = 3.

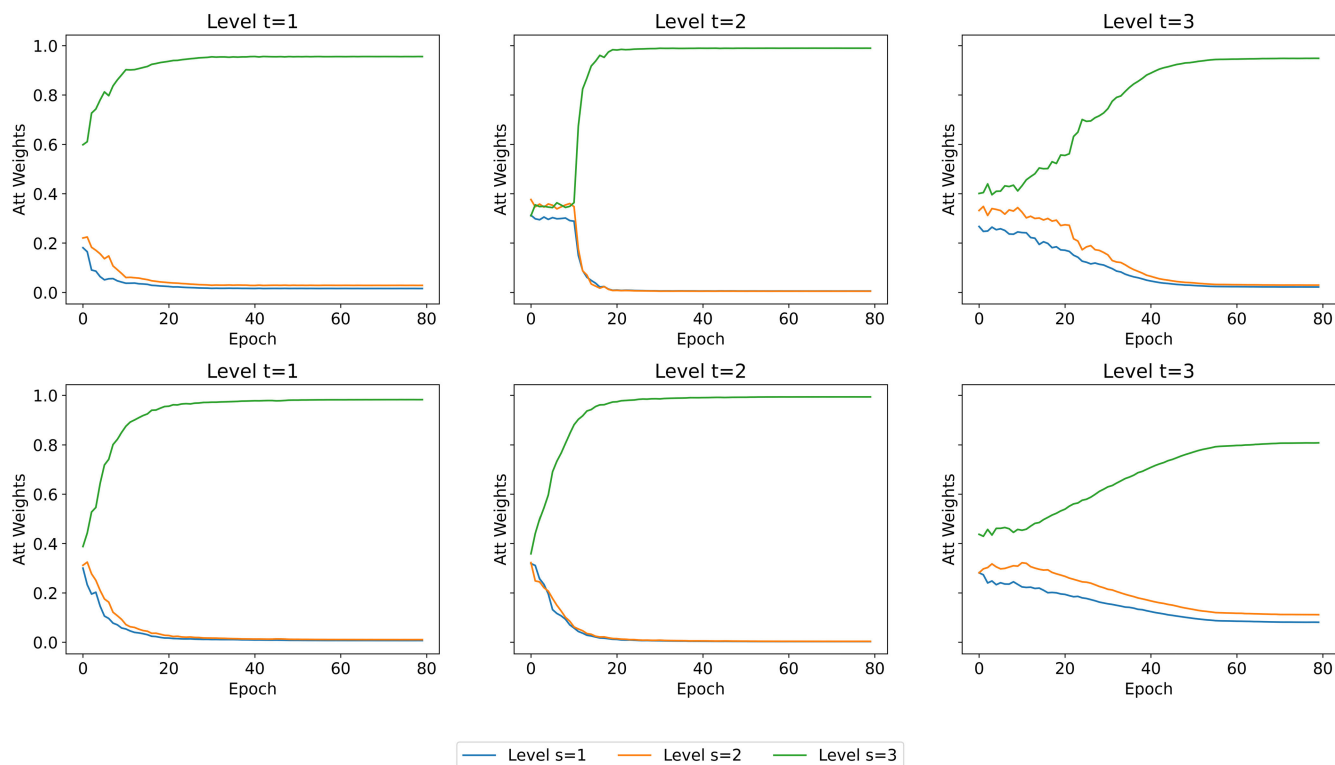


FIGURE 21. Attention weights per level on CIFAR-100. On the first row, using BT-Strategy, and in the second row without it.

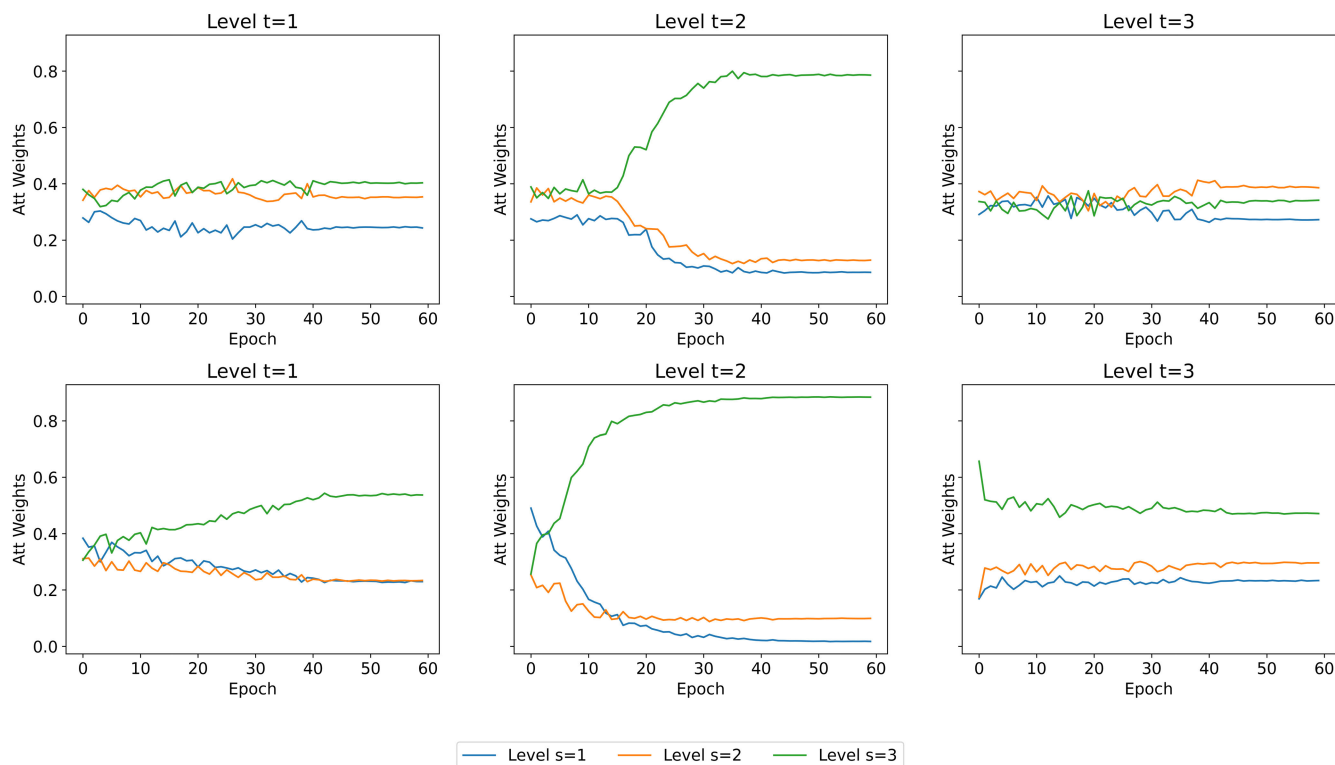


FIGURE 22. Attention weights per level on fashion MNIST. On the first row, using BT-Strategy, and in the second row without it.

Fig. 14, 15, 16, 17 and 18 show the relations for CIFAR-10 obtained using the proposed model (BA-CNN), B-CNN,

H-CNN, Add-net and Concat-net respectively. We can see that the proposed model more precisely reconstructs the

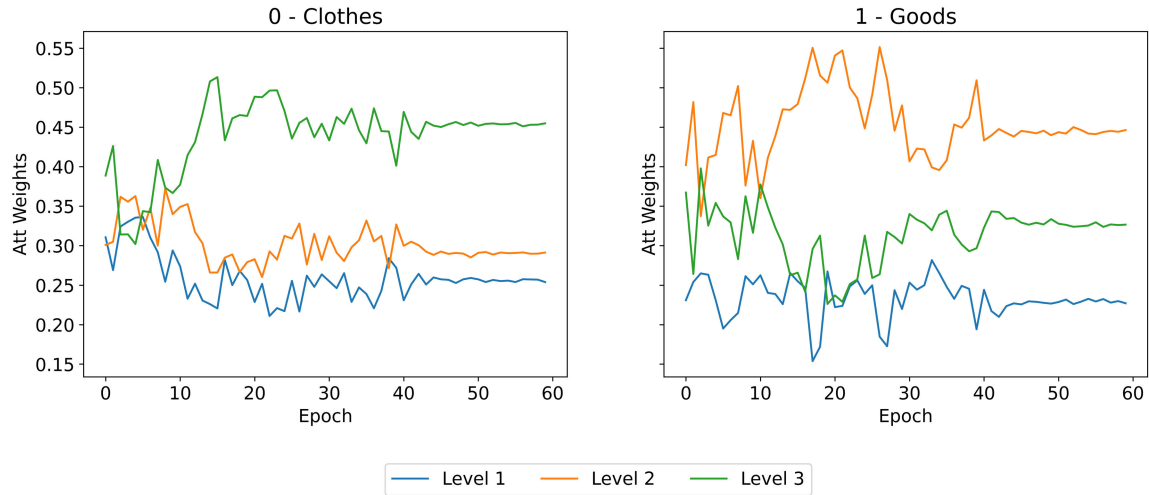


FIGURE 23. Attention weights for level 1 on Fashion MNIST.

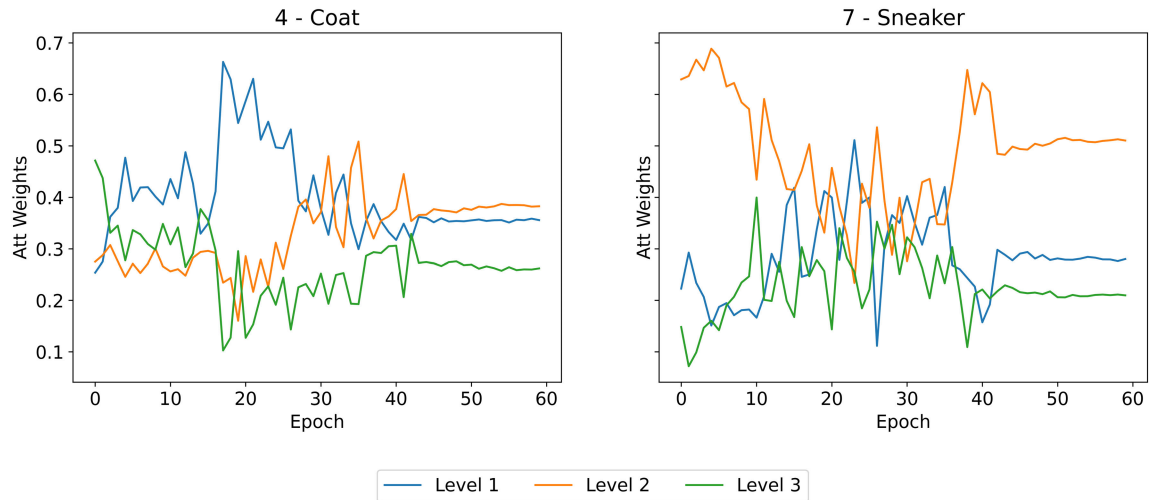


FIGURE 24. Attention weights for class “Coat” and “Sneaker” from Fashion MNIST dataset.

correct relations in the hierarchy. Each edge of the graph was drawn using an alpha value proportional to the number of predictions that included that edge.

G. ATTENTION WEIGHTS

In addition to allowing performance gains in many tasks, attention has become a popular method for gaining insights into the internals of deep learning models and their predictions [47]. Indeed, as attention allows the model to focus on a subset of features when predicting an output, attention weights have been used as proxies of feature importance since their early introduction [48]. Bahdanau et al. [40] used attention weights to analyze the alignment between words learned by the model. Xu et al. [49] showed that the attended regions in an image-captioning task aligned well with human intuition. More recently, Voita et al. [50] demonstrated that

the maximum attention weight of a head in Transformer models agrees reasonably well with its relevance to machine-translation tasks. Although some recent studies have suggested that attention scores can be inconsistent with other feature importance measures [51], [52], subsequent experiments have shown that attention weights can still help analyze the decision and learning mechanisms of otherwise black-box models [47], [53], [54].

Here, we briefly analyze how the proposed model attends to different levels of the class hierarchy during training. For a given level s , we compute the attention weight α_{st} given to branch s when classifying data at level t and report the average among the training instances. Fig. 19 depicts the resulting scores for CIFAR-10. We can see that the coarser levels ($t = 1$ and $t = 2$) give much more importance to features extracted by branch $s = 3$. Attention weights

TABLE 8. P-values of the Friedman test, by groups on CIFAR-10.

| Model | Accuracy by Level | | | Hierarchical Metrics | | |
|---------------------|-------------------|---------|---------|----------------------|-------------|--------|
| | Level 1 | Level 2 | Level 3 | Accuracy | Consistency | F1 |
| With BT-Strategy | 0.0000 | 0.0000 | 0.0000 | 0.0000 | 0.0000 | 0.0000 |
| Without BT-Strategy | 0.0000 | 0.0000 | 0.0000 | 0.0000 | 0.0000 | 0.0000 |

TABLE 9. P-values of the Wilcoxon Test on CIFAR-10.

| Model w/BT-Strategy | Accuracy by Level | | | Hierarchical Metrics | | |
|-----------------------|-------------------|---------|---------------|----------------------|-------------|--------|
| | Level 1 | Level 2 | Level 3 | Accuracy | Consistency | F1 |
| B-CNN | 0.0293 | 0.0293 | 1.0000 | 0.0293 | 0.0293 | 0.0293 |
| BA-CNN | 0.0293 | 0.0293 | 0.1230 | 0.0293 | 0.0293 | 0.0293 |
| H-CNN | 0.0293 | 0.0293 | 0.0410 | 0.0293 | 0.0293 | 0.0293 |
| Add-net | 0.0293 | 0.0293 | 0.0410 | 0.0293 | 0.0293 | 0.0293 |
| Concat-net | 0.0293 | 0.0293 | 0.0410 | 0.0293 | 0.0293 | 0.0293 |
| Model w/o BT-Strategy | | | | | | |
| B-CNN | 0.0195 | 0.0195 | 0.0586 | 0.0195 | 0.0195 | 0.0195 |
| H-CNN C | 0.0195 | 0.0195 | 0.0879 | 0.0195 | 0.0195 | 0.0195 |
| Add-net | 0.0195 | 0.0195 | 1.0000 | 0.0195 | 0.0195 | 0.0195 |
| Concat-net | 0.0195 | 0.0195 | 0.4102 | 0.0195 | 0.0195 | 0.0195 |

TABLE 10. P-values of the Friedman test, by groups on CIFAR-100.

| Model | Accuracy by Level | | | Hierarchical Metrics | | |
|---------------------|-------------------|---------|---------|----------------------|-------------|--------|
| | Level 1 | Level 2 | Level 3 | Accuracy | Consistency | F1 |
| With BT-Strategy | 0.0000 | 0.0000 | 0.0000 | 0.0000 | 0.0000 | 0.0000 |
| Without BT-Strategy | 0.0000 | 0.0000 | 0.0000 | 0.0000 | 0.0000 | 0.0000 |

TABLE 11. P-values of the Wilcoxon test against BA-CNN on CIFAR-100.

| Model w/BT-Strategy | Accuracy by Level | | | Hierarchical Metrics | | |
|-----------------------|-------------------|---------------|---------|----------------------|---------------|---------------|
| | Level 1 | Level 2 | Level 3 | Accuracy | Consistency | F1 |
| B-CNN | 0.0195 | 0.0195 | - | 0.0195 | 0.0195 | 0.0195 |
| H-CNN | 0.0195 | 0.0195 | - | 0.0195 | 0.0195 | 0.0195 |
| Add-net | 0.0195 | 0.0195 | - | 0.0195 | 0.0195 | 0.0195 |
| Concat-net | 0.0195 | 0.0195 | - | 0.0195 | 0.0195 | 0.0195 |
| Model w/o BT-Strategy | | | | | | |
| B-CNN | 0.0293 | 0.0293 | - | 0.0293 | 0.0293 | 0.0293 |
| BA-CNN | 1.0000 | 0.5566 | - | 0.1624 | 1.0000 | 0.1465 |
| H-CNN | 0.0293 | 0.0293 | - | 0.0293 | 0.0293 | 0.0293 |
| Add-net | 0.0293 | 0.0293 | - | 0.0293 | 0.0293 | 0.0293 |
| Concat-net | 0.0293 | 0.0293 | - | 0.0293 | 0.0293 | 0.0293 |

are not constant during training but quickly converge to a configuration in which features extracted from coarser branches have negligible importance. This result suggests that *bottom-up feedback* between branches is useful for learning a class taxonomy, a type of internal feedback that solutions have systematically ignored, thus favoring top-down communication patterns. Fig. 21 confirms this result for CIFAR-100. Although the convergence is slightly slower, Fig. 22 shows that the same pattern is also observed for level $t = 2$ on Fashion MNIST. Interestingly, however, level $t = 1$ spreads attention more evenly among the branches, suggesting that the attention distribution is case-dependent: fixed attention weights can work in some datasets but produce unsatisfactory results in others. Fig. 19, 21, and 22 show that attention at level $t = 3$ often behaves differently from the other levels and varies more among the datasets. In addition, these figures show that the effect of the BT-Strategy is to systematically increase the attention the branches give to the finer level of the class hierarchy, which is

TABLE 12. P-values of the Wilcoxon test against H-CNN on CIFAR-100 (Complementary table of Table 11).

| Model w/BT-Strategy | Accuracy on Level 3 |
|-----------------------|---------------------|
| B-CNN | 0.0293 |
| BA-CNN | 0.0293 |
| Add-net | 0.0293 |
| Concat-net | 0.0293 |
| Model w/o BT-Strategy | |
| B-CNN | 0.0410 |
| BA-CNN | 0.0410 |
| BA-CNN* | 0.0410 |
| Add-net | 0.0410 |
| Concat-net | 0.0410 |

expected because the BT-Strategy gradually focuses training on that branch.

In Fig. 20, 23, and 24, we show how the attention distribution varies among different classes. For this experiment, we averaged the attention weights by considering only (training) instances belonging to the same class at a given

TABLE 13. P-values of the Friedman test, by groups on Fashion MNIST.

| Model | Accuracy by Level | | | Hierarchical Metrics | | |
|---------------------|-------------------|---------|---------|----------------------|-------------|--------|
| | Level 1 | Level 2 | Level 3 | Accuracy | Consistency | F1 |
| With BT-Strategy | 0.1860 | 0.0020 | 0.0033 | 0.0000 | 0.0000 | 0.0110 |
| Without BT-Strategy | 0.6225 | 0.0020 | 0.0014 | 0.0000 | 0.0000 | 0.0027 |

TABLE 14. P-values of the Wilcoxon test against BA-CNN on Fashion MNIST.

| Model w/BT-Strategy | Accuracy by Level | | | Hierarchical Metrics | | |
|-----------------------|-------------------|---------------|---------------|----------------------|-------------|---------------|
| | Level 1 | Level 2 | Level 3 | Accuracy | Consistency | F1 |
| B-CNN | 1.0000 | 1.0000 | 1.0000 | 0.0293 | 0.0293 | 1.0000 |
| BA-CNN | 1.0000 | 1.0000 | 1.0000 | 0.2930 | 0.0293 | 1.0000 |
| H-CNN | 1.0000 | 1.0000 | 0.7793 | 0.0293 | 0.0293 | 1.0000 |
| Add-net | 1.0000 | 0.2051 | 1.0000 | 0.0293 | 0.0293 | 1.0000 |
| Concat-net | 1.0000 | 0.0879 | 1.0000 | 0.0293 | 0.0293 | 1.0000 |
| Model w/o BT-Strategy | | | | | | |
| B-CNN | 1.0000 | - | 0.2930 | 0.0195 | 0.0195 | - |
| H-CNN | 1.0000 | - | 1.0000 | 0.0195 | 0.0195 | - |
| Add-net | 1.0000 | - | 0.5566 | 0.0195 | 0.0195 | - |
| Concat-net | 1.0000 | - | 0.9668 | 0.0195 | 0.0195 | - |

level t of the hierarchy. We selected cases that depicted the variety of attention patterns we observed in our experiments. For CIFAR-10, for instance, as shown in Fig. 20, there are three classes at level $t = 3$, which contradicts the global attention distribution previously discussed. For these classes, the attention mechanism weights more the feature vector received from the same branch ($s = 3$). We can also observe in Fig. 20 that the attention given to the previous branch ($s = 2$) varies significantly between classes. In Fig. 23 and 24, we present the attention weights for two classes of the first level and two classes at the last level. We can see that the relevance of branch $s = 2$ for the predictions at level $t = 1$ can respect the global pattern (class “Clothes”) or be significantly higher than the average curves suggest (class “Goods”). Similarly, at level $t = 3$, we can find classes for which the attention given to the feature vector from the same branch ($s = 3$) is the lowest, contradicting the general behavior. We offer these results as evidence of the fact that the context information that the model needs to disambiguate a class varies from case to case. This supports the hypothesis that an attention mechanism among branches can be more effective for hierarchical classification than a fixed connectivity pattern or a handcrafted rule to feed the branches.

VII. CONCLUSION

This paper presented a novel technique for deep hierarchical classification that recombines feature maps extracted at different depths of a convolutional neural network to classify data into the different levels of the class hierarchy. The novelty of our approach lies in introducing an attention mechanism devised to learn how the predictions at different hierarchy levels must influence each other during training and inference. The attention module allows feature maps to be exchanged and fused in a dynamic data-driven way that supports, in particular, top-down and bottom-up cross-level interactions.

TABLE 15. P-values of the Wilcoxon test against H-CNN on Fashion MNIST without BT-Strategy (Complementary table of Table 14).

| Model | Accuracy on Level 2 | Hierarchical F1 |
|------------|---------------------|-----------------|
| B-CNN | 0.0977 | 0.1367 |
| BA-CNN | 1.0000 | 1.0000 |
| Add-net | 0.0195 | 0.0977 |
| Concat-net | 0.0391 | 0.1953 |

We assessed the proposed methodology using well-known hierarchical benchmarks based on the CIFAR-10, CIFAR-100, and Fashion MNIST datasets. Experiments demonstrated that, in all cases, our algorithm improves the hierarchical accuracy and hierarchical consistency of four state-of-the-art models in a statistically significant way. Regarding per-level performance, we found that the coarser taxonomic levels are typically those that benefit most from the attention mechanism. At the deeper level of the hierarchy, we observed more mixed results, but statistical tests reveal that competitive predictions for that specific level also come out very often. For instance, in the CIFAR-100 dataset, the proposed method consistently outperformed existing methods by a margin of around 10% accuracy at the first level of the hierarchy, at the cost of only about 4% accuracy at the deeper hierarchy level. This result is remarkable because CIFAR-100 has a large number of classes.

To conclude, we designed experiments to assess the hypothesis that our method is better than baselines for learning hierarchical dependencies between concepts. Experiments showed that our technique outperforms state-of-the-art models in hierarchical consistency: it typically traverses the hierarchy by choosing consistent paths from the root to the leaves. In particular, we observed the most considerable improvement in the CIFAR-100 dataset, where the proposed method achieved a 25% improvement on current methods.

In view of the experimental results, we thus conclude that our method advances current research on deep hierarchical

TABLE 16. P-values of the Friedman test, by groups on CIFAR-10 for Base-B models.

| Model | Accuracy by Level | | | Hierarchical Metrics | | |
|---------------------|-------------------|---------|---------|----------------------|-------------|--------|
| | Level 1 | Level 2 | Level 3 | Accuracy | Consistency | F1 |
| With BT-Strategy | 0.0000 | 0.0000 | 0.0000 | 0.0000 | 0.0000 | 0.0000 |
| Without BT-Strategy | 0.0000 | 0.0006 | 0.0465 | 0.0000 | 0.0000 | 0.0001 |

TABLE 17. P-values of the Wilcoxon Test, against BA-CNN on CIFAR-10 for Base-B models.

| Model w/BT-Strategy | Accuracy by Level | | | Hierarchical Metrics | | |
|-----------------------|-------------------|---------------|---------|----------------------|-------------|---------------|
| | Level 1 | Level 2 | Level 3 | Accuracy | Consistency | F1 |
| B-CNN | 0.0195 | 0.0195 | - | 0.0195 | 0.0195 | 0.0195 |
| H-CNN | 0.0195 | 0.0195 | - | 0.0195 | 0.0195 | 0.0195 |
| Add-net | 0.0195 | 0.0195 | - | 0.0195 | 0.0195 | 0.0195 |
| Concat-net | 0.0195 | 0.0195 | - | 0.0195 | 0.0195 | 0.0195 |
| Model w/o BT-Strategy | | | | | | |
| B-CNN | 0.0195 | 0.0195 | - | 0.0195 | 0.0195 | 0.0195 |
| H-CNN | 0.0195 | 1.0000 | - | 0.0195 | 0.0195 | 0.1367 |
| Add-net | 0.0195 | 0.0586 | - | 0.0195 | 0.0195 | 0.0391 |
| Concat-net | 0.0195 | 0.0977 | - | 0.0195 | 0.0195 | 0.0586 |

classification providing predictions more consistent with the class hierarchy. Furthermore, the analysis of the attentional weights revealed that hierarchical classification benefits from the dynamic connection between branches provided by the attention mechanism. The neural net updates the importance of the input feature maps as training progresses and uses different weights to classify different input patterns. This flexibility contrasts with current solutions' rigid and sometimes cumbersome connectivity patterns. In addition, the ablation study revealed that this flexibility allows the model to use a low number of parameters to obtain the current advantages, which translates only in a slight impact in terms of time and space complexity.

In current methods (including ours), the depth at which a branch emerges is handcrafted by the modeler. Future work will explore the idea that the network can automatically and dynamically make that choice. Specifically, we plan to devise a secondary attention mechanism that can select and aggregate visual features from different depths of the central block to feed the branches. We plan to investigate also different implementations of the primary attention mechanism. In particular, incorporating recent advances in deep multi-head attention is an exciting direction to explore.

APPENDIX

This Appendix includes all the p -values from the Friedman and Wilcoxon tests performed on the related experiments. It is important to note that the Wilcoxon test was performed against the proposed model. When our model did not obtain the best numerical result, we performed that test against the winner and presented the results in complementary tables. The above occurs in Tables 11, 14, and 17. P-values are highlighted when the statistical test shows significant differences at 5% level (p -value > 0.05).

TABLE 18. P-values of the Wilcoxon Test, against BA-CNN on CIFAR-10 for Base-B models (Complementary table of Table 17).

| Model w/BT-Strategy against B-CNN | Accuracy on Level 3 |
|-------------------------------------|---------------------|
| BA-CNN | 1.0000 |
| H-CNN | 1.0000 |
| Add-net | 0.0293 |
| Concat-net | 0.0293 |
| Model w/o BT-Strategy against H-CNN | |
| B-CNN | 0.5566 |
| BA-CNN | 0.1465 |
| Add-net | 1.0000 |
| Concat-net | 1.0000 |

REFERENCES

- [1] C. N. Silla and A. A. Freitas, "A survey of hierarchical classification across different application domains," *Data Mining Knowl. Discovery*, vol. 22, nos. 1–2, pp. 31–72, 2011.
- [2] J. Hernández, L. E. Sucar, and E. F. Morales, "Multidimensional hierarchical classification," *Expert Syst. Appl.*, vol. 41, no. 17, pp. 7671–7677, Dec. 2014.
- [3] C.-C. Lee, E. Mower, C. Busso, S. Lee, and S. Narayanan, "Emotion recognition using a hierarchical binary decision tree approach," *Speech Commun.*, vol. 53, nos. 9–10, pp. 1162–1171, Nov. 2011.
- [4] S. Dumais and H. Chen, "Hierarchical classification of web content," in *Proc. 23rd Annu. Int. ACM SIGIR Conf. Res. Develop. Inf. Retr.*, Jul. 2000, pp. 256–263.
- [5] G. An, M. Akiba, K. Omodaka, T. Nakazawa, and H. Yokota, "Hierarchical deep learning models using transfer learning for disease detection and classification based on small number of medical images," *Sci. Rep.*, vol. 11, no. 1, pp. 1–9, Mar. 2021.
- [6] Z. Barutcuoglu, R. E. Schapire, and O. G. Troyanskaya, "Hierarchical multi-label prediction of gene function," *Bioinformatics*, vol. 22, no. 7, pp. 830–836, Jan. 2006.
- [7] R. Cerri, R. C. Barros, and A. C. P. L. F. de Carvalho, "Hierarchical classification of gene ontology-based protein functions with neural networks," in *Proc. Int. Joint Conf. Neural Netw. (IJCNN)*, Jul. 2015, pp. 1–8.
- [8] X.-L. Wang, H. Zhao, and B.-L. Lu, "A meta-top-down method for large-scale hierarchical classification," *IEEE Trans. Knowl. Data Eng.*, vol. 26, no. 3, pp. 500–513, Mar. 2014.
- [9] S. Kiritchenko, S. Matwin, and A. F. Famili, "Functional annotation of genes using hierarchical text categorization," *Proc. ACL Workshop Linking Biol. Literature, Ontologies Databases, Mining Biolog. Semantics*, 2005, pp. 1–6.

- [10] S. Kiritchenko, S. Matwin, R. Nock, and A. F. Famili, "Learning and evaluation in the presence of class hierarchies: Application to text categorization," in *Advances in Artificial Intelligence*, L. Lamontagne and M. Marchand, Eds. Berlin, Germany: Springer, 2006, pp. 395–406.
- [11] A. Krizhevsky, I. Sutskever, and G. E. Hinton, "ImageNet classification with deep convolutional neural networks," in *Proc. Adv. Neural Inf. Process. Syst.*, vol. 25, 2012, pp. 1097–1105.
- [12] A. Bilal, A. Jourabloo, M. Ye, X. Liu, and L. Ren, "Do convolutional neural networks learn class hierarchy?" *IEEE Trans. Vis. Comput. Graphics*, vol. 24, no. 1, pp. 152–162, Jan. 2018.
- [13] R. Waseem and Z. Wang, "Deep convolutional neural networks for image classification: A comprehensive review," *Neural Comput.*, vol. 29, no. 9, pp. 2352–2449, Aug. 2017.
- [14] Z. Yan, H. Zhang, R. Piramuthu, V. Jagadeesh, D. DeCoste, W. Di, and Y. Yu, "HD-CNN: Hierarchical deep convolutional neural networks for large scale visual recognition," in *Proc. IEEE Int. Conf. Comput. Vis. (ICCV)*, Dec. 2015, pp. 2740–2748.
- [15] X. Zhu and M. Bain, "B-CNN: Branch convolutional neural network for hierarchical classification," 2017, *arXiv:1709.09890*.
- [16] M. Inoue, C. H. Forster, and A. C. dos Santos, "Semantic hierarchy-based convolutional neural networks for image classification," in *Proc. Int. Joint Conf. Neural Netw. (IJCNN)*, Jul. 2020, pp. 1–8.
- [17] X. Zhang, L. Tang, H. Luo, S. Zhong, Z. Guan, L. Chen, C. Zhao, J. Peng, and J. Fan, "Hierarchical bilinear convolutional neural network for image classification," *IET Comput. Vis.*, vol. 15, no. 3, pp. 197–207, Apr. 2021.
- [18] B. Kolisnik, I. Hogan, and F. Zulkernine, "Condition-CNN: A hierarchical multi-label fashion image classification model," *Expert Syst. Appl.*, vol. 182, Nov. 2021, Art. no. 115195.
- [19] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, and I. Polosukhin, "Attention is all you need," in *Proc. Adv. Neural Inf. Process. Syst.*, vol. 30. Red Hook, NY, USA: Curran Associates, 2017, pp. 5998–6008.
- [20] F. Wu, J. Zhang, and V. Honavar, "Learning classifiers using hierarchically structured class taxonomies," in *Abstraction, Reformulation Approximation*, J.-D. Zucker and L. Saitta, Eds. Berlin, Germany: Springer, 2005, pp. 313–320.
- [21] A. Sun and E.-P. Lim, "Hierarchical text classification and evaluation," in *Proc. IEEE Int. Conf. Data Mining*, Nov. 2001, pp. 521–528.
- [22] A. A. Freitas and A. C. de Carvalho, *Research and Trends in Data Mining Technologies and Applications, Chapter A Tutorial on Hierarchical Classification With Applications in Bioinformatics*. Seoul-t'ukpyolsi, South Korea: Idea Group Pub., 2007.
- [23] J. J. Burred and A. Lerch, "A hierarchical approach to automatic musical genre classification," in *Proc. 6th Int. Conf. Digit. Audio Effects*, 2003, pp. 8–11.
- [24] J. G. S. Barbedo and A. Lopes, "Automatic genre classification of musical signals," *EURASIP J. Adv. Signal Process.*, vol. 2007, no. 1, pp. 1–12, Dec. 2006.
- [25] Z. Kuang, J. Yu, Z. Li, B. Zhang, and J. Fan, "Integrating multi-level deep learning and concept ontology for large-scale visual recognition," *Pattern Recognit.*, vol. 78, pp. 198–214, Jun. 2018.
- [26] R. Babbar, I. Partalas, E. Gaussier, and M. R. Amini, "On flat versus hierarchical classification in large-scale taxonomies," in *Proc. Adv. Neural Inf. Process. Syst.*, vol. 26, 2013, pp. 1824–1832.
- [27] Y. Labrou and T. Finin, "Yahoo! As an ontology: Using Yahoo! Categories to describe documents," in *Proc. 8th Int. Conf. Inf. Knowl. Manage.*, Nov. 1999, pp. 180–187.
- [28] J. J. Rocchio, "Relevance Feedback in Information Retrieval. Englewood Cliffs, NJ, USA: Prentice-Hall, 1971, pp. 313–323.
- [29] C. Vens, J. Struyf, L. Schietgat, S. Džeroski, and H. Blockeel, "Decision trees for hierarchical multi-label classification," *Mach. Learn.*, vol. 73, no. 2, pp. 185–214, Nov. 2008.
- [30] A. Kosmopoulos, I. Partalas, E. Gaussier, G. Paliouras, and I. Androutsopoulos, "Evaluation measures for hierarchical classification: A unified view and novel approaches," *Data Mining Knowl. Discovery*, vol. 29, no. 3, pp. 820–865, May 2015.
- [31] Y. Seo and K.-S. Shin, "Hierarchical convolutional neural networks for fashion image classification," *Expert Syst. Appl.*, vol. 116, pp. 328–339, Feb. 2019.
- [32] R. Sali, S. Adewole, L. Ehsan, L. A. Denson, P. Kelly, B. C. Amadi, L. Holtz, S. A. Ali, S. R. Moore, S. Syed, and D. E. Brown, "Hierarchical deep convolutional neural networks for multi-category diagnosis of gastrointestinal disorders on histopathological images," in *Proc. IEEE Int. Conf. Healthcare Informat. (ICHI)*, Nov. 2020, pp. 1–6.
- [33] T.-Y. Lin, A. RoyChowdhury, and S. Maji, "Bilinear CNN models for fine-grained visual recognition," in *Proc. IEEE Int. Conf. Comput. Vis. (ICCV)*, Dec. 2015, pp. 1449–1457.
- [34] P. Aggarwal. (2019). *Fashion Product Images Dataset*. Kaggle. [Online]. Available: <https://www.kaggle.com/paramaggarwal/fashion-product-images-dataset>
- [35] R. La Grassa, I. Gallo, and N. Landro, "Learn class hierarchy using convolutional neural networks," *Int. J. Speech Technol.*, vol. 51, no. 10, pp. 6622–6632, Oct. 2021.
- [36] K. Kowsari, R. Sali, M. Khan, W. Adorno, S. Ali, S. Moore, B. Amadi, P. Kelly, S. Syed, and D. Brown, "Diagnosis of celiac disease and environmental enteropathy on biopsy images using color balancing on convolutional neural networks," *Proc. Future Technol. Conf. (FTC)*, 2019, pp. 750–765.
- [37] D. Gao, "Deep hierarchical classification for category prediction in e-commerce system," in *Proc. 3rd Workshop e-Commerce (NLP)*, 2020, pp. 64–68.
- [38] J. Chorowski, D. Bahdanau, D. Serdyuk, K. Cho, and Y. Bengio, "Attention-based models for speech recognition," in *Proc. 28th Int. Conf. Neural Inf. Process. Syst.*, vol. 1. Cambridge, MA, USA: MIT Press, 2015, pp. 577–585.
- [39] M. Guo, T. Xu, J. Liu, Z. Liu, P. Jiang, T. Mu, S. Zhang, R. R. Martin, M. Cheng, and S. Hu, "Attention mechanisms in computer vision: A survey," *Comput. Vis. Media*, vol. 8, no. 3, pp. 331–368, 2022.
- [40] D. Bahdanau, K. Cho, and Y. Bengio, "Neural machine translation by jointly learning to align and translate," in *Proc. 3rd Int. Conf. Learn. Represent., (ICLR)*, 2015.
- [41] T. Luong, H. Pham, and C. D. Manning, "Effective approaches to attention-based neural machine translation," in *Proc. Conf. Empirical Methods Natural Lang. Process.*, Sep. 2015, pp. 1412–1421.
- [42] K. Atkinson, T. Bench-Capon, and D. Bollegala, "Explanation in AI and law: Past, present and future," *Artif. Intell.*, vol. 289, Dec. 2020, Art. no. 103387.
- [43] J. Lee, J.-H. Shin, and J.-S. Kim, "Interactive visualization and manipulation of attention-based neural machine translation," in *Proc. Conf. Empirical Methods Natural Lang. Process., Syst. Demonstrations*, 2017, pp. 121–126.
- [44] A. Krizhevsky and G. Hinton, "Learning multiple layers of features from tiny images," Univ. Toronto, Toronto, ON, USA, Tech. Rep., 2009.
- [45] H. Xiao, K. Rasul, and R. Vollgraf, "Fashion-MNIST: A novel image dataset for benchmarking machine learning algorithms," 2017, *arXiv:1708.07747*.
- [46] K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image recognition," in *Proc. 3rd Int. Conf. Learn. Represent., (ICLR)*, Y. Bengio and Y. LeCun, Eds. San Diego, CA, USA, May 2015.
- [47] A. Galassi, M. Lippi, and P. Torrioni, "Attention in natural language processing," *IEEE Trans. Neural Netw. Learn. Syst.*, vol. 32, no. 10, pp. 4291–4308, Oct. 2021.
- [48] Z. Niu, G. Zhong, and H. Yu, "A review on the attention mechanism of deep learning," *Neurocomputing*, vol. 452, pp. 48–62, Sep. 2021.
- [49] K. Xu, J. Ba, R. Kiros, K. Cho, A. Courville, R. Salakhudinov, R. Zemel, and Y. Bengio, "Show, attend and tell: Neural image caption generation with visual attention," in *Proc. Int. Conf. Mach. Learn.*, 2015, pp. 2048–2057.
- [50] E. Voita, D. Talbot, F. Moiseev, R. Sennrich, and I. Titov, "Analyzing multi-head self-attention: Specialized heads do the heavy lifting, the rest can be pruned," in *Proc. 57th Annu. Meeting Assoc. Comput. Linguistics*, 2019, pp. 5797–5808.
- [51] S. Jain and B. C. Wallace, "Attention is not explanation," in *Proc. Conf. North Amer. Chapter Assoc. Comput. Linguistics, Hum. Lang. Technol.*, vol. 1, 2019, pp. 3543–3556.
- [52] S. Serrano and A. N. Smith, "Is attention interpretable," in *Proc. 57th Conf. Assoc. Comput. Linguistics, (ACL)*, vol. 1, A. Korhonen, D. R. Traum, and L. Màrquez, Eds. Florence, Italy, Aug. 2019, pp. 2931–2951.

- [53] S. Wiegreffe and Y. Pinter, "Attention is not not explanation," in *Proc. Conf. Empirical Methods Natural Lang. Process. 9th Int. Joint Conf. Natural Lang. Process. (EMNLP-IJCNLP)*, K. Inui, J. Jiang, V. Ng, and X. Wan, Eds. Hong Kong: Association for Computational Linguistics, 2019, pp. 11–20.
- [54] M. Tutek and J. Snajder, "Toward practical usage of the attention mechanism as a tool for interpretability," *IEEE Access*, vol. 10, pp. 47011–47030, 2022.



RICARDO ÑANCULEF received the Ph.D. degree in computer science from Federico Santa María Technical University, Chile, in December 2011. His Ph.D. thesis was on the use of computational geometry methods for training support vector machines.

From 2012 to 2013, he held a Postdoctoral Research position with the Intelligent Systems Laboratory, University of Bristol, U.K., working on the application of online learning methods to real-time text tagging. In 2014, he was an Assistant Professor with Federico Santa María Technical University, teaching computational statistics and machine learning. His current research interests include ensemble learning, deep learning, pattern recognition, and applications to time series and biomedical data analysis. He collaborates on these topics with research groups from the University of Bologna, Italy, and the University of Barcelona, Spain.



IVÁN PIZARRO received the bachelor's degree in naval electronic engineering from Naval Polytechnical Academy, Chile, in 2017. He is currently pursuing the master's degree in computer science with Federico Santa María Technical University, Chile.

Since 2017, he has been with the Cybersecurity Department, Chilean Navy. His research interests include machine learning, deep neural networks, and image classification.



CARLOS VALLE received the Ph.D. degree in computer science from Federico Santa María Technical University, Chile, in 2014.

From 2014 to 2018, he was a Researcher with the Computer Science Department, Federico Santa María Technical University. Since 2018, he has been an Associate Professor with the Department of Computer Science and Informatics, University of Playa Ancha, Chile. His research interests include machine learning, ensemble learning, deep neural networks, and pattern recognition with applications to time series.

• • •