# Learning Decision Trees Recurrently Through Communication

Stephan Alaniz[1,2]     Diego Marcos[3]     Bernt Schiele[2]     Zeynep Akata[1,2,4]

[1]University of Tübingen          [2]MPI for Informatics          [3]Wageningen University          [4]MPI for Intelligent Systems
Tübingen AI Center     Saarland Informatics Campus     Wageningen, Netherlands          Max Planck Campus

## Abstract

*Integrated interpretability without sacrificing the prediction accuracy of decision making algorithms has the potential of greatly improving their value to the user. Instead of assigning a label to an image directly, we propose to learn iterative binary sub-decisions, inducing sparsity and transparency in the decision making process. The key aspect of our model is its ability to build a decision tree whose structure is encoded into the memory representation of a Recurrent Neural Network jointly learned by two models communicating through message passing. In addition, our model assigns a semantic meaning to each decision in the form of binary attributes, providing concise, semantic and relevant rationalizations to the user. On three benchmark image classification datasets, including the large-scale ImageNet, our model generates human interpretable binary decision sequences explaining the predictions of the network while maintaining state-of-the-art accuracy.*

## 1. Introduction

The decision mechanism of deep Convolutional Neural Networks (CNNs) is often hidden from the user, hindering their employment in critical applications such as healthcare, where a thorough understanding of this mechanism may be required. The aim for analyzing the decision mechanism, i.e. *introspection*, is to reveal the internal process of the decision maker to a machine learning practitioner or user [48]. However, models offering explanations through introspection may result in a performance loss [18, 45].

Incorporating recent advances in multi-agent communication [15], we formulate the decision process as an iterative decision tree and embed its structure into the memory representation of a Recurrent Neural Network (RNN). Our model uses message-passing [20] with discrete symbols from a vocabulary. A tunable parameter controls whether to learn this vocabulary from scratch or to map it to human-understandable attributes assigning a meaning to every decision to improve its interpretability. Further, encoding the decision tree into the memory of an RNN retains the flexibility and performance of CNNs while being scalable. In-
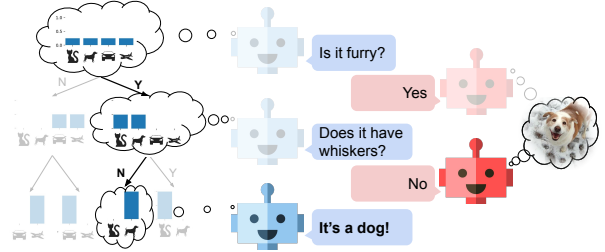


Figure 1: Our *Recurrent Decision Tree (RDT)* (left) asks questions, *Attribute-based Learner (AbL)* (right) answers with a yes/no s.t. the accuracy improves after each step.

stead of requiring an exponential number of tree nodes with increasing depth, our model learns orders of magnitude fewer nodes with a constant number of model parameters for an arbitrary tree depth. After training, our neural model can be converted exactly into a standard decision tree, being computationally efficient at test time.

Our framework (see Figure 1) exposes a decision path in the form of an explainable decision chain by breaking down the decision process into multiple binary decisions. Our recurrent decision-tree (RDT) (blue) does not see the image, and has to infer the image class, e.g. *dog*, by recurrently asking binary questions, e.g. *does it have whiskers?* Our attribute-based learner (AbL) (red) answers these questions with yes/no by looking at the image, allowing the RDT to update the class probabilities and the memory representation of the previous questions and answers. This is repeated until the RDT reaches a final decision and the decision tree becomes easily understandable as it associates the binary answers with semantic attributes, e.g. *has whiskers*.

Our contributions are: 1) We propose a recurrent decision tree model (RDTC) with hard node splits and overcome current limitations of decision trees in terms of depth scalability and flexibility; 2) We predict attributes in an end-to-end manner allowing human-interpretable explanations; 3) We showcase on three datasets that our model generates explainable decision trees more efficiently than related methods while retaining the performance of non-explainable CNNs. Our code is publicly available at: https://github.com/ExplainableML/rdtc.

## 2. Related Work

**Decision Trees with Neural Networks.** Decision trees are used across many machine learning tasks and applications, including medical diagnosis [1, 31], remote sensing [16, 19] and judicial decision making [30]. They make no assumptions on the data, and are inherently interpretable [26].

To improve their performance, combining decision trees with neural networks has been explored by building hierarchical classifiers [46, 4, 65, 64], by transferring models [25, 55, 17, 23], and through regularization [61]. Recently, [32, 56, 59] have proposed learning decision trees directly with neural networks. NBDT [59] constructs trees in the weight space of a neural network and Adaptive Neural Trees [56] directly model the neural network as a decision tree, where each node and edge correspond to one or more network modules. The prior work closest to ours is the dNDF [32], which first uses a CNN to determine the routing probabilities on each node and then combines nodes to an ensemble of decision trees that jointly make the prediction. Our method differs in that 1) we focus on explainability by explicitly only considering a hard binary decision and 2) the depth and branching structure of our decision trees is learned by an RNN instead of being fixed a priori.

**Multi-Agent Communication.** Learning to communicate in a multi-agent setting has gained interest with the emergence of deep reinforcement learning [15, 20, 37, 5, 28, 12, 9, 38, 36]. Most works focus on establishing a novel communication protocol from scratch. [15] and [5] train multiple agents to maximize a shared utility by establishing their own language. However, large scale multi-agent settings can suffer from too much communication, as valuable information comes with extensive computations [28]. Targeted communication focuses on key information and allows iterative exchange of information before performing a task that can improve both performance and interpretability [12].

Image reference games are used to study the emergence of language [37] and effectiveness in communication also when concepts are being misunderstood [9]. [20] propose an agent that composes a message of categorical symbols to another agent that uses the information in these messages to solve a referential game. Our model in contrast allows both to learn a communication protocol from scratch or use human-understandable concepts as a vocabulary.

**Attributes.** Attributes are human understandable visual properties of objects that are shared between classes. Attributes have been used for image description [14, 6], caption generation [47], face recognition [7], image retrieval [33, 54], action recognition [67, 63], novelty detection [57] and object classification [35, 53, 42, 52, 8]. In this work we propose to use attributes as explanations, i.e. they label the branches in the learned decision tree, allowing users to easily inspect the reasoning encoded by the tree.

**Explainability through sparsity.** Optimizing representations to be sparse [66] when seeking interpretability [60] draws some resemblance with the working memory of humans [40], which is limited to a handful of items at the same time. [13] hypothesizes that the nature of these items (they need to be understandable *per se*), their number and the structure in which they are presented all impact the interpretability of a representation. Furthermore, interpretability can be achieved by regularizing neural networks such that their representations, not only to become sparse [43], but also adopt the structure of a decision tree [61].

Although both sparsity and tree depth have been used as proxies for interpretability in decision trees, human studies suggest that the best proxy is problem-dependent [34]. Beyond explainable ML, a sparse representation is considered to be essential for moving towards hybrid deep learning-symbolic models [11, 44] and for obtaining representations that are closer to conscious reasoning [2]. Indeed, a recent model of how human brains work postulate a conceptualization step, linked to dimensionality reduction, followed by an attention mechanism that sparsely selects concepts [10].

## 3. **RDTC** Framework

Our Recurrent Decision Tree via Communication framework is a sequential interaction between the Recurrent Decision Tree (RDT) and Attribute-based Learner (AbL) models trained to classify images by communicating (see Figure 2). RDT learns a decision path allowing introspection and AbL provides attribute-based rationales to make the communication human-understandable.

### 3.1. Communication between RDT and AbL

For any single image $x$, our RDT model iteratively aggregates information into an explicit memory $\mathcal{M}$ that is sufficient to predict the correct class label $y \in \mathcal{Y}$. Initially, it starts with no prior information $\mathcal{M}^{(0)}$. To gather more information, the RDT agent sends a query message $c^{(t)}$ to the AbL agent. The AbL answers the query $c^{(t)}$ with a binary response $d^{(t)} \in \{0, 1\}$ that RDT uses to update its explicit memory $\mathcal{M}^{(t)} = \mathcal{M}^{(t-1)} \oplus (c^{(t)}, d^{(t)})$ to improve its class prediction. This constitutes one iteration $t$ of the agent-to-agent communication. The interaction repeats until a maximum number of steps is reached or until convergence.

**Communication Protocol.** The vocabulary size $|A|$ is set to the total number of attributes for every dataset. RDT and AbL learn to communicate with the set of tokens provided by the vocabulary in an end-to-end manner. Note that, the AbL agent attaches a human-understandable meaning to these tokens when annotated attribute data is available.

At each communication step $t$, RDT chooses one attribute $a_{c^{(t)}}$ from the vocabulary, identified by its index $c^{(t)}$, and requests its presence or absence in the image. AbL then provides its binary prediction of this attribute, i.e. $d^{(t)}$. We
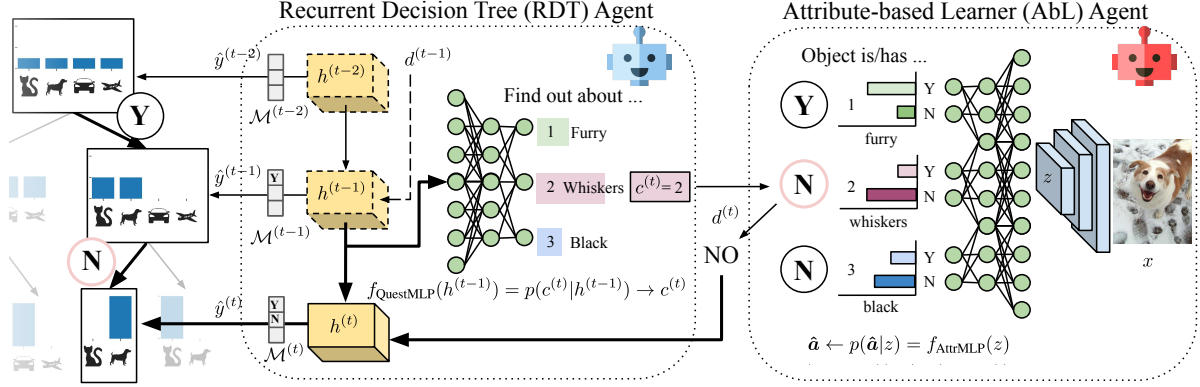
Figure 2: A single communication step between the RDT (left) and AbL (right) in our `RDTC` framework. RDT uses the hidden state $h^{(t-1)}$ of its LSTM (yellow) to requests a single attribute $a_{c^{(t)}}$ by selecting it through its $f_{\text{QuestMLP}}$. AbL uses its $f_{\text{AttrMLP}}$ to predict a binary response $d^{(t)} = \hat{a}_{c^{(t)}}$ indicating the presence/absence of the attribute. Finally, RDT updates its state $h^{(t)}$ and explicit memory $\mathcal{M}^{(t)}$ with the binary response to improve its classification prediction $\hat{y}^{(t)}$.

deliberately limit the messages of AbL to be binary as clear yes/no answers are easier to interpret.

**Discrete Messages.** RDT asks for the attribute via the index $c^{(t)}$ and the AbL responds with a binary $d^{(t)}$. The Gumbel-softmax estimator [27, 41] allows to sample from a discrete categorical distribution via the reparameterization trick [29, 50] to obtain the gradients of this sampling process. We sample $g_i$ from a Gumbel distribution and then compute a continuous relaxation of the categorical distribution:

$$\text{GumbelSoftmax}(\log \boldsymbol{\pi})_i = \frac{\exp((\log \pi_i + g_i)/\tau)}{\sum_{j=1}^{K} \exp((\log \pi_j + g_j)/\tau)} \quad (1)$$

where $\log \pi$ are the unnormalized log-probabilities of the categorical distribution, $\tau$ is the temperature that parameterizes the discrete approximation. When $\tau \approx 0$, the output is a one-hot vector and otherwise, it is a continuous signal.

Stochasticity is important for exploring all possible indices $c^{(t)}$ of vocabulary $A$ to find the most relevant attribute at each step $t$. Therefore, we use Gumbel-softmax with $K = |A|$ to sample the attribute index $c^{(t)}$ for RDT. As each $d^{(t)}$ corresponds to the presence or absence of an attribute in $x$, a deterministic prediction is beneficial. By introducing a temperature $\tau$ to a regular softmax [23] in AbL, we approximate the $\arg \max$ function deterministically as $\tau$ approaches 0:

$$\text{TempSoftmax}(\log \boldsymbol{\pi})_i = \frac{\exp(\log \pi_i/\tau)}{\sum_{j=1}^{K} \exp(\log \pi_i/\tau)} \quad (2)$$

Since we use binary attributes, in this case $K = 2$. Popular training strategies include (a) annealing $\tau$ over time and (b) augmenting the soft approximation with an $\arg \max$ that discretizes the activation in the forward pass and results in the identity function in the backward pass. Using (b) guarantees the communication to be always discrete.

## 3.2. Recurrent Decision Tree (RDT) Model

RDT consists of three parts: an explicit memory $\mathcal{M}$, an LSTM [24], and a question-decoder module, *Question MLP* (see Figure 2 (left)). $\mathcal{M}^{(t)}$ contains all the binary attributes, i.e. the responses of AbL $d_{1:t}$ up to step $t$. The LSTM keeps track of the attribute order with its hidden state $h^{(t)}$ to encode the current point in the decision tree and decide on the next question. RDT decodes its last hidden state $h^{(t-1)}$ into a categorical distribution via $f_{\text{QuestMLP}}$:

$$\log p(c^{(t)}|h^{(t-1)}) = f_{\text{QuestMLP}}(h^{(t-1)}) \quad (3)$$

where $p(c^{(t)}|h^{(t-1)})$ indicates the likelihood of requesting a particular attribute. We denote the attribute index $c^{(t)} \in \{1, \ldots, |A|\}$ sampled by:

$$c^{(t)} = \text{GumbelSoftmax}(f_{\text{QuestMLP}}(h^{(t-1)})). \quad (4)$$

After each iteration of the communication loop, RDT updates its explicit memory $\mathcal{M}^{(t)} = \mathcal{M}^{(t-1)} \oplus (c^{(t)}, d^{(t)})$. Concretely, $\mathcal{M} \in \{0,1\}^{|A| \times 2}$ is initialized with all zeros and at each time step, we set $\mathcal{M}_{c^{(t)},d^{(t)}} := 1$. Encoding the attribute in a one-hot vector helps to indicate missing information with all zeros. $\mathcal{M}^{(t)}$ keeps track of already observed attributes and their values. RDT updates $h^{(t)}$ with:

$$h^{(t)} = \text{LSTM}(h^{(t-1)}, \mathcal{M}^{(t)}, c^{(t)}, d^{(t)}). \quad (5)$$

and at each time step $\mathcal{M}$ is used to predict the class label:

$$\hat{y}^{(t)} = f_{\text{ClassMLP}}(\mathcal{M}^{(t)}). \quad (6)$$

Since the primary objective of RDT is to maximize the classification performance, we minimize the CE loss between the predicted and the true class probabilities:

$$\mathcal{L} = \frac{1}{T}\sum_{t=1}^{T} \mathcal{L}_{CE}(y, \hat{y}^{(t)}) = -\frac{1}{T}\sum_{t=1}^{T}\sum_{i} y_i \log \hat{y}_i^{(t)}. \quad (7)$$

By averaging the $\mathcal{L}_{CE}$ over all $T$ time steps, RDT predicts the correct class in a small number of communication steps which also allows it to be evaluated at any intermediate step unlike most other decision tree models that classify only at the leaf nodes (see supplementary for a comparison).

Since our decision tree can be evaluated after every communication step, the depth of the tree is not a fixed hyperparameter, but can be adaptively chosen at test time. This provides a flexible model that can be tuned for higher interpretability (shallow tree) or higher performance (deeper tree) at test time without the need for retraining.

### 3.3. Attribute-based Learner (AbL) Agent

The AbL feeds its CNN image features $z$ to $f_{\text{AttrMLP}}$ to predict a set of learned binary attributes queried by the RDT (Figure 2, right) where softmax with temperature gives us binary attributes $\hat{a} \in \{0,1\}^{|A|}$, the discretization of $p(\hat{a}|z)$:

$$\hat{a} = \text{TempSoftmax}(f_{\text{AttrMLP}}(z)). \qquad (8)$$

When the RDT requests the attribute with the index $c^{(t)}$, the AbL simply returns the binarized response about the attribute using $c^{(t)}$, i.e. $d^{(t)} = \hat{a}_{c^{(t)}}$. The attributes are either discovered in an end-to-end manner by optimizing the loss in Equation 7 (RDTC, i.e. Recurrent Decision Tree via Communication) or they are predicted as human-interpretable concepts using an attribute loss (aRDTC, i.e. attribute-based Recurrent Decision Tree via Communication).

**Attribute Loss.** Minimizing the classification loss at each time step is equivalent to finding a binary data split that reduces the class-distribution entropy the most, i.e. information gain in classical decision trees. However, a split that best separates the data is not always easy to interpret, especially when the features used for this split result from a non-linear transformation as in a CNN.

We propose to integrate further interpretability by learning $\hat{a}$ that align with class-level human-annotated attributes $\alpha$ using a second cross-entropy term weighted by $\lambda$:

$$\mathcal{L} = \frac{1}{T}\sum_{t=1}^{T}\Big[(1-\lambda)\mathcal{L}_{CE}(y,\hat{y}^{(t)}) + \lambda\mathcal{L}_{CE}(\alpha_{y,c^{(t)}},\hat{a}_{c^{(t)}})\Big]. \qquad (9)$$

Note that the attribute loss is imposed only on those attributes employed by the model. If an attribute is deemed not to be useful, e.g., if an attribute is weak or hard to predict, our RDT model learns to ignore that attribute.

When $\lambda > 0$, our model (aRDTC) learns to use ground-truth attributes and gives the binary splits a semantic meaning. For instance, the question of RDT for attribute with index $c^{(t)}$ can be interpreted as "does it have a black beak?" with $a_{c^{(t)}}$: "has black beak". When $\lambda = 0$, RDTC does not use any human-annotated attributes and automatically discovers them. Either of these settings may be desirable given the application as we show empirically.

**Algorithm 1** RDTC decision tree distillation

**Input:** Training images $X$
 Stopping *threshold*
**Output:** Decision tree DT
1: DT = empty decision tree
2: **for** $x$ in $X$ **do**
3:     DT.reset_to_root_node()
4:     $\hat{a} = \text{AbL}(x)$                 # attributes from image
5:     **for** $t = 1$ to $n$ **do**
6:         $\hat{y}^{(t)}, c^{(t)} = \text{RDT.step}(d^{(t)})$  # class/attribute of node
7:         **if** not DT.node_exists() **then**
8:             DT.add_node($\hat{y}^{(t)}, c^{(t)}$)
9:         **end if**
10:         **if** $\max_i \hat{y}_i^{(t)} > threshold$ **then**
11:             break                 # prune when confident
12:         **end if**
13:         $d^{(t)} = \hat{a}_{c^{(t)}}$                 # attribute yes/no
14:         DT.to_next_branch($d^{(t)}$)     # $1 \rightarrow$ left; $0 \rightarrow$ right
15:     **end for**
16: **end for**
17: **return** DT

### 3.4. Decision Tree Distillation

The RDT and AbL are trained end-to-end since the communication between these two models is differentiable. At test time we distill the RDT into an explicit decision tree, i.e., the global structure of nodes, including splitting feature and threshold. The distilled decision tree then models the trained neural network $f_{\text{RDT}} \equiv f_{\text{DT}}$.

Algorithm 1 describes the procedure of extracting a decision tree from RDT. The decision nodes of the decison tree make hard splits based on the presence/absence of an attribute. GumbelSoftmax adds stochasticity to RDT, which is useful for training, but at test time deterministically choosing the attribute with highest probability is essential for improving the performance and learning a static tree. Hence, it is replaced with TempSoftmax.

We start with an empty decision tree and fill it with nodes as we run the training data through the whole RDTC model (lines 1-2). Whenever a previously unseen node is discovered, we add it to the tree including information about the attribute ($c^{(t)}$), where the next node is added, i.e., left of the current node if $d^{(t-1)}$ is 1 or to the right if $d^{(t-1)}$ is 0 (line 14), and the current prediction of the class labels $\hat{y}^{(t)}$ (lines 7-9). We prune the distilled decision tree, i.e. we stop adding nodes to the tree once $\hat{y}$ is greater than a threshold (=0.95) for a class (lines 10-12). The end result is a decision tree that outputs the same class predictions as our trained neural network RDT given the attribute prediction from our AbL, while being fully explainable by matching learned attributes with human-annotated attribute data.

## 4. Experiments

**Datasets and attributes.** We validate our model on the large-scale ImageNet [51] with 1.2M images from 1K classes. In addition, we use AWA2 [35, 62] and CUB [58], i.e. two medium-scale benchmark attribute datasets. AWA2 comprises 37K images from 50 animal classes with 85 attributes, while CUB contains 11K images from 200 fine-grained bird species with 312 attributes. Since our model considers splits on hard decisions, we binarize the attributes on all datasets with a threshold at 0.5, i.e., an attribute is present if more than 50% of the annotations agree. When an official classification test set is not provided, for all experiments across the datasets, we randomly assign 20% of each class as test data and 10% of the training data as a validation set to tune hyperparameters.

**Architecture and parameters.** The MLPs consist of two layers with a ReLU non-linearity. We learn the temperature $\tau$ of the Gumbel-softmax estimator jointly with the network from an initial value for $\tau$. During training, we always roll out the decision sequence to a maximum number of steps. At test time, we apply our decision tree distillation and stop as soon as the RDT reaches a confidence level specified by a *threshold* parameter (or once the maximum number of decisions is reached). We report the mean per-class accuracy over 5 runs to avoid bias towards highly populated classes.

### 4.1. Comparing with the State of the Art

We compare our aRDTC and RDTC with classical decision trees (aDT and DT) as baselines, ResNet (ResNet [21] and aResNet) and Deep Neural Decision Forests (dNDF) [32] as the state of the art.

**ResNet and aResNet.** ResNet-152 pre-trained on ImageNet and fine-tuned on each of the datasets including its softmax classifier serves as non-explainable deep neural network (ResNet). Augmented with attribute data, we train aResNet by first predicting the attributes with the same architecture as our AbL model and then a linear layer on top.

**Our aRDTC and RDTC.** Our attribute-based recurrent decision tree (aRDTC) (Section 3.3) uses the attribute loss to associate a human-understandable meaning to the binary decisions. On the other hand, our recurrent decision tree (RDTC) does not use an attribute loss ($\lambda = 0$), and therefore purely optimizes classification performance.

**dNDF.** The dNDF explicitly models the decision tree by mapping each inner node to an output neuron with sigmoid activation. These nodes define the routing probabilities of the input to the leaves through exhaustive tree traversal where each leaf node stores a class distribution. The final prediction is the averaged class prediction weighted by the routing probabilities of every leaf. As using multiple randomized trees weakens the interpretability, for a fair comparison, we use a single tree instead of random forests.

| Model | AWA2 | CUB | ImageNet |
|---|---|---|---|
| ResNet [21] | $98.2 \pm 0.0$ | $79.0 \pm 0.2$ | $73.0 \pm 0.1$ |
| aResNet | $98.3 \pm 0.0$ | $77.3 \pm 0.5$ | N/A |
| DT | $92.3 \pm 0.4$ | $43.5 \pm 0.3$ | $55.2 \pm 1.0$ |
| dNDF [32] | $97.6 \pm 0.2$ | $73.8 \pm 0.3$ | $72.6 \pm 0.1$ |
| RDTC (Ours) | $\mathbf{98.0} \pm 0.1$ | $\mathbf{78.1} \pm 0.2$ | $\mathbf{72.8} \pm 0.1$ |
| aDT | $97.9 \pm 0.9$ | $70.6 \pm 1.3$ | N/A |
| aRDTC (Ours) | $\mathbf{98.1} \pm 0.0$ | $\mathbf{77.9} \pm 0.6$ | N/A |

Table 1: Comparing our aRDTC ($\lambda = 0.2$) and RDTC ($\lambda = 0$) to the decision tree (aDT and DT), closely related dNDF [32], and ResNet [21] (aResNet, i.e. ResNet with attribute prediction). As ImageNet do not have attributes, aResNet, aRDTC and aDT are not applicable (over 5 runs).

**aDT and DT.** The classical decision tree (DT) is learned on top of the same image features $z$ by the perceptual module. At each time step, the dataset is split using a single dimension of $z$ until a leaf node only contains samples of the same class or a regularization strategy leads to early stopping. We incorporate attributes into the DT baseline, i.e. Attribute Decision Tree (aDT). First, we train a MLP on top of the image features $z$ to predict class attributes using a binary cross-entropy loss analogously to the attribute loss of our aRDTC model. Second, we fit a decision tree on these predicted attributes for each image to determine the class. Both DT and aDT are learned using the CART algorithm [3] and the Gini impurity index as splitting criterion due to its computational advantage over entropy-based methods [49].

**Classification results.** As observed in Table 1, compared to the Decision Tree baselines of their kind, our model variants achieve significantly higher accuracy across all datasets, e.g. RDTC vs DT achieves 98.0% vs 92.3% and aRDTC vs aDT achieves 98.1% vs 97.9% on AWA2 because our model scales better and reaches consistent results through gradient-based optimization. Moreover, although RDTC and aRDTC work with constrained single-bit communications to improve explainability, they succeed in maintaining the accuracy of the non-explainable state-of-the-art across all datasets, e.g. 72.8% vs 73.0% on ImageNet.

Fine-grained decision splits are extremely challenging to explain because objects are visually similar to each other and the distinguishing factor is nuanced. Despite this challenge on CUB, the classification accuracy of RDTC is almost twice as high than classical decision trees that use the same deep features, i.e., 78.1% vs 43.5% DT. On the other hand, our RDTC not only outperforms dNDF (78.1% vs 73.8%), our model exhibits improved interpretability, because we use hard instead of soft binary splits. As it is hard for non-experts to judge the correctness of the predictions, explanations in this domain are particularly im-
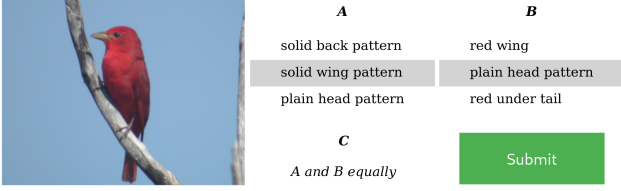
Figure 3: The user picks which set of 3 attributes best fit the image or if they match equally well (attributes come from 2 models out of aRDTC, aDT, aResNet at a time).



Figure 4: User study results. We show how often the attributes of one model were preferred over any other and when both were found equal (middle).

portant. Typically, associating a semantic meaning to the decision path improves human interpretability with a significant loss in accuracy, e.g. aResNet vs ResNet (77.3% vs 79.0%). On the other hand, on our model this trade-off is less pronounced. When trained with the attribute loss, i.e. `aRDTC` achieves a higher accuracy compared to aDT (77.9% vs 70.6% on CUB) as well as aResNet in addition giving a semantic meaning to the splits.

**User study.** The use of named attributes enables humans to understand the decision of the model. However, if all attributes are allowed to be used simultaneously, such as in aResNet, this decision becomes less comprehensible. In contrast, aRDTC provides a sparse solution that considers only a subset of attributes for each prediction. To quantify the relevance of the selected attributes, we perform a user study with aRDTC, aDT and aResNet on CUB. Since our aRDTC predicts the class label at each step, we select the attributes that change the class probability the most to determine the most critical attributes for the decision. For aDT we use the Mean Decrease Impurity (MDI) [39] to find features of maximum importance and for aResNet we select the attributes with the highest weight for the output class.

The user is prompted with an image as well as two sets of three attributes, i.e., the three most relevant attributes from two models at a time. As some attributes are difficult to recognize, e.g. cone beak, we provide attribute icons with their names and a bird anatomy sketch. The task is to select the set of three attributes that best match the image (see Figure 3). The user can also report that both sets of attributes fit the image equally well. We repeat the study on 600 randomly selected images from the CUB test such that each model is compared 200 times against every other model.
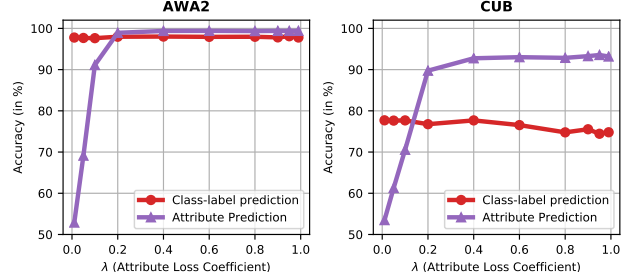


Figure 5: Explainability trade-off with `aRDTC` on AWA2 and CUB. We vary $\lambda$ of the attribute loss and report image classification accuracy of RDT (red) and attribute prediction accuracy of AbL (purple). $\lambda \in [0.01, 0.99]$.

We measure how often the attributes of each model are chosen over the other models. Since we only show attributes of two models at a time, we obtain a direct comparison for all pairs of models. Our results in Figure 4 indicate that decision tree models select more relevant attributes than aResNet. The attributes of aRDTC are preferred much more often than aResNet (48% vs. 17%). Similarly, aDT is selected more often than aResNet (46.5% vs. 21%). When comparing the two decision tree models, our aRDTC is slightly favored at 29% over aDT at 26.5% with the majority of users finding them produce equally fitting attributes (44.5%). These results suggest that the tree structure of the decision making also helps in isolating more relevant attributes by putting more weight on individual attributes selected early by the decision tree rather than spreading the contribution among all attributes.

### 4.2. Evaluating The Model Components

In this section, we evaluate several aspects of our model such as its behavior towards accuracy-explainability trade-off, ablating its memory mechanism and scalability.

**Accuracy and Explainability Trade-Off** The trade-off between the classification loss and the attribute loss in our `aRDTC` model can be measured by varying $\lambda \in [0.01, 0.99]$. Our results on AWA2 and CUB in Figure 5 show a slight decrease in the overall classification accuracy (red curve), when $\lambda$ approaches to 1.0 which gives more weight to the attribute prediction as opposed to class label prediction. Indeed, RDTC achieves a higher accuracy than `aRDTC` that is trained with the attribute loss indicating a tradeoff between explainability and accuracy. Increasing $\lambda$ leads to a slight decrease in classification accuracy, and generally similar to that of fully optimizing class prediction when $\lambda = 0$.

Furthermore, we measure the effect of $\lambda$ to the attribute prediction accuracy of the decision tree as compared to their ground-truth (purple curve). We observe a high attribute prediction accuracy even with a small $\lambda$, e.g. $\lambda = 0.2$. As we increase $\lambda$ in the range of 0.2 to 1.0, there is only a slight
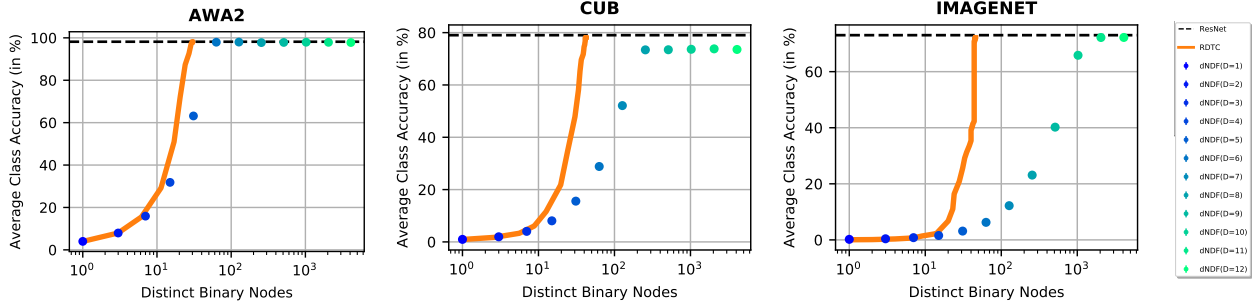
Figure 6: Accuracy with increasing number of nodes in `RDTC` and dNDF on AWA2, CUB and ImageNet. As `RDTC` can reuse learned nodes in the tree and has adaptive tree depth, we train it once and evaluate it at different depths. dNDF needs to be retrained for every depth hyperparameter (D) and the number of nodes scales exponentially with tree depth.

| Model | AWA2 (# att) | CUB (# att) | ImageNet (# att) |
|---|---|---|---|
| RDTC-L | 97.7 (19) | 73.0 (50) | 60.8 (159) |
| RDTC-M | 97.9 (57) | 77.2 (93) | 71.6 (82) |
| RDTC | **98.0** (30) | **78.1** (42) | **72.8** (46) |
| aRDTC-L | 97.9 (29) | 69.1 (32) | N/A |
| aRDTC-M | 98.0 (37) | 76.4 (52) | N/A |
| aRDTC | **98.1** (34) | **77.9** (38) | N/A |

Table 2: Ablating the memory mechanism of `aRDTC` ($\lambda = 0.2$) and `RDTC` ($\lambda = 0$). Tree state is encoded as either only the LSTM (`L`), only the explicit memory (`M`) or both (+ median number of distinct attributes the model learns).

increase in attribute prediction accuracy, indicating that our `aRDTC` is robust against the choice of $\lambda$ across datasets as long as it is chosen to be at least 0.2.

**Ablating the Memory Mechanism** The LSTM state $h$ and explicit memory $\mathcal{M}$ in RDT contains previously observed decision nodes and the current decision. While the LSTM state allows to encode the attribute order, the explicit memory serves as a more direct representation of all gathered information about the image. We ablate our RDT model with respect to its tree encoding-types.

Table 2 shows the classification accuracy of the following configurations: `aRDTC-L`, i.e. with only LSTM and attribute loss, and `aRDTC-M`, i.e. with only explicit memory, (vs `RDTC-L` and `RDTC-M` without the attribute loss). We observe that `aRDTC-M` consistently performs better than `aRDTC-L`, e.g. on CUB (76.4% vs. 69.1%). Moreover, combining the two in our full model generally improves the performance (up to 1.5% on CUB). These results indicate that the explicit memory is important for accuracy.

The median number of distinct attributes in parenthesis shows that `aRDTC-L` retains fewer decision nodes than `aRDTC-M`. For instance, `aRDTC` learns to only use 38 out of all 312 attributes of CUB ($\approx 12\%$) and on ImageNet `RDTC` uses only 46 learned binary attributes as opposed to the

1000 continuous features commonly used in ResNet. This increases sparsity of our model in the attribute space and improves interpretability by using fewer nodes when using the LSTM. We conclude that combining the two memory types in our RDTC model provides the best of both worlds, a high classification accuracy in few binary decisions such that the explanations of our model are concise and accurate. **Scalability of the Learned Decision Trees.** For our `RDTC`, increasing the tree depth simply translates to increasing the number of binary decisions, i.e., time steps of the model-to-model communication. Hence, `RDTC` scales linearly with the depth of the tree while the number of weights stays constant. On the other hand, DT and dNDF grow exponentially in their number of parameters with the depth of the tree. When the same attribute is needed at different locations in the tree, our model learns the meaning of this attribute once and reuses it, while DT and dNDF would have to relearn the split. Finally, `RDTC` does not require finetuning a depth parameter. Hence, we have the flexibility of changing the tree depth at test time without retraining.

We compare the classification accuracy of `RDTC` and dNDF with an increasing number of distinct tree nodes on three datasets. As shown in Figure 6, `RDTC` (orange line) is trained only once and evaluated at different tree depths at test time while we have to retrain dNDF for each depth parameter. While the number of nodes of dNDF scales exponentially with depth (note the log-scale on the x-axis), our model adaptively learns the number of binary attributes needed to solve these classification tasks. Hence, it stops using more attributes when no further distinct nodes are necessary. We observe that `RDTC` uses up to an order of magnitude fewer tree nodes on AWA2, CUB and ImageNet to achieve the same or better performance. At the same time, `RDTC` only needs to be trained once and can be adaptively reduced in tree depth at test time.

### 4.3. Qualitative Results

Zooming into the decision process of misclassifications on CUB, we investigate how our model treats counterfac-
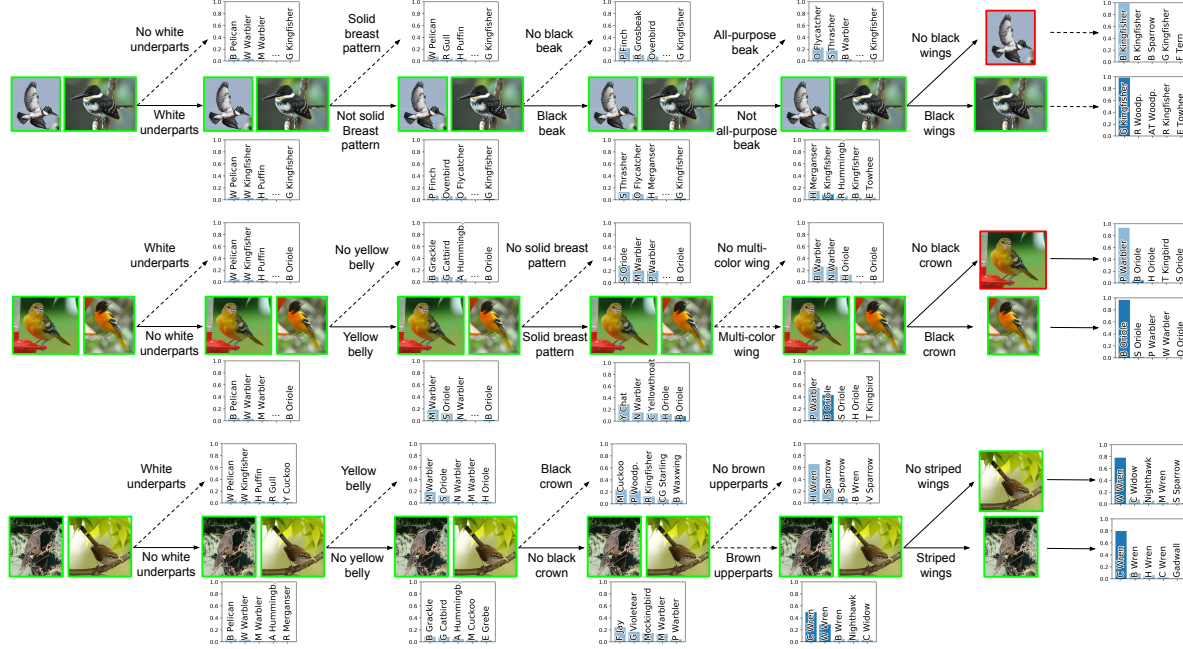
Figure 7: Top: Two "Green Kingfisher" images follow the same path except for "black wings", i.e. the flying bird gets misclassified as a "belted kingfisher" as black wings are not visible. Middle: Baltimore Oriole image (left) gets incorrectly classified as Prothonotary Warbler because of the missing "black crown" in the female bird. Such discrepancies, e.g. per-class attributes not reflecting the image content, make CUB difficult. Bottom: Cactus Wren (left) and Bewick Wren (right) share many characteristics except from "striped wings" which our model uses to split these classes.

tual classes which is useful as explanations are often contrastive [22]. We provide further qualitative examples revealing the decision tree of our aRDTC model on the fine-grained CUB and the decision tree of our RDTC model on ImageNet without the attributes in the supplementary.

In Figure 7 (Top), we inspect the point in the tree where the error occurred. The lower path corresponds to the most probable path taken for birds of class Green Kingfisher. Both images follow the same path for four decisions, the error occurs in the fifth decision. For the flying bird, our model decides that it "does not have black wings" and incorrectly classifies it as a Belted Kingfisher, a closely related class to Green Kingfisher, but without black wings. In addition, our model depicts its current belief of the correct class at any time during the process, i.e., probability plots at every branch which reveals some critical binary decisions, when the predicted class changes drastically, such as the "black wings" decision. This way, a user inspecting our explainable decision tree can make a more informed decision on the value of the prediction of the model.

In Figure 7 (Middle), the Baltimore Oriole image on the left gets incorrectly classified as Prothonotary Warbler because of the missing male-specific "black crown" attribute in the female bird. Such discrepancies, e.g. per-class attributes not reflecting the image content, make CUB an ex-

tremely challenging dataset. In Figure 7 (Bottom), the Cactus Wren image on the left and Bewick Wren image on the right share many characteristics except from "striped wings". The decision path is common until then where our model uses this attribute to split these classes.

## 5. Conclusion

In this work, we propose to learn a decision tree recurrently through communication between two-agents. Our RDTC framework adaptively changes tree depth at test time, allows to reuse of the learned decision nodes and improves scalability. It also uses human understandable attributes and hard binary splits for easier interpretation. Our experiments show that combining an explicit memory and an LSTM is important to obtain good performances with few inquiries. Our model maintains the accuracy of non-explainable deep models and outperforming the state-of-the-art deep decision tree learners. Qualitatively inspecting individual examples demonstrates the reasoning behind the failure and other challenging fine-grained cases, while a user study shows that RDTC selects more visually relevant attributes than a comparable linear semantic bottleneck model.

# References

[1] Ahmad Taher Azar and Shereen M El-Metwally. Decision tree classifiers for automated medical diagnosis. *Neural Computing and Applications*, 23(7-8), 2013. 2

[2] Yoshua Bengio. The consciousness prior. *CoRR:1709.08568*, 2017. 2

[3] Leo Breiman, Jerome Friedman, Charles J Stone, and Richard A Olshen. *Classification and regression trees*. CRC press, 1984. 5

[4] Clemens Alexander Brust and Joachim Denzler. Integrating domain knowledge: Using hierarchies to improve deep classifiers. In *Pattern Recognition*. Springer, 2020. 2

[5] Kris Cao, Angeliki Lazaridou, Marc Lanctot, Joel Z. Leibo, Karl Tuyls, and Stephen Clark. Emergent communication through negotiation. In *ICLR*, 2018. 2

[6] Huizhong Chen, Andrew Gallagher, and Bernd Girod. Describing clothing by semantic attributes. In *ECCV*, 2012. 2

[7] H. Chen, A. Gallagher, and B. Girod. What's in a name: First names as facial attributes. In *CVPR*, 2013. 2

[8] Runjin Chen, Hao Chen, Ge Huang, Jie Ren, and Quanshi Zhang. Explaining neural networks semantically and quantitatively. In *ICCV*, 2019. 2

[9] Rodolfo Corona, Stephan Alaniz, and Zeynep Akata. Modeling conceptual understanding in image reference games. In *NeurIPS*, 2019. 2

[10] Aurelio Cortese, Benedetto De Martino, and Mitsuo Kawato. The neural and cognitive architecture for learning from a small sample. *Current opinion in neurobiology*, 55, 2019. 2

[11] Miles Cranmer, Alvaro Sanchez-Gonzalez, Peter Battaglia, Rui Xu, Kyle Cranmer, David Spergel, and Shirley Ho. Discovering symbolic models from deep learning with inductive biases. *CoRR:2006.11287*, 2020. 2

[12] Abhishek Das, Théophile Gervet, Joshua Romoff, Dhruv Batra, Devi Parikh, Mike Rabbat, and Joelle Pineau. TarMAC: Targeted multi-agent communication. In *ICML*, 2019. 2

[13] Finale Doshi-Velez and Been Kim. Towards a rigorous science of interpretable machine learning. *CoRR:1702.08608*, 2017. 2

[14] Vittorio Ferrari and Andrew Zisserman. Learning visual attributes. In *NeurIPS*. 2008. 2

[15] Jakob N. Foerster, Yannis M. Assael, Nando de Freitas, and Shimon Whiteson. Learning to communicate with deep multi-agent reinforcement learning. In *NeurIPS*, 2016. 1, 2

[16] Mark A Friedl and Carla E Brodley. Decision tree classification of land cover from remotely sensed data. *Remote Sensing of Environment*, 61(3), 1997. 2

[17] Nicholas Frosst and Geoffrey Hinton. Distilling a neural network into a soft decision tree. *CoRR:1711.09784*, 2017. 2

[18] David Gunning and David W Aha. DARPA's explainable artificial intelligence program. *AI Magazine*, 40(2), 2019. 1

[19] Matthew C Hansen, Ruth S DeFries, John RG Townshend, and Rob Sohlberg. Global land cover classification at 1 km spatial resolution using a classification tree approach. *International Journal of Remote Sensing*, 21(6-7), 2000. 2

[20] Serhii Havrylov and Ivan Titov. Emergence of language with multi-agent games: Learning to communicate with sequences of symbols. In *NeurIPS*, 2017. 1, 2

[21] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *CVPR*, 2016. 5

[22] Lisa Anne Hendricks, Ronghang Hu, Trevor Darrell, and Zeynep Akata. Grounding visual explanations. In *ECCV*, 2018. 8

[23] Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. Distilling the knowledge in a neural network. *CoRR:1503.02531*, 2015. 2, 3

[24] Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural Computation*, 9(8), 1997. 3

[25] K. D. Humbird, J. L. Peterson, and R. G. Mcclarren. Deep neural network initialization with decision trees. *IEEE Trans. Neural Netw. Learn. Syst*, 30(5), 2019. 2

[26] Johan Huysmans, Karel Dejaeger, Christophe Mues, Jan Vanthienen, and Bart Baesens. An empirical evaluation of the comprehensibility of decision table, tree and rule based predictive models. *Decision Support Systems*, 51(1):141–154, 2011. 2

[27] Eric Jang, Shixiang Gu, and Ben Poole. Categorical reparameterization with gumbel-softmax. In *ICLR*, 2017. 3

[28] Jiechuan Jiang and Zongqing Lu. Learning attentional communication for multi-agent cooperation. In *NeurIPS*, 2018. 2

[29] Diederik P. Kingma and Max Welling. Auto-encoding variational bayes. In *ICLR*, 2014. 3

[30] Jon Kleinberg, Himabindu Lakkaraju, Jure Leskovec, Jens Ludwig, and Sendhil Mullainathan. Human decisions and machine predictions. *The Quarterly Journal of Economics*, 133(1), 2018. 2

[31] Igor Kononenko. Machine learning for medical diagnosis: history, state of the art and perspective. *Artificial Intelligence in Medicine*, 23(1), 2001. 2

[32] Peter Kontschieder, Madalina Fiterau, Antonio Criminisi, and Samuel Rota Bulò. Deep neural decision forests. In *IJCAI*, 2016. 2, 5

[33] Neeraj Kumar, Peter Belhumeur, and Shree Nayar. Facetracer: A search engine for large collections of images with faces. In David Forsyth, Philip Torr, and Andrew Zisserman, editors, *ECCV*, 2008. 2

[34] Isaac Lage, Andrew Ross, Samuel J Gershman, Been Kim, and Finale Doshi-Velez. Human-in-the-loop interpretability prior. In *NeurIPS*, 2018. 2

[35] C. H. Lampert, H. Nickisch, and S. Harmeling. Attribute-based classification for zero-shot visual object categorization. *IEEE TPAMI*, 36(3), 2014. 2, 5

[36] Angeliki Lazaridou and Marco Baroni. Emergent multi-agent communication in the deep learning era. *CoRR:2006.02419*, 2020. 2

[37] Angeliki Lazaridou, Karl Moritz Hermann, Karl Tuyls, and Stephen Clark. Emergence of linguistic communication from referential games with symbolic and pixel input. In *ICLR*, 2018. 2

[38] Angeliki Lazaridou, Anna Potapenko, and Olivier Tieleman. Multi-agent communication meets natural language: Synergies between functional and structural language learning. In *ACL*, 2020. 2

[39] Gilles Louppe, Louis Wehenkel, Antonio Sutera, and Pierre Geurts. Understanding variable importances in forests of randomized trees. In Christopher J. C. Burges, Léon Bottou, Zoubin Ghahramani, and Kilian Q. Weinberger, editors, *NeurIPS*, 2013. 6

[40] Wei Ji Ma, Masud Husain, and Paul M Bays. Changing concepts of working memory. *Nature Neuroscience*, 17(3), 2014. 2

[41] Chris J. Maddison, Andriy Mnih, and Yee Whye Teh. The concrete distribution: A continuous relaxation of discrete random variables. In *ICLR*, 2017. 3

[42] Dhruv Kumar Mahajan, Sundararajan Sellamanickam, and Vinod Nair. A joint learning framework for attribute models and object descriptions. In *ICCV*, 2011. 2

[43] Diego Marcos, Ruth Fong, Sylvain Lobry, Rémi Flamary, Nicolas Courty, and Devis Tuia. Contextual semantic interpretability. In *ACCV*, 2020. 2

[44] Gary Marcus. The next decade in ai: four steps towards robust artificial intelligence. *CoRR:2002.06177*, 2020. 2

[45] Raphaël Marée, Pierre Geurts, Justus Piater, Louis Wehenkel, K-S Hong, and Z Zhang. A generic approach for image classification based on decision tree ensembles and local sub-windows. In *ACCV*, 2004. 1

[46] V. N. Murthy, V. Singh, T. Chen, R. Manmatha, and D. Comaniciu. Deep decision network for multi-class image classification. In *CVPR*, 2016. 2

[47] Vicente Ordonez, Girish Kulkarni, and Tamara L. Berg. Im2text: Describing images using 1 million captioned photographs. In *NeurIPS*. 2011. 2

[48] Dong Huk Park, Lisa Anne Hendricks, Zeynep Akata, Anna Rohrbach, Bernt Schiele, Trevor Darrell, and Marcus Rohrbach. Multimodal explanations: Justifying decisions and pointing to the evidence. In *CVPR*, 2018. 1

[49] Laura Elena Raileanu and Kilian Stoffel. Theoretical comparison between the gini index and information gain criteria. *Annals of Mathematics and Artificial Intelligence*, 41(1), 2004. 5

[50] Danilo Jimenez Rezende, Shakir Mohamed, and Daan Wierstra. Stochastic backpropagation and approximate inference in deep generative models. In *ICML*, 2014. 3

[51] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, Alexander C. Berg, and Li Fei-Fei. ImageNet Large Scale Visual Recognition Challenge. *IJCV*, 115(3), 2015. 5

[52] Viktoriia Sharmanska, Novi Quadrianto, and Christoph H. Lampert. Augmented attribute representations. In *ECCV*, 2012. 2

[53] Zhiyuan Shi, Yongxin Yang, Timothy M. Hospedales, and Tao Xiang. Weakly supervised learning of objects, attributes and their associations. In *ECCV*, 2014. 2

[54] B. Siddiquie, R. Feris, and L. Davis. Image ranking and retrieval based on multi-attribute queries. In *CVPR*, 2011. 2

[55] Chapman Siu. Transferring tree ensembles to neural networks. In *ICONIP*, 2019. 2

[56] Ryutaro Tanno, Kai Arulkumaran, Daniel Alexander, Antonio Criminisi, and Aditya Nori. Adaptive neural trees. In *ICML*, 2019. 2

[57] Catherine Wah and Serge Belongie. Attribute-based detection of unfamiliar classes with humans in the loop. In *CVPR*, 2013. 2

[58] C. Wah, S. Branson, P. Welinder, P. Perona, and S. Belongie. The Caltech-UCSD Birds-200-2011 Dataset. Technical Report CNS-TR-2011-001, California Institute of Technology, 2011. 5

[59] Alvin Wan, Lisa Dunlap, Daniel Ho, Jihan Yin, Scott Lee, Henry Jin, Suzanne Petryk, Sarah Adel Bargal, and Joseph E Gonzalez. NBDT: Neural-backed decision trees. *CoRR:2004.00221*, 2020. 2

[60] John Wright, Yi Ma, Julien Mairal, Guillermo Sapiro, Thomas S Huang, and Shuicheng Yan. Sparse representation for computer vision and pattern recognition. *Proceedings of the IEEE*, 98(6), 2010. 2

[61] Mike Wu, Michael C Hughes, Sonali Parbhoo, Maurizio Zazzi, Volker Roth, and Finale Doshi-Velez. Beyond sparsity: Tree regularization of deep models for interpretability. In *AAAI*, 2018. 2

[62] Yongqin Xian, Christoph H. Lampert, Bernt Schiele, and Zeynep Akata. Zero-shot learning - A comprehensive evaluation of the good, the bad and the ugly. *IEEE Trans. Pattern Anal. Mach. Intell.*, 41(9), 2019. 5

[63] Bangpeng Yao, Xiaoye Jiang, Aditya Khosla, Andy Lai Lin, Leonidas J. Guibas, and Fei-Fei Li. Human action recognition by learning bases of action attributes and parts. In *ICCV*, 2011. 2

[64] Quanshi Zhang, Ruiming Cao, Ying Nian Wu, and Song-Chun Zhu. Growing interpretable part graphs on convnets via multi-shot learning. In Satinder P. Singh and Shaul Markovitch, editors, *AAAI*, 2017. 2

[65] Quanshi Zhang, Ruiming Cao, Ying Nian Wu, and Song-Chun Zhu. Mining object parts from cnns via active question-answering. In *CVPR*, 2017. 2

[66] Yuchen Zhang, Jason D Lee, and Michael I Jordan. L1-regularized neural networks are improperly learnable in polynomial time. In *ICML*, 2016. 2

[67] Qiang Zhou and Gang Wang. Atomic action features: A new feature for action recognition. In *ECCVW*, 2012. 2

# Learning Decision Trees Recurrently Through Communication
#
-
#
## Supplementary Material

## A. Ablating Loss and Maximum Steps T

When training our `RDTC` model with a cross-entropy loss for image classification, we found that training is more efficient and yields better results when applying the loss term at every step $t$ in the communication loop up to the maximum step $T$.

$$\mathcal{L} = \frac{1}{T} \sum_{t=1}^{T} \mathcal{L}_{CE}(y, \hat{y}^{(t)}) \tag{10}$$

One natural alternative to this approach is to apply the loss only at step $T$, the leaf node, essentially removing the sum Equation 10. In Figure 8, we show the difference of applying the loss at every time step (full loss) or only at the end (leaf loss) on CUB and AWA2. The final performance of the decision tree is the same, however, applying the loss at every time step produces a tree that has a better performance when evaluated at intermediate steps and results in a smaller tree after pruning, i.e., fewer tree nodes are used for the final tree to obtain best performance.

Moreover, we found that when hyperparameter $T$ is chosen sufficiently high, we are able to reach this maximum performance while our tree distillation process ensures that the tree size does not increase past the point where the classifier achieves the highest accuracy. Figure 9 shows classification accuracy with increasing tree depth. Accuracy does not decrease past some value for $T$ where the model performs best, and choosing any value bigger results in an equally explainable tree after pruning.

## B. Decision Trees and Explanations of CUB and AWA2

Illustrating the decision making process helps the user get an explainable overview of the internal decision process of the whole classifier. We point to the tree branch into which a certain class (indicated by an example image from this class) falls along with the attribute associated with that branch. We inspect the learned structure of the decision tree by illustrating the splits from our `aRDTC` model on CUB in Figure 11, and on AWA2 in Figure 13. Here, the left and right sub-tree indicates that the attribute is present or absent respectively. For instance, on CUB, the first decision deals with identifying bird with white underparts, sep-
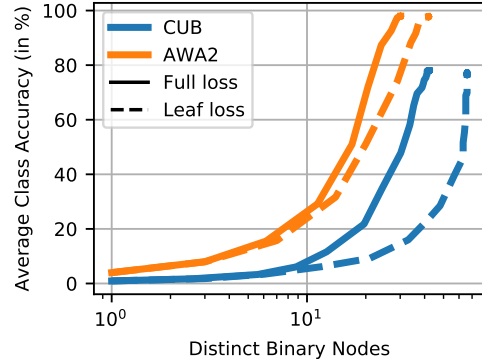


Figure 8: Accuracy and number of distinct nodes in `RDTC` on AWA2, CUB comparing our full loss at each time step (solid line) with a loss only applied at leaf nodes (dashed line). The full loss uses fewer nodes, i.e., a smaller tree, to achieve the same accuracy.
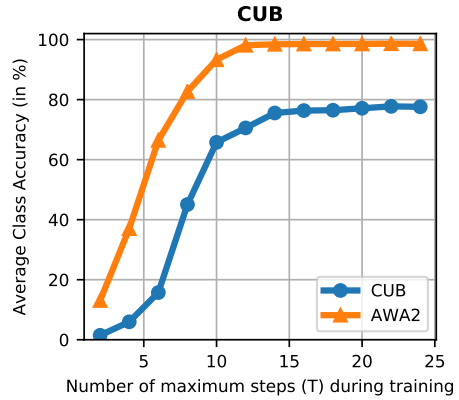


Figure 9: Training `RDTC` while varying hyperparameter $T$. As $T$ increases, the model achieve a better accuracy up to a value of $T$ where a plateau is reached. When increasing $T$ further, the final tree size of `RDTC` does not increase due to pruning during tree distillation.

arating these from birds with any other color. These categories get further refined with each binary split via a hierarchical clustering that reveals the decision tree structure of

our `aRDTC` framework. These serve as additional examples of introspection, showing that our model allows to make a more informed decision about the trustworthiness of the network's prediction.

In Figure 7, we illustrate a qualitative example of the classification of two images of Scarlet Tanagers made by our model trained on CUB. Both images follow the same path for the first decisions, before diverging when it comes to the decision whether the bird has black wings. The top bird actually does not have black wings and, thus, is classified as a Summer Tanager, a bird species with the same appearance as Scarlet Tanager except for having red instead of black wings.

Equivalently in Figure 14, we illustrate a qualitative example of the classification of two images of tigers made by our model trained on AWA2. Again, both images follow the same decision until, for the white tiger, our model wrongly predicts "no stripes" and incorrectly classifies it as a lion. Together with the full decision trees, these explanations allow for detailed introspections into the global decision process our `aRDTC` model.

## C. Explanations without Attributes

When working with datasets that do not provide annotated attributes, we can train our `RDTC` with $\lambda = 0$, which still exposes the decision tree structure. This allows introspection into the intermediate class splits of the model revealing a hierarchy that can reveal semantics. When applies on CIFAR-10, our `RDTC` model not only retains ResNet performance (93.1% vs. 93.3%), it also semantically clusters the data even though there is no attribute guidance. Figure 10 shows the resulting decision tree of `RDTC` on CIFAR-10. In the first binary split, we observe that `RDTC` separates the animal classes from the vehicles. Subsequently, vehicles are clustered into motor vehicles (car, truck) and the rest (airplane, ship). For animals, our model also finds reasonable clusters such as grouping cat and dog, as well as grouping horse and deer. ImageNet is a more challenging dataset, where we observe similar behaviour. In Figure 15, we show the decision tree of the first decisions on ImageNet with a randomly selected subset of classes, each represented by one representative image. Our model separates animals from inanimate objects in the first tree split following the data semantics. In the later decisions of the tree, there are clusters of dogs/cats, birds, monkeys on one side of the tree and clusters of furniture and electrical appliances on the other. These example show that, even when no additional attribute information is given, tree splits often follow semantics that are exposed by the decision tree learned by our `RDTC`.
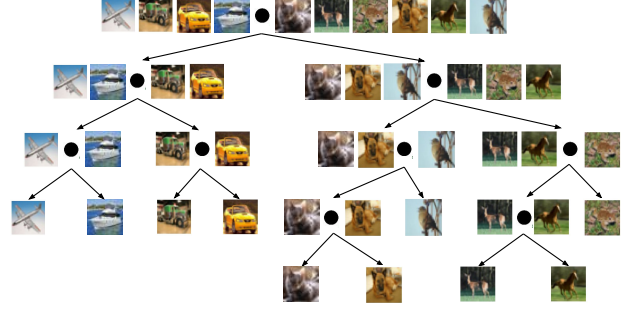


Figure 10: Our `RDTC` learns the decision tree on CIFAR10 without attribute data. While decision nodes do not have ground truth attributes, we can still interpret the decision, e.g., the first node separates animals from vehicles.

---

**Algorithm 2** RDTC training

**Input:** Image $x$, label $y$
Max # of decisions $T$, Attribute data $\alpha$
**Output:** Predicted label $\hat{y}$
Binary decision sequence $d^{(1)}, \ldots, d^{(T)}$

1: $z = \text{CNN}(x)$
2: $\hat{a} = \text{TempSoftmax}(f_{\text{AttrMLP}}(z))$
3: init $\mathcal{M}^0, h_0$
4: $\mathcal{L} = 0$
5: **for** $t = 1$ to $T$ **do**
6: $\quad c^{(t)} = \text{GumbelSoftmax}(f_{\text{QuestMLP}}(h^{(t-1)}))$
7: $\quad d^{(t)} = \hat{a}_{c^{(t)}}$
8: $\quad \mathcal{M}^{(t)} = \mathcal{M}^{(t-1)} \oplus (c^{(t)}, d^{(t)})$
9: $\quad h^{(t)} = \text{LSTM}(h^{(t-1)}, \mathcal{M}^{(t)}, c^{(t)}, d^{(t)})$
10: $\quad \hat{y}^{(t)} = f_{\text{ClassMLP}}(\mathcal{M}^{(t)})$
11: $\quad \mathcal{L}^{(t)} = \frac{1}{T}\left[(1-\lambda)\mathcal{L}_{CE}(y, \hat{y}^{(t)}) + \lambda\mathcal{L}_{CE}(\alpha_{y,c^{(t)}}, \hat{a}_{c^{(t)}})\right]$
12: $\quad \mathcal{L} = \mathcal{L} + \mathcal{L}^{(t)}$
13: **end for**
14: gradient update with $\mathcal{L}$
15: **return** $\hat{y}^{(T)}; d^{(1)}, \ldots, d^{(T)}$

---

## D. RDTC Training Algorithm

For a concise representation of the `RDTC` training algotithm, we present a summary in Algorithm 2 including both components, RDT and AbL, iterative loss calculation and gradient updates using the terminology of the main paper.
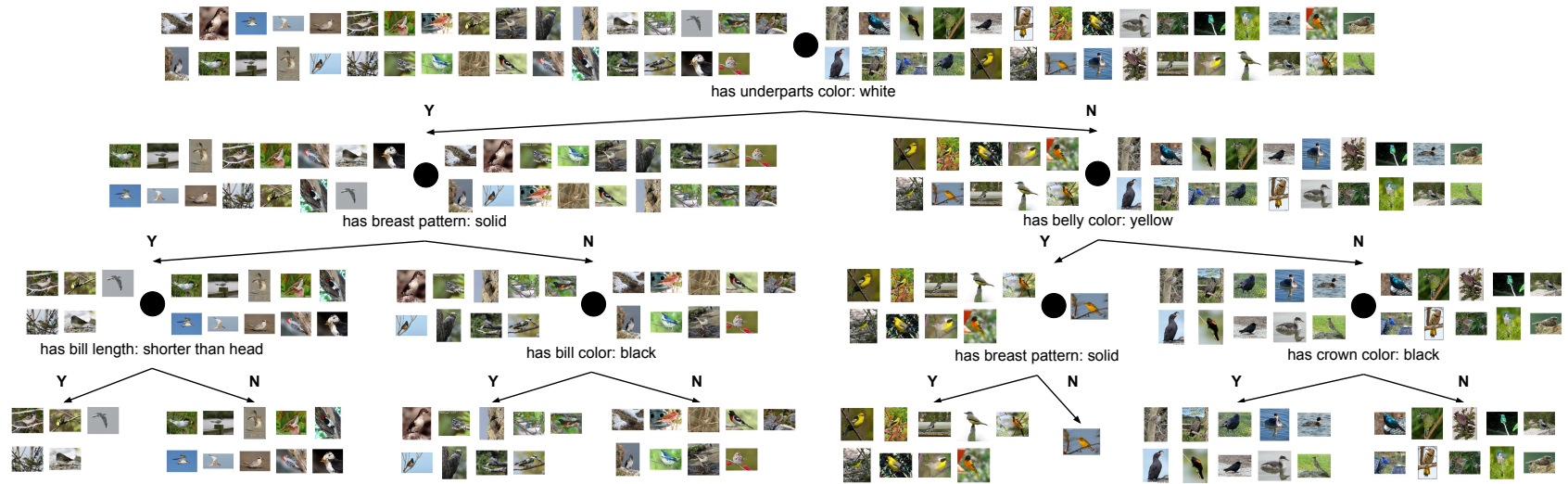
Figure 11: Our aRDTC learns explainable decisions via the decision tree showing the path for each class and it also gives each decision a human-understandable meaning. Here we show the first three decisions for a subset of the 200 classes of birds in CUB where a randomly selected image from a class represents each class.
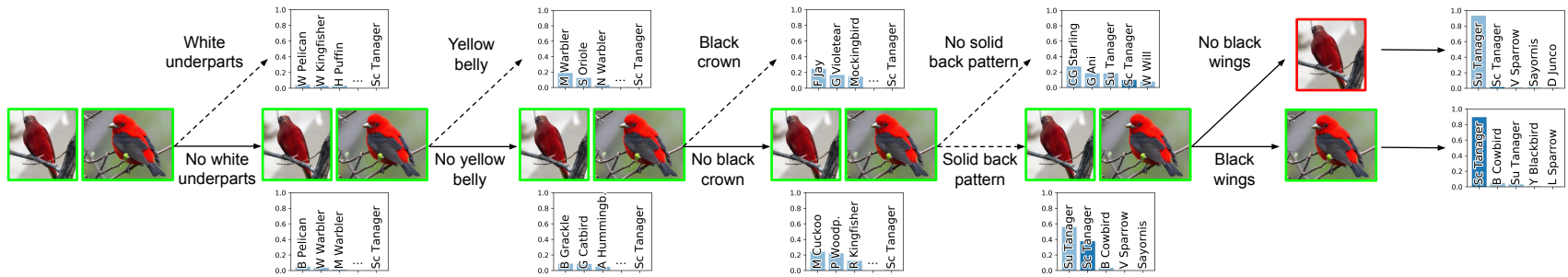


Figure 12: Our aRDTC points to the reasoning behind a wrong decision. Here we illustrate two images from the "Green Kingfisher" class. The lower path lead to a correct classification. Both images follow the same path except for the decision of "black wings". The flying bird gets classified as a "Belted Kingfisher" incorrectly because the black wings are not visible.
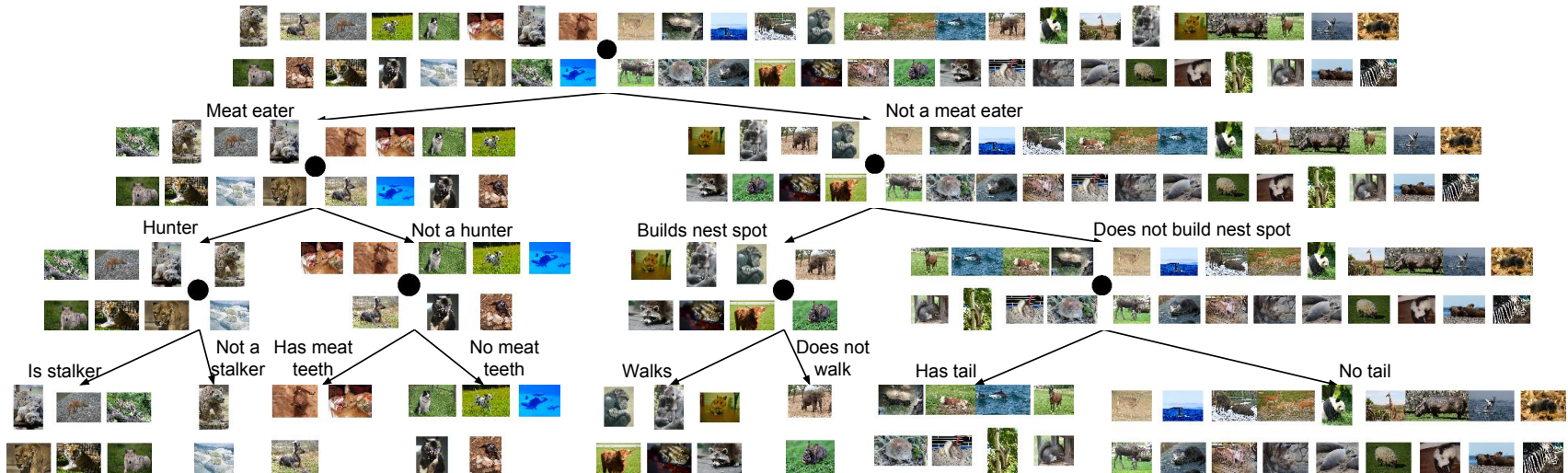
Figure 13: Learned explainable decisions on AWA2 by our aRDTC model. We show the decision tree of the most likely path for each class, *i.e.*, introspection, and give each decision a human-understandable meaning, *i.e.*, rationalization. The tree exposes the thought process of our model, *e.g.*, it decides to separate meat-eating animals from all other animals in the first step.
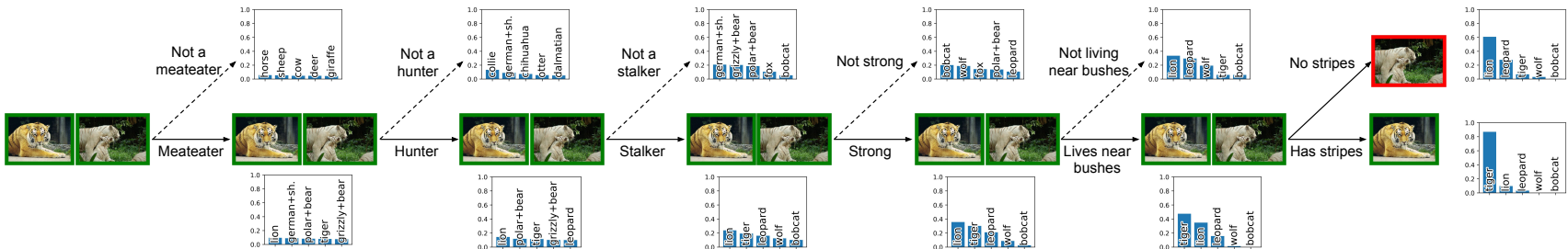


Figure 14: Decision process for two tiger images in AWA2 along with the current label prediction at each step. The lower (upper) path is taken when the attribute is present (absent) for a given class. Both images follow the same path except for the last decision, "has stripes". Since these are absent in the white tiger, it gets classified incorrectly as a lion.
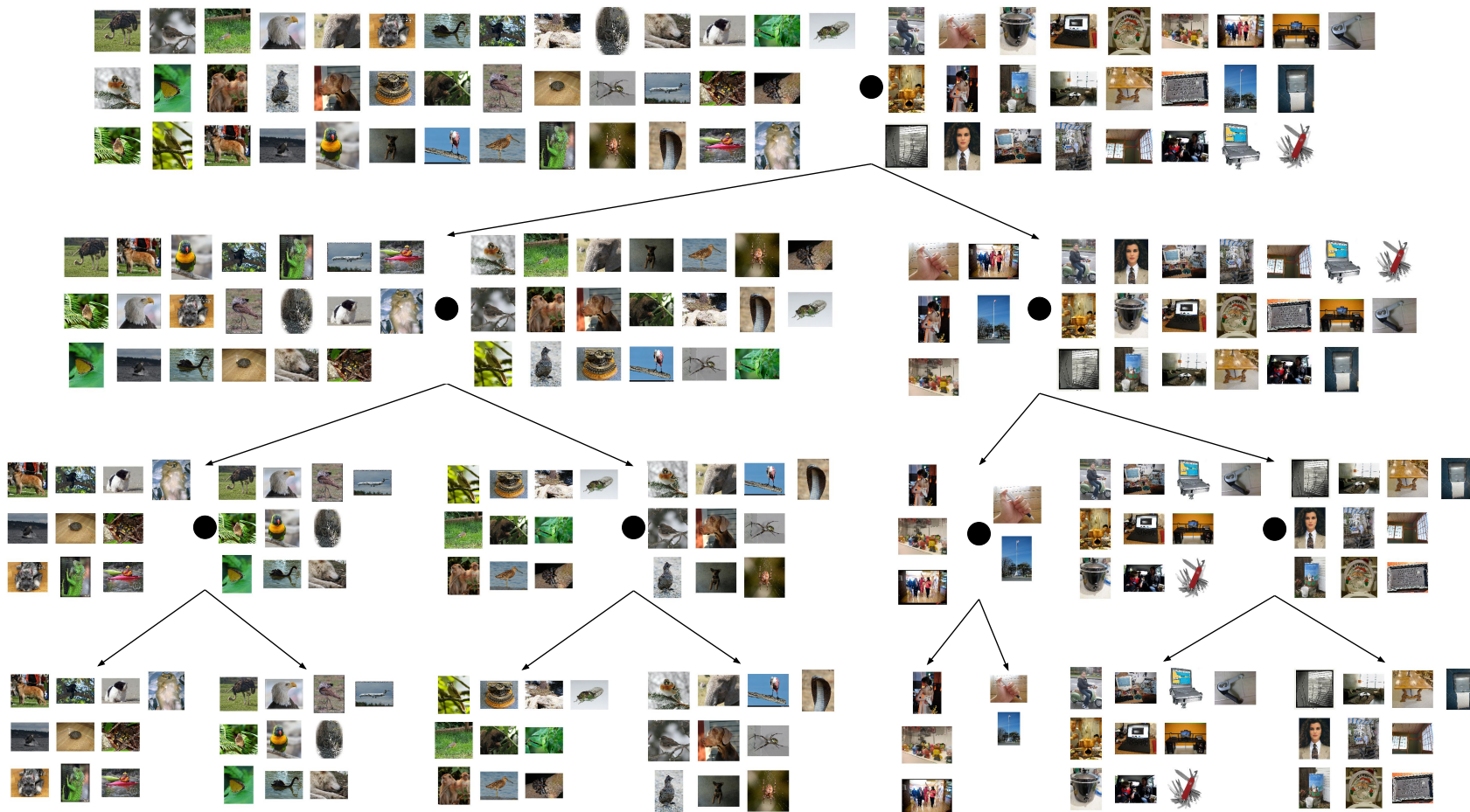
Figure 15: Learned binary decisions on ImageNet by our RDTC model. A subset of randomly chosen classes are shown by one representative image of each class. Our tree reveals a clustering as decision splits narrow down towards a specific subset of classes.