

Towards Deep Machine Reasoning: a Prototype-based Deep Neural Network with Decision Tree Inference

Plamen Angelov, Eduardo Soares

Abstract—In this paper we introduce the DMR – a prototype-based method and network architecture for deep learning which is using a decision tree (DT)- based inference and synthetic data to balance the classes. It builds upon the recently introduced xDNN method [1] addressing more complex multi-class problems, specifically when classes are highly imbalanced. DMR moves away from a direct decision based on all classes towards a layered DT of pair-wise class comparisons. In addition, it forces the prototypes to be balanced between classes regardless of possible class imbalances of the training data. It has two novel mechanisms, namely i) using a DT to determine the winning class label, and ii) balancing the classes by synthesizing data around the prototypes determined from the available training data. As a result, we improved significantly the performance of the resulting fully explainable DNN as evidenced by the best reported result on the well know benchmark problem Caltech-101 surpassing our own recently published "world record". Furthermore, we also achieved another "world record" for another very hard benchmark problem, namely Caltech-256 as well as surpassed the results of other approaches on Faces-1999 problem. In summary, we propose a new approach specifically advantageous for imbalanced multi-class problems that achieved two world records on well known hard benchmark problems and the best result on another problem in terms of accuracy. Moreover, DMR offers full explainability, does not require GPUs and can continue to learn from new data by adding new prototypes preserving the previous ones but not requiring full retraining.

I. INTRODUCTION

In this paper we introduce a new deep neural network (DNN) which departs from the amorphous and highly abstract, "black box" model structure towards deep machine reasoning (DMR) architecture. This is based on the following principle differences from the traditional approach: i) use of prototypes as the core of the method; ii) use of a DT for decision making (class labeling) instead of a flat "winner takes all" type function; iii) using similarity as a measure of association to prototypes; iv) possibility to express the method in a form of human-interpretable IF-THEN rules with partial degree of satisfaction and to visualise by Voronoi tessellation or by prototypes.

The staggering increase of the amount and complexity of the data sets and streams led to a move from rule-based systems (fuzzy, Bayesian inference, Markov decision processes, using Q tables in reinforcement learning, case base reasoning, etc.) towards DNN which have proven their efficiency in a number of problems ranging from speech,

image recognition and language translation to games [2]. This abundance of data led, however, to the temptation to shortcut from data to the solutions driven entirely by the accuracy and ignoring the depth of understanding the problem at hand, and getting insights.

In DRM we make use of the strong properties of the DNN and add new mechanisms to address their shortcomings. For example, the DNN are very efficient feature extractors, especially for image processing problems [2]. We use this in DRM and we also use layered structure/architecture. We further benefit from the transfer learning approach [3].

Traditional classifiers assume balanced classes, but in practice classes are usually (highly) imbalanced. For example, in fault detection and identification the amount of data about the faulty cases are usually significantly smaller than the amount of data for "normal" operation [4]. In social applications, for example, this leads to possible un-fairness [5] when the data is highly imbalanced with dominating class(es) and minority class(es).

Finally, traditional statistical modelling is heavily influenced by averages and starts with assumptions about the data distributions which are then put to a test by parametrisation [6]. We take the opposite approach starting with the observed data samples and generalise from these local densities and global multivariate generative distributions. These empirically derived distributions have discrete and continuous form [7]. Their discrete forms are exact while the continuous forms which are needed for the inference are local estimates.

Prototype-based models have demonstrated their high efficiency, e.g. the discriminative models such as kNN [8], SVM [9], less so RBF [10] and LVQ [11]. The latter two are also good in terms of explainability [1]. Explainability is undoubtedly, the Achilles heel of the DNN and the solution we propose is to have a synergy between reasoning and learning rather than the current dichotomy.

In this paper we offer a new deep learning architecture and method that builds upon our recently introduced xDNN [1] method by adding two important novelties, namely: i) using a DT to determine the winning class label, and ii) balancing the classes by synthesising data around the prototypes determined from the available training data.

We validated the new DMR method on three well known benchmark problems, namely Faces-1999, Caltech-101 and Caltech-256. Both Caltech problems are very hard and there is a public record of the best results achieved so far [12]. We surpassed one of them (Caltech-101) with xDNN already [1].

With DMR we surpass our own xDNN "world record". Furthermore, we also surpassed the best record on Caltech-256 as well as on Faces-1999 problems. Moreover, DMR does not require GPUs, computationally lean and can continue to train for new data without the need for full re-training.

The remainder of the paper is organised as follows: Section II introduces the concept and novelties of the proposed approach. Section III presents the proposed architecture used during the training phase. Section IV outlines the learning procedure, section V introduces the architecture of the DMR used during the validation phase. Section VI illustrates explainability of DMR in terms of IF-THEN rules. Numerical experiments are presented in the Section VII, results are analysed in Section VIII and the paper is concluded with Section IX.

II. CONCEPT AND NOVELTIES OF THE PROPOSED APPROACH

The problem we consider in this paper is to design a classifier with deep architecture that is *explainable-by-design* due to the use of prototypes [1]. Prototypes are a small subset of the training data that are highly representative. This is because they are the local peaks of the distribution [7].

Let us denote the training data set of points by $x = \{x_1, \dots, x_N\} \in \mathbb{R}^n$ with corresponding class labels $y_1, \dots, y_C \in \{1, \dots, C\}$. Here, N is the number of training data samples and n is their dimensionality (number of features); C is the number of classes. DMR starts by selecting a set of descriptive prototypes $\pi \in X \subset P$ for each class/per class, M_j is the total number of prototypes of class j ; $M_j = |P_j|$; $M = \sum_{j=1}^C M_j$. Notice that $M_j > 1$ for $\forall j$, i.e. we usually consider more than a single prototype per class. The prototype extraction process (which can be both, offline and online) is described in more detail in [7],[1]. At the heart of practically all prototype-based methods is the concept that the prototypes of class C are designed to be close to many training points of class C and far from training points of the other classes. As pointed out in [1] "This idea captures the sense in which the word *prototypical* is commonly used".

The power of prototype-based approaches stems from the fact that they are *explainable-by-design* [5], easy to understand by the users because they represent samples of the training data, e.g. images. They can be used for classification. Any new data sample with unknown label, $x \in \mathbb{R}^n$ can be associated with the nearest prototype from the sets P_1, P_2, \dots, P_C ; $P = P_1 \cup P_2 \cup \dots \cup P_C$.

$$L(x) = \underset{x \in X}{\operatorname{argmin}} \min_{\pi \in P} d(x, \pi). \quad (1)$$

A. Decision Tree layer

In traditional DNN, the decision is flat, *en bloc* in the form of a single stage "winner takes all" function as in eq. (1) and is the last layer of the network. In xDNN [1] we also followed this popular decision concept, but split it into two stages: i) per class winner, and ii) across classes global decision. In DMR, similarly to xDNN [1] the decision mechanism is part of the architecture used for validation of the results because

the training is per class and no decision for the class label is needed during the training. In this paper, the proposed DMR is using a multi-layer DT formed by pairwise comparison of top two classes in terms of minimum error in training as detailed in Section V and Fig. (4). The reason the result is significantly different is that the Voronoi tessellation regions of the data clouds that are formed around each prototype (local zones of influence) are significantly different when binary decisions are made.

B. Balancing classes through synthesising training data strategically

The second innovation of the proposed method is related to the balancing of the classes. We achieve this by synthetic data augmentation. In this paper we propose a different approach from our recently published one [13] for synthesising data for highly imbalanced classification problems. The differences are that in this paper we synthesise data around prototypes which makes these synthetic data more likely to have the same class as the prototype. The method starts by identifying a population of pairwise neighbouring data samples from minority classes around prototypes. Then, it imposes a Gaussian disturbance on these data samples, and, finally, it generates synthetic samples by creating linear interpolations between these extrapolations. A further difference from our recent method [13] is that in this paper we use the standard deviation, σ as a radius of influence around the prototype rather than absolute distance of first order. We then augment the training data set with this synthetically generated data set as shown in Fig. (1), see the augmented prototypes layer.

III. ARCHITECTURE OF THE PROPOSED DMR APPROACH (DURING THE TRAINING PHASE)

The architecture of the proposed classifier can be represented as a multi-layered DNN with a very clear semantic and functional meaning by design. The architecture for the training and for the validation phases are different as detailed in Figs. 1 and 4. The training phase is performed per class (except the last layer) and includes the following layers:

1) Input (features) layer

This is the first layer which defines the data space. The number of inputs is determined by the nature of the problem that the data describe. In many problems these are clearly known physical or biomedical variables, e.g. velocities, pressure, temperature, etc. In image processing problems traditionally size, shape of objects or HoG [14] were used as well as more abstract methods like GIST [15]. More recently, convolutional neural networks (CNN) like AlexNet [16], VGG-VD-16 [17], Inception [18], ResNet [19], Inception-Resnet [20] have proven to be very efficient to encode images and represent them as a highly abstract vector of the outputs from the Fully Connected Layer (FCL). The proposed DMR architecture is agnostic to the source of the features vector that the input layer represents. It can be any of the above. In this paper without any loss of generality we use a 1×4096 dimensional vector formed by the outputs

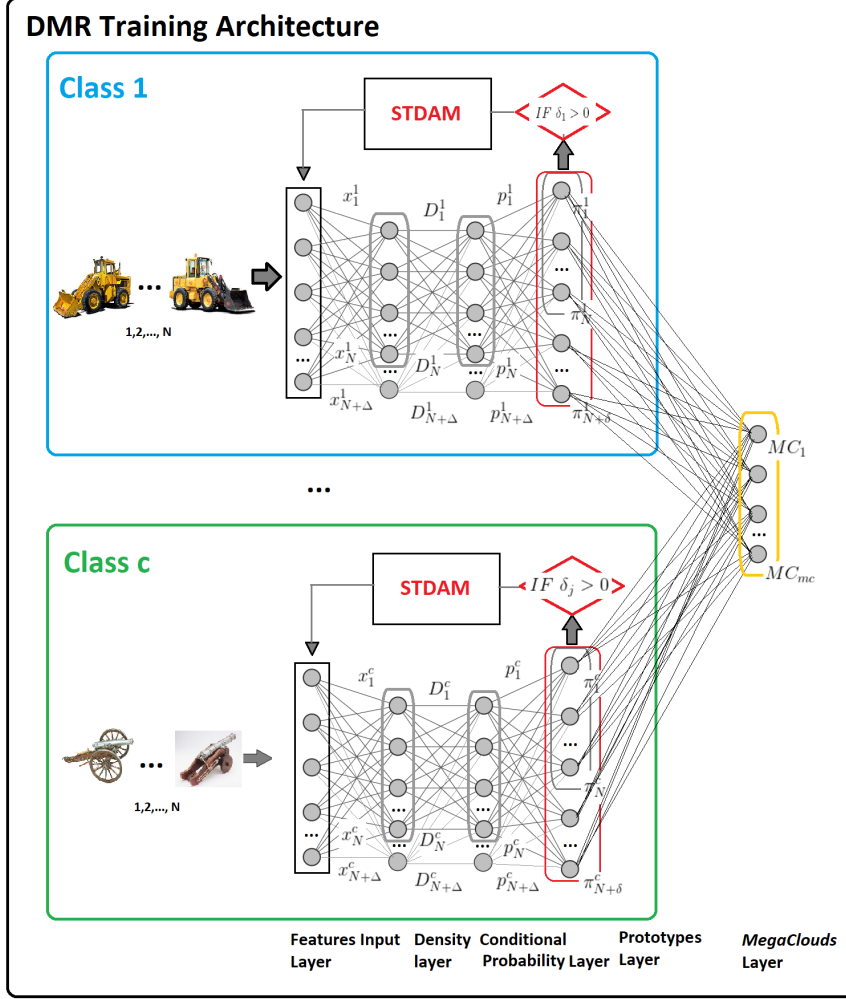


Fig. 1. DMR Architecture during the training phase (STDAM stands for Synthetic Training Data Augmentation Mechanism).

from the first FCL from a VGG–VD–16 pre-trained on Imagenet [21].

2) Data density layer

This layer is composed of neurons whose activation function represents the data density, D defined by a Cauchy function [7]:

$$D(x) = \frac{1}{1 + \frac{\|x - \mu\|^2}{\sigma^2}}, \quad (2)$$

where D is the density, μ is the global mean, and σ is the variance. In [7] it was demonstrated theoretically that starting from the mutual proximity of the data samples in the data space and using Euclidean (or Mahalanobis) type distance D takes the form of a Cauchy function. Moreover, data density can be updated recursively as detailed in [22]. The value of the data density, D represents the closeness to the mean and is in the range $0 < D \leq 1$. It obtains its maximum (of 1) when $x = \mu$.

D is indicative for the centrality of a data sample and its suitability to be a prototype due to its proximity to other data samples.

3) Conditional probability layer

The conditional probability can be estimated from the empirically observed data as described in [7] where it is also called *typicality* τ . It can be given by eq. (3). The integral of $\int_{-\infty}^{\infty} p(C|x) dx = 1$ same as for the pdf [7], but it is multi-modal:

$$p(C|x) = \frac{\sum_{i=1}^C N_i D(x^i)}{\sum_{i=1}^C N_i \int_{-\infty}^{\infty} D(x^i) dx} \quad (3)$$

where N_i denotes the number of data samples associated with (support of) the i -th data cloud, $\sum_{i=1}^C N_i = N$. Notice that since $p(C|x)$ is empirically derived [7] it is not constrained by any *prior* assumptions about the data distribution type or even about the random or deterministic nature of the data. This is clearly more

realistic in comparison with the common approach which (for theoretical convenience) assumes randomness and independence of the features of the experimentally observed data which is usually far from the reality.

4) Prototypes layer

The next layer consists of prototypes, π . This is the core layer of the proposed DMR architecture. This layer is responsible to provide *explainable-by-design* model. Prototypes are the local peaks of the data density (and, respectively, local peaks of the conditional probability, eq. (3)) identified in the previous layers/stages. The proposed DMR algorithm absorbs the new data samples by assigning them to the nearest prototype:

$$j^* = \underset{i=1,\dots,N; j=1,\dots,M}{\operatorname{argmin}} |x_i - \pi_j| \quad (4)$$

In this way, each prototype forms a *cloud* of data that it represents. These "data clouds" form Voronoi tessellation, illustrated in Fig.2

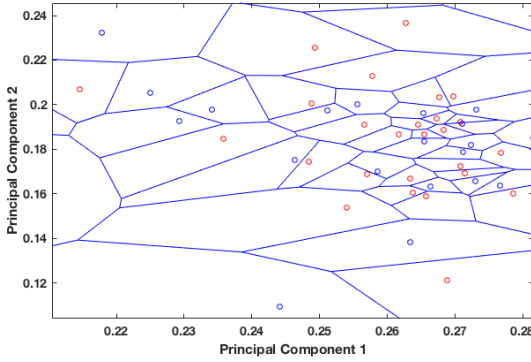


Fig. 2. Identified prototypes – Voronoi Tessellation. The blue circle represents "class 1" and the red circle denotes "class 2".

The prototypes are independent from each other. Therefore, one can change the structure by adding a new prototype without influencing the other already existing prototypes. In other words, the proposed DMR network is highly parallelizable and suitable for dynamically evolving applications with non-stationary data streams and evolving data patterns where new prototypes may be added if the data pattern requires this. The proposed DMR network is trained *per class* forming a set of prototypes *per class*. Therefore, all the calculations are done for each class separately. New prototypes are added to this layer when the following condition is met [7]:

$$\begin{aligned} & \text{IF } (D(x) \geq \max_{j=1,\dots,M} D(\pi_j)) \\ & \text{OR } (D(x) \leq \min_{j=1,\dots,M} D(\pi_j)) \quad (5) \\ & \text{THEN (add a new data cloud } (j \leftarrow j + 1)) \end{aligned}$$

If that is the case, then the vector of features of the current training data sample becomes a new prototype, π_{j+1} forms a new *data cloud* [1].

5) Synthetic data augmentation

This mechanism is not a separate layer, but a feedback process that gets information from the prototypes layer, augments the training data set (in the form of synthetically added features vectors close to the existing prototypes) and expands the size of the prototypes layer by balancing the amount of prototypes per class. This mechanism is one of the two novelties of the proposed approach in comparison with our recent xDNN [1] method. The rationale for and the main functionality of this mechanism has been described in Section II.B. In fact, this is an augmentation of the amounts of training data (by augmenting N to $N + \Delta$ made by feeding back the information from the prototypes layer. As a result, the size of the prototypes layer is expanded (by δ) so that the number of prototypes per class is being balanced. This is visualised in Fig. (1) where the red solid rectangle includes the black dotted one (original prototypes) but also adds prototypes which result from adding synthetic training data.

6) MegaClouds layer

This is the final layer of the training architecture. Unlike the previous layers it is cross-class. At this layer prototypes from all classes are put together and once this is done all the adjacent *data clouds* that have the same class label are combined into *mega-clouds*, see Fig.(3). Notice that the number of *megaclouds*, $i = 1, 2, \dots, MG$ is significantly smaller than the number of prototypes, ($MG \ll M$) and the interpretability improves significantly.

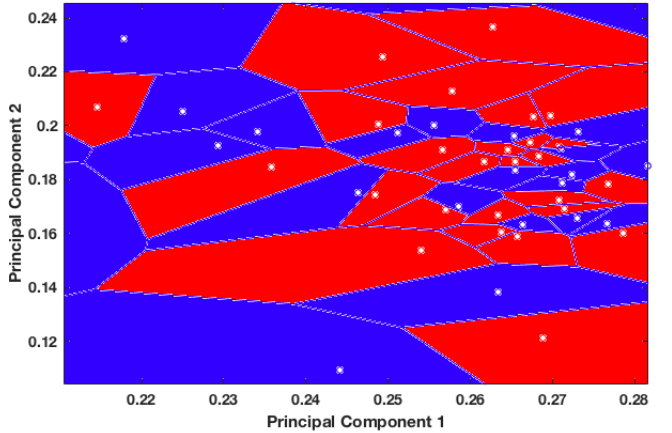


Fig. 3. *Mega-Clouds* are result of merging adjacent *data clouds* which has the same class label.

IV. LEARNING PROCEDURE

The learning of DMR is summarised below by the following pseudo-code. The proposed architecture is feed-forward with the exception of the synthetic data augmentation mechanism which feeds back from the prototype layer back to the input layer. The proposed method can work both, in a batch mode as well as on a per sample basis, online.

DMR: Learning Procedure

- 1: Read the first feature vector sample x_i of class c ;
 - 2: Standardise and normalise the data as detailed in [7]
 - 3: Set $i \leftarrow 1; j \leftarrow 1; \pi_1 \leftarrow x_i; \mu \leftarrow x_1; N \leftarrow 1$
 - 4: **FOR** $i = 2, \dots$
 - 5: Read x_i ;
 - 6: Calculate $D(x_i)$ and $D(\pi_j)$ ($j = 1, 2, \dots, M$) according to eq. (2);
 - 7: **IF** eq. (5) holds
 - 8: Create new prototype: $j \leftarrow j+1; \pi_j \leftarrow x_i; N \leftarrow N+1$
 - 9: **ELSE**
 - 10: Search for the nearest prototype according to eq. (4);
 - 11: Update the nearest prototype as:
 $N \leftarrow N + 1;$
 $\pi_j \leftarrow \frac{N_j}{N_j+1}\pi_j + \frac{N_j}{N_j+1}x_i;$
 - 12: Balance the number of prototypes through synthetic data augmentation mechanism detailed below;
 - 13: **END**
 - 14: **END**
-

Synthetic Data Generation

- 1: **FOR** $j = 1, 2, \dots, C$ **DO**
- 2: Calculate the amount of synthetic data samples needed to balance the pair of classes j and $j+1$: $\delta = M_j - M_{j+1}$.
- 3: **UNTIL** $\delta = 0$ **DO**
- 4: $k = 1$
- 5: Randomly select a pair of neighbouring data samples $(p_k, q_k)^*$ from the 0.3σ zone around the prototype from the minority class;
- 6: Apply Gaussian disturbance to $(p_k, q_k)^*$ by eq. (6) and obtain $(\hat{p}_k, \hat{q}_k)^*$ [23];

$$(\hat{p}_k, \hat{q}_k)^* = (p_k + g_p, q_k + g_q)_k^* \quad (6)$$

where $g_p = [g_{p,1}, g_{p,2}, \dots, g_{p,R}]^T$ and $G_q = [g_{q,1}, g_{q,2}, \dots, g_{q,R}]^T$ are two R dimensional randomly generated vectors sampled from the Gaussian distributions, $g_{p,l}, g_{q,l} \sim \mathbb{N}(0, \sigma)$ ($l = 1, 2, \dots, R$) with σ being the standard deviation.

- 7: Create random interpolation ρ_k between $(\hat{p}_k, \hat{q}_k)^*$ as follows [13]:

$$\rho_k = \alpha_k^T \hat{p}_k + (1 - \alpha_k)^T \hat{q}_k \quad (7)$$

where $\alpha_k = [\alpha_{k,1}, \alpha_{k,2}, \dots, \alpha_{k,R}]^T$ is a R dimensional random vector, elements of which follows the uniform distribution within the range $[0,1]$.

- 8: $k \leftarrow k + 1$
 - 9: **END UNTIL**
 - 10: **END FOR**
-

The architecture of DMR for the validation phase (see Fig. 4) has the following layers.

- 1) **Input (features) layer** The first layer is exact the same as in the training phase and has been described in section III.

- 2) **Ranked prototypes layer**

In this layer we rank order all the prototypes in terms of minimum error during the training. Then we organise them in overlapping pairs: we start with the top two prototypes (providing smaller error) and then the pair of the second best and the third; further on, the pair of the third and the fourth, etc. In this way, all prototypes take part twice except the best one and the worst one, see Fig. (4). The output of this layer is the degree of similarity, S between the unlabeled data sample and the respective prototype. The activation functions of the neurons of this layer are defined as follows:

$$S_j = S(x_i, \pi_j) = \frac{1}{1 + \frac{(x - \pi_j)^2}{\sigma_j^2}}, \quad (8)$$

where $j = 1, 2, \dots, M; i = 1, 2, \dots, N$. It is easy to see that for similarity we use the same Cauchy function as the data density, eq. (2).

- 3) **Maximum similarity layer**

Each neuron of this layer is performing a simple max operation over the pair of similarity values that are coming from the previous layer, namely:

$$S_{j,j-1}^* = \max(S_{j-1}, S_j) \quad (9)$$

The winner goes forward.

- 4) **Pair-wise confidence checks layer**

In this layer we check if the confidence in the best of the two potential outcomes is high enough. In this paper we use a threshold, $Thr=0.9$, which means 90% similarity of the new, unlabeled data sample to any prototype. The neurons of this layer are linked between each other forming a competitive layer. This link is activated if the confidence check fails (see Fig. 2). The flow of the information to the next layer is conditional on the outcome from the confidence check. First, the top two pairs of prototypes are checked. If the winner surpasses Thr it is the winner. Otherwise, the flow goes down to the next pair (in the same layer of the network, the key Fig. 4 is closed) and so on.

IF ($\min(S_{j,j-1}^*, S_{j-1,j-2}^*) \geq Thr$) *THEN* (Step 4)
ELSE (Step 3)

- 5) **Pair-wise winners layer**

Pair-wise decisions are made to determine the winning prototype from the candidate pair (S_{j-1}^*, S_j^*) , which passed the confidence check in the proceeding layer.

$$Label = \operatorname{argmax}(S_{j,j-1}^*, S_{j-1,j-2}^*) \quad (10)$$

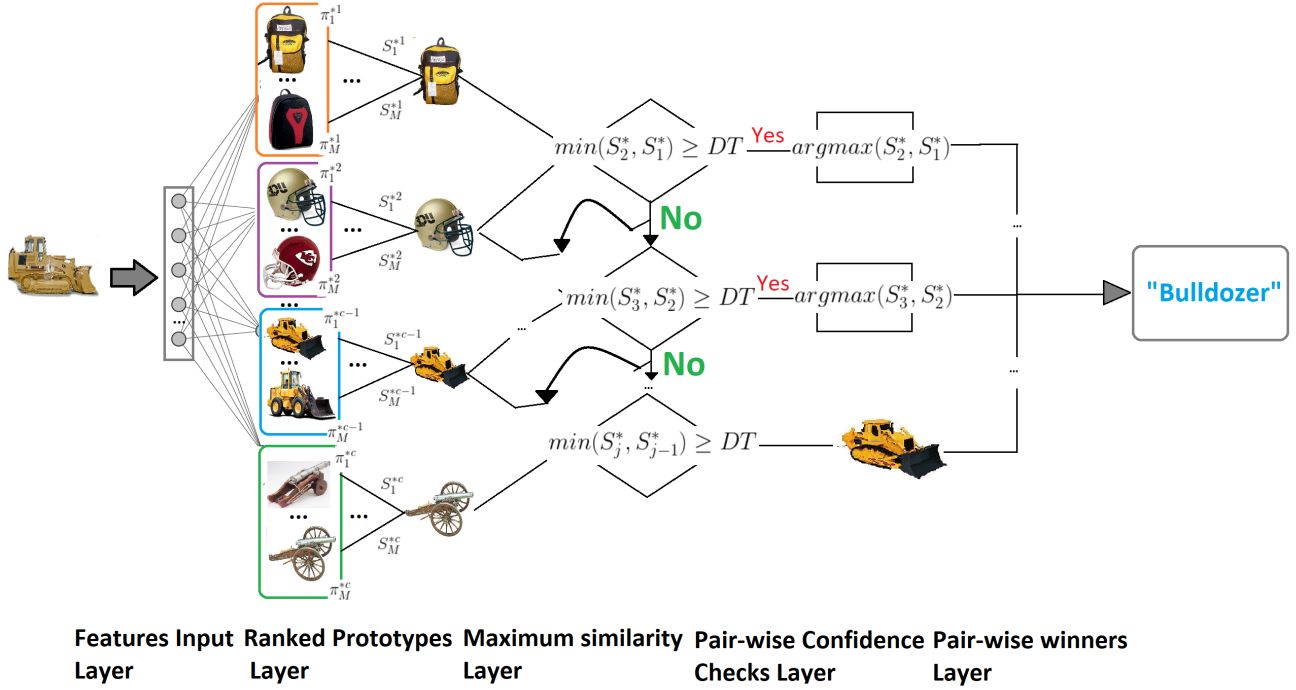





Fig. 4. Architecture for the validation process of the proposed DMR approach.

VI. EXPLAINING THE DMR NETWORK AS A SET OF IF...THEN RULES

One of the main advantages of the proposed DMR approach is that it is *explainable-by-design* and can be represented, for example, in the form of IF...THEN rules [22]. People can easily understand rules and prototypes. These are often easy to visualise, e.g. in case of images and can also be expressed as a set of linguistic rules as follows:

IF (Image \sim ) THEN "Bulldozer"

where \sim denotes "similar to"; it can also be seen as a fuzzy degree of membership. One rule per prototype can be formed. All rules per class can be combined together using logical OR, also known as disjunction or S-norm:

IF (Image \sim ) OR (Image \sim ) OR
 ... OR (Image \sim ) THEN "Bulldozer"

VII. NUMERICAL EXPERIMENTS

We validated our proposed approach, DMR using several complex, well-known image classification benchmark data

sets (Faces-1999, Caltech-101, and Caltech-256). Description of the data sets are given below:

1) *Faces-1999*: The Faces-1999 data set [24] contains 450 frontal real faces images from 27 different people. This data set is highly unbalanced.

2) *Caltech-101*: The Caltech-101 data set [25] contains 9144 images in divided into 102 categories(one background). The Caltech-101 dataset is highly unbalanced and is widely used as bench marking data set.

A. Caltech-256

Caletch-256 has 30,607 images divided into 257 object categories (one of which is the background) [25].

B. Performance Evaluation

The performance of the classification methods is usually evaluated based on their accuracy index which is defined as follows:

$$ACC(\%) = \frac{TP + TN}{TP + FP + TN + FN}, \quad (11)$$

where TP, FP, TN, FN denote true and false, negative and positive, respectively.

All the experiments were conducted with MATLAB 2018a using a personal computer with a 1.8 GHz Intel Core i5 processor, 8-GB RAM, and MacOS operating system. The classification experiments were executed using 10-fold cross validation under the same ratio of training-to-testing (80% to 20%) sample sets.

VIII. RESULTS AND ANALYSIS

Computational simulations were performed to assess the accuracy of the proposed explainable tree-based deep learning method (DMR), against other state-of-the-art approaches.

A. Faces Data set

Table I shows that the proposed DMR method provides the best result in terms of classification accuracy than its state-of-the-art competitors. The number of model parameters for DMR (and xDNN) is, strictly speaking, zero, because the 2 parameters (mean, μ and standard deviation, σ) per prototype (*data cloud*) are derived from the data and are not algorithmic parameters or user-defined parameters. However, the tree-based structure of the proposed DMR and the mechanism for balancing the classes allow the result to surpass all others. The proposed deep reasoning through a layered pair-wise DT is exploiting and benefiting from the old principle of *divide et impera*.

TABLE I
PERFORMANCE COMPARISON: FACES-1999 DATA SET

Method	Accuracy
DMR	96.71%
VGG-VD-16	96.32%
xDNN	96.15%
VGG-VD-16	96.32%
SVM	95.51%
KNN	88.54%
DT	61.53%

B. Caltech-101 Data set

Table II shows the results considering the challenging Caltech-101 data set. It is possible to note through Table II that the proposed DMR method provides the best result in terms of classification accuracy. The proposed Caltech-101 is hugely unbalanced, and the inner data augmentation mechanism of the proposed DMR method favour the balance of the data, consequently, it improves the final classification result. Moreover, the intelligent tree-based structure of the proposed method allows interpretability and also favours the improvement in the classification accuracy of the given model.

The proposed explainable tree-based DNN surpasses in terms of accuracy the state-of-the-art VGG-VD-16 algorithm which is a well-established convolutional deep neural network. Moreover, it could also surpass other state-of-art approaches.

TABLE II
PERFORMANCE COMPARISON: CALTECH-101 DATA SET

Method	Accuracy
DMR	94.31%
SPP-net	91.44%
xDNN	90.62%
VGG-VD-16	90.32%
KNN	85.65%
DT	54.42%

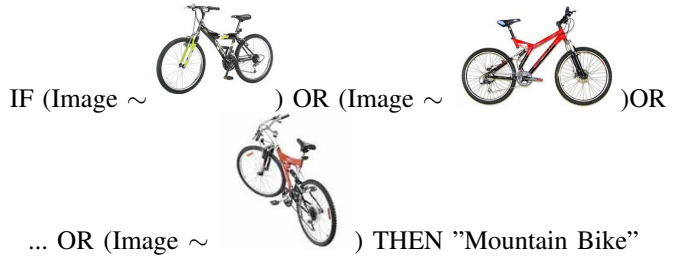
C. Caltech-256 Data set

Results for Caltech-256 are presented in Table III.

TABLE III
PERFORMANCE COMPARISON: CALTECH-256 DATA SET

Method	Accuracy
DMR	77.54%
xDNN [1]	75.41%
SVM(1) [26]	24.6%
SVM(2) [26]	39.6%
SVM(3) [26]	46.0%
SVM(4) [26]	51.3%
SVM(5) [26]	65.6%
SVM(7) [26]	71.7%
Softmax(5)[26]	65.7%
Softmax(7) [26]	74.2%

These results demonstrate that the proposed DMR approach obtains the best classification accuracy ever reported for this complex problem, namely, 77.54%. The proposed approach not only surpasses all published competitors but also offers a clearly explainable model.



DMR even surpasses the recently introduced by us xDNN approach [1], which reported the world best result on 5 December 2019 for this classification problem.

IX. CONCLUSION

In this paper we introduce the DMR – a prototype-based explainable DNN with DT inference and balanced amount of prototypes per class regardless of the possible imbalances of the training data. The proposed method offers two main novelties, namely: i) using a DT to determine the winning class label, and ii) balancing the classes by synthesising data around the prototypes determined from the available training data. It demonstrates excellent performance surpassing three well known benchmark problems (Caltech-101, Caltech-256 and Faces-1999) where the first two are the the best results published. The proposed approach is explainable-by-design, computationally efficient (no need for GPUs, high degree of parallelization possible, no iterative search procedures and parameter optimisation). Furthermore, it offers the ability to learn continuously (live-long) adapting smoothly to new data patterns. It is a step towards bringing closer machine learning and automated reasoning into what we call *deep machine reasoning* aiming not only high levels of accuracy but also deeper understanding and insight.

REFERENCES

- [1] P. Angelov and E. Soares, "Towards explainable deep neural networks (xDNN)," 2019.
- [2] I. Goodfellow, Y. Bengio, and A. Courville, *Deep learning*. MIT press, 2016.
- [3] L. Y. Pratt, "Discriminability-based transfer between neural networks," in *Advances in neural information processing systems*, 1993, pp. 204–211.
- [4] B. S. J. Costa, P. P. Angelov, and L. A. Guedes, "Fully unsupervised fault detection and identification based on recursive density estimation and self-evolving cloud-based classifier," *Neurocomputing*, vol. 150, pp. 289–303, 2015.
- [5] E. Soares and P. Angelov, "Fair-by-design explainable models for prediction of recidivism," *arXiv preprint arXiv:1910.02043*, 2019.
- [6] T. Hastie, R. Tibshirani, and J. Friedman, *The elements of statistical learning: data mining, inference, and prediction*. Springer Science & Business Media, 2009.
- [7] P. P. Angelov and X. Gu, *Empirical approach to machine learning*. Springer, 2019.
- [8] S. Zhang, X. Li, M. Zong, X. Zhu, and D. Cheng, "Learning k for knn classification," *ACM Transactions on Intelligent Systems and Technology (TIST)*, vol. 8, no. 3, pp. 1–19, 2017.
- [9] J. A. Suykens and J. Vandewalle, "Least squares support vector machine classifiers," *Neural processing letters*, vol. 9, no. 3, pp. 293–300, 1999.
- [10] M. J. Er, S. Wu, J. Lu, and H. L. Toh, "Face recognition with radial basis function (rbf) neural networks," *IEEE transactions on neural networks*, vol. 13, no. 3, pp. 697–710, 2002.
- [11] T. Kohonen, J. Hynninen, J. Kangas, J. Laaksonen, and K. Torkkola, "Lpq pak: The learning vector quantization program package," Technical report, Tech. Rep., 1996.
- [12] K. He, X. Zhang, S. Ren, and J. Sun, "Spatial pyramid pooling in deep convolutional networks for visual recognition," *IEEE transactions on pattern analysis and machine intelligence*, vol. 37, no. 9, pp. 1904–1916, 2015.
- [13] X. Gu, P. P. Angelov, and E. A. Soares, "A self-adaptive synthetic over-sampling technique for imbalanced classification," *arXiv preprint arXiv:1911.11018*, 2019.
- [14] K. Mizuno, Y. Terachi, K. Takagi, S. Izumi, H. Kawaguchi, and M. Yoshimoto, "Architectural study of hog feature extraction processor for real-time object detection," in *2012 IEEE Workshop on Signal Processing Systems*. IEEE, 2012, pp. 197–202.
- [15] B. Solmaz, S. M. Assari, and M. Shah, "Classifying web videos using a global video descriptor," *Machine vision and applications*, vol. 24, no. 7, pp. 1473–1485, 2013.
- [16] A. Krizhevsky, I. Sutskever, and G. E. Hinton, "Imagenet classification with deep convolutional neural networks," in *Advances in neural information processing systems*, 2012, pp. 1097–1105.
- [17] K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image recognition," *arXiv preprint arXiv:1409.1556*, 2014.
- [18] C. Szegedy, W. Liu, Y. Jia, P. Sermanet, S. Reed, D. Anguelov, D. Erhan, V. Vanhoucke, and A. Rabinovich, "Going deeper with convolutions," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2015, pp. 1–9.
- [19] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 770–778.
- [20] C. Szegedy, S. Ioffe, V. Vanhoucke, and A. A. Alemi, "Inception-v4, inception-resnet and the impact of residual connections on learning," in *Thirty-First AAAI Conference on Artificial Intelligence*, 2017.
- [21] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei, "Imagenet: A large-scale hierarchical image database," in *2009 IEEE conference on computer vision and pattern recognition*. Ieee, 2009, pp. 248–255.
- [22] P. Angelov, *Autonomous learning systems: from data streams to knowledge in real-time*. John Wiley & Sons, 2012.
- [23] J. S. Freudenberg, R. H. Middleton, and V. Solo, "Stabilization and disturbance attenuation over a gaussian communication channel," *IEEE Transactions on Automatic Control*, vol. 55, no. 3, pp. 795–799, 2010.
- [24] M. Weber and M. Weber, "Caltech frontal face database," *California Institute of Technology*, 1999.
- [25] G. Griffin, A. Holub, and P. Perona, "Caltech-256 object category dataset," 2007.
- [26] M. D. Zeiler and R. Fergus, "Visualizing and understanding convolutional networks," in *European conference on computer vision*. Springer, 2014, pp. 818–833.