

Towards Explainable Deep Neural Networks (xDNN)

Plamen Angelov^a, Eduardo Soares^a

^a*School of Computing and Communications, LIRA Research Centre, Lancaster University,
Lancaster, LA1 4WA, UK*

E-mail: p.angelov@lancaster.ac.uk; e.almeidasoares@lancaster.ac.uk

Abstract

1 In this paper, we propose an elegant solution that is directly addressing the
2 bottlenecks of the traditional deep learning approaches and offers an explainable
3 internal architecture that can outperform the existing methods, requires very
4 little computational resources (no need for GPUs) and short training times
5 (in the order of seconds). The proposed approach, xDNN is using prototypes.
6 Prototypes are actual training data samples (images), which are local peaks of
7 the empirical data distribution called *typicality* as well as of the data density.
8 This generative model is identified in a closed form and equates to the pdf but
9 is derived automatically and entirely from the training data with no user- or
10 problem-specific thresholds, parameters or intervention. The proposed xDNN
11 offers a new deep learning architecture that combines reasoning and learning in
12 a synergy. It is non-iterative and non-parametric, which explains its efficiency
13 in terms of time and computational resources. From the user perspective, the
14 proposed approach is clearly understandable to human users. **We tested it**
15 **on challenging problems as the classification of different lighting conditions for**
16 **driving scenes (iROADS), object detection (Caltech-256, and Caltech-101), and**
17 **SARS-CoV-2 identification via computed tomography scan (COVID CT-scans**
18 **dataset). xDNN outperforms the other methods including deep learning in**
19 **terms of accuracy, time to train and offers an explainable classifier.**

Keywords:

Explainable AI, Interpretability, Prototype-based Models, Deep-Learning.

20 1. Introduction

21 Deep learning has demonstrated ability to achieve highly accurate results in
22 different application domains such as speech recognition (Xiong et al., 2018),
23 image recognition (He et al., 2016), and language translation (LeCun et al.,
24 2015) and other complex problems (Goodfellow et al., 2016). It attracted
25 the attention of media and the wider public (Sejnowski, 2018). It has also
26 proven to be very valuable and efficient in automating the usually laborious
27 and sometimes controversial pre-processing stage of feature extraction. The
28 main criticism towards deep learning is usually related to its ‘black-box’ nature
29 and requirements for huge amount of labeled data, computational resources
30 (GPU accelerators as a standard), long times (hours) of training, high power
31 and energy requirements (Rudin, 2019). Indeed, a traditional deep learning
32 (e.g. convolutional neural network) algorithm involves hundreds of millions of
33 weights/coefficients/parameters that require iterative optimization procedures.
34 In addition, these hundreds of millions of parameters are abstract and detached
35 from the physical nature of the problem being modelled. However, the auto-
36 mated way to extract them is very attractive in high throughput applications of
37 complex problems like image processing where the human expertise may simply
38 be not available or very expensive.

39 Feature extraction is an important pre-processing stage, which defines the
40 data space and may influence the level of accuracy the end result provides.
41 Therefore, we consider this very useful property of the traditional deep learn-
42 ing and step on it combined with another important recent result in the deep
43 learning domain, namely, the transfer learning. This concept postulates that
44 knowledge in the form of a model architecture learned in one context can be
45 re-used and useful in another context (Hu et al., 2015). Transfer learning helps
46 to considerably reduce the amount of time used for training. Moreover, it also
47 may help to improve the accuracy of the models (Zhuang et al., 2015).

48 Stepping on the two main achievements of the deep learning - top accuracy
49 combined with an automatic approach for feature extraction for complex prob-

50 lems, such as image classification, we try to address its deficiencies such as the
51 lack of explainability (Rudin, 2019), computational burden, power and energy
52 resources required, ability to self-adapt and evolve (Soares and Angelov, 2019).
53 Interpretability and explainability are extremely important for high stake appli-
54 cations, such as autonomous cars, medical or court decisions, etc. For example,
55 it is extremely important to know the reasons why a car took some action,
56 especially if this car is involved in an accident (Doshi-Velez and Kim, 2017).

57 The state-of-the-art classifiers offer a choice between higher explainability
58 for the price of lower accuracy or vice versa (Figure 1). Before deep learning
59 (Schmidhuber, 2015), machine-learning and pattern-recognition required sub-
60 stantial domain expertise to model a feature extractor that could transform
61 the raw data into a feature vector which defines the data space within which
62 the learning subsystem could detect or classify data patterns (LeCun et al.,
63 2015). Deep learning offers new way to extract abstract features automatically.
64 Moreover, pre-trained structures can be reused for different tasks through the
65 transfer learning technique (Hu et al., 2015). Transfer learning helps to consid-
66 erably reduce the amount of time used for training, moreover, it also **may help**
67 to improve the accuracy of the models (Zhuang et al., 2015). In this paper,
68 we propose a new approach, xDNN that offers both, high level of explainability
69 combined with the top accuracy.

70 The proposed approach, xDNN offers a new deep learning architecture that
71 combines reasoning and learning in a synergy. It is based on prototypes and
72 the data density (Angelov and Gu, 2019) as well as *typicality* - an empirically
73 derived pdf (Angelov et al., 2017). It is non-iterative and non-parametric, which
74 explains its efficiency in terms of time and computational resources. From the
75 user perspective, the proposed approach is clearly understandable to human
76 users. We tested it on some well-known benchmark data sets such as iRoads
77 (Rezaei and Terauchi, 2013) and Caltech-256 (Griffin et al., 2007) and xDNN
78 outperforms the other methods including deep learning in terms of accuracy,
79 time to train, moreover, offers an explainable classifier. In fact, the result on
80 the very hard Caltech-256 problem (which has 257 classes) represents a world

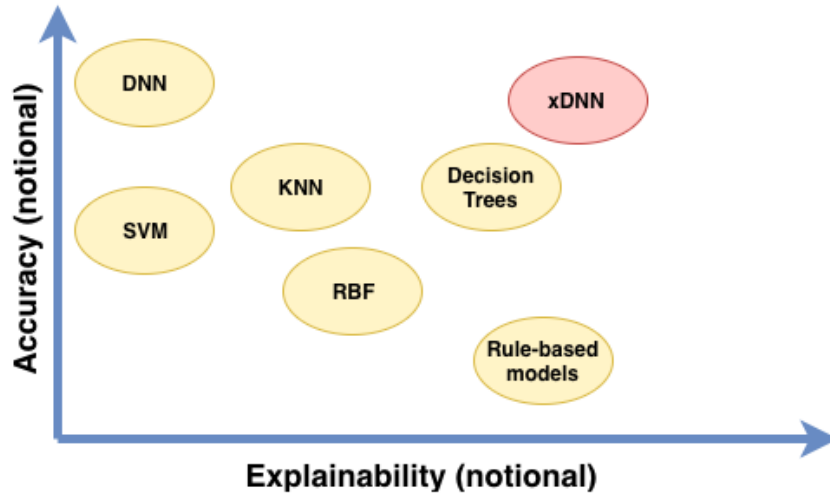


Figure 1: Trade-off between accuracy and explainability.

81 record (He et al., 2015).

82 The remainder of this paper is organized as follows: The next section intro-
 83 duces a brief literature review. The proposed explainable deep learning approach
 84 is presented in Section III. The data employed in the analysis is presented in Sec-
 85 tion IV, and the results are presented in Section V. The discussion is presented
 86 in the last section of this paper.

87 2. Brief Literature Review

88 Deep Neural Networks have often been designed purely for accuracy. The
 89 decisions made by these networks are at best interpreted by *post hoc* techniques
 90 (Li et al., 2018) or not interpreted at all. That is, the first step is the selection
 91 of the network architecture by the human and the attempt to interpret the
 92 trained model and the learned high-level features follows. Therefore, the *post*
 93 *hoc* interpretability analysis requires a separate modeling effort (Saralajew et al.,
 94 2018) and is an approximation rather than a deep explanation of the cause-effect
 95 relations and reasoning. One of the problems with *post hoc* approach is that

96 the explanations can change for different models used. In other words, it is easy
97 to create multiple conflicting yet convincing explanations for how the network
98 would classify a single object.

99 Prototypes-based classifiers are a reasoning process that do not consider *post*
100 *hoc* analysis (Biehl et al., 2016). They rely on the similarity (proximity in the
101 feature space) of a data sample to a given prototype (Biehl et al., 2016, 2013).
102 Different works have different meanings for the word "prototype" (Biehl et al.,
103 2016, 2013, Saralajew et al., 2018), in our case we consider prototypes to be the
104 most representative data samples of the training set (the data samples which
105 have local peaks of the density (Angelov and Gu, 2019)). In other cases, a
106 prototype can be considered as a convex combination of several observations,
107 and not necessarily required to be close to any data sample of the training set
108 or even to be feasible (Oyedotun and Khashman, 2017, Liu et al., 2018).

109 Our work is closely aligned with other prototype classification techniques
110 in machine learning. Prototype classification is a classical form of case-based
111 reasoning (Li et al., 2018); however, as (Li et al., 2018) uses neural networks,
112 the distance measure between prototypes and observations is measured in a
113 latent space. (Li et al., 2018) uses an auto encoder to create a latent low-
114 dimensional space, and distances to prototypes are computed in that latent
115 space. Other works also use Euclidean distance calculation can be expressed in
116 terms of convolution operations in the neural network sense (Nebel et al., 2017,
117 Biehl et al., 2013). This and the computation of the Euclidean distance in terms
118 of a dot product are essential steps towards efficient computational schemes for
119 prototype-based neural network layers.

120 In contrast, the proposed method uses local densities and global multivari-
121 ate generative distributions based on an empirically derived form of the prob-
122 ability distribution function (Angelov and Gu, 2019). Furthermore, differently
123 from other prototype-based classifiers, the presented method is non-iterative
124 and non-parametric as it is using recursive calculations and no search proce-
125 dures. Moreover, the proposed algorithm can learn continuously without full
126 re-training.

127 **3. Explainable Deep Neural Network**

128 *3.1. Architecture and Training of the proposed xDNN*

129 The proposed explainable deep neural network (xDNN) classifier is formed
130 of several layers with a very clear semantic and functional meaning. In addition
131 to the internal clarity and transparency it also offers a very clear from the user
132 point of view set of prototype-based *IF...THEN* rules. Prototypes are selected
133 data samples (images) that the user can easily view, understand and appreciate
134 the similarity to other validation images. xDNN offers a synergy between the
135 statistical learning and reasoning bringing both together. In most of the other
136 approaches there is a dichotomy and preference of one over the other. We
137 advocate and demonstrate that both, learning and reasoning can work together
138 in a synergy and produce very impressive results. Indeed, the proposed xDNN
139 method outperforms all published results (Rezaei and Terauchi, 2013, He et al.,
140 2015, Angelov and Gu, 2018) in terms of accuracy. Moreover, in terms of time
141 for training, computational simplicity, low power and energy required it is also
142 far ahead. The proposed approach can be described as a feedforward neural
143 network which has an incremental learning algorithm that autonomously self-
144 develops and evolves its structure adding new prototypes to reflect the possibly
145 changing (dynamically evolving) data pattern (Soares and Angelov, 2019). As
146 shown in Figure 3, xDNN is composed of the following layers–

- 147 1. Features layer;
- 148 2. Density layer;
- 149 3. Typicality layer;
- 150 4. Prototypes layer;
- 151 5. *MegaClouds* layer;

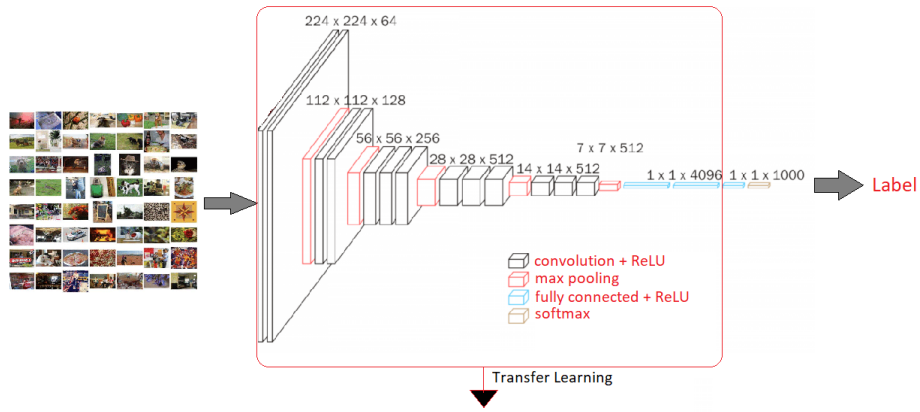


Figure 2: Pre-training a traditional deep neural network (weights of the network are being optimized/trained). Using the transfer learning concept this architecture with the weights are used as feature extractor (the last fully connected layer is considered as a feature vector). Adapted from (Simonyan and Zisserman, 2014).

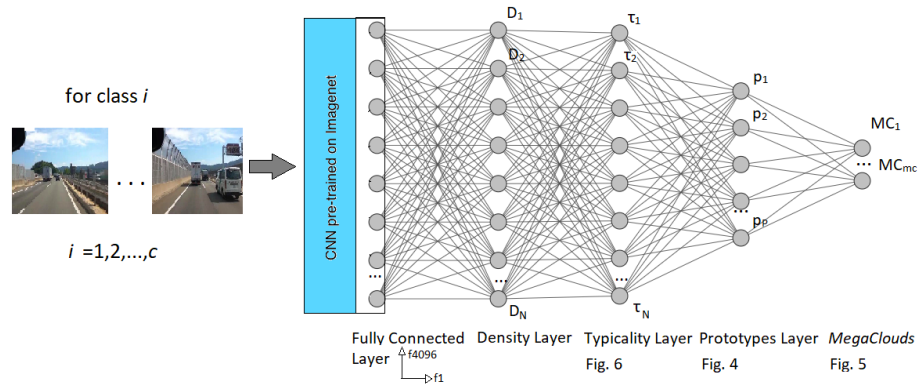


Figure 3: xDNN training architecture (per class).

- 152 1. **Features layer:** (Defines the data space)
- 153 The Feature Layer is the first phase of the proposed xDNN method. This
- 154 layer is in charge of extracting global features vector from the images.
- 155 This first layer can be formed by more traditional ‘handcrafted’ meth-

156 ods such as GIST (Solmaz et al., 2013) or HoG (Mizuno et al., 2012).
 157 Alternatively, it can be formed by the fully connected layer (FCL) of
 158 the pre-trained convolutional neural network approaches such as AlexNet
 159 (Krizhevsky et al., 2012), VGG-VD-16 (Simonyan and Zisserman, 2014),
 160 and Inception (Szegedy et al., 2015), residual neural networks such as
 161 Resnet (He et al., 2016) or Inception-Resnet (Szegedy et al., 2017), etc.
 162 Using pre-trained deep neural network approach allows automatic extrac-
 163 tion of more abstract and discriminative high-level features. In this paper,
 164 pre-trained VGG-VD-16 DCNN is employed for feature extraction. Ac-
 165 cording to (Ren et al., 2016), VGG-VD-16 has a simple structure and
 166 it can achieve a better performance in comparison with other pre-trained
 167 deep neural networks. The first fully connected layer from VGG-VD-16
 168 provides a 1×4096 dimensional vector.

169 a) The values are then standardized using the following equation (1):

$$\hat{x}_{i,j} = \frac{x_{i,j} - \mu(x_{i,j})}{\sigma(x_{i,j})} \quad (1)$$

170 where \hat{x} denotes a standardized features vector x of the image I (x are
 171 the values provided by the FCL), $i = 1, 2, \dots, N$ denotes the time stamp
 172 or the ID of the image, $j = 1, 2, \dots, n$ refers to the number of features of
 173 the given x in our case $n = 4096$.

174 b) The standardized values are normalised to bring them to the range
 175 $[0;1]$:

$$\bar{x}_{i,j} = \frac{\hat{x}_{i,j} - \min_i(\hat{x}_{i,j})}{\max_i(\hat{x}_{i,j}) - \min_i(\hat{x}_{i,j})} \quad (2)$$

176 where \bar{x} denotes the normalized value of the features vector. For clarity
 177 in the rest of the paper we will use x instead of \bar{x} .

178 **Initialization:**

179 Meta-parameters for the xDNN are initialized with the first observed data
 180 sample (image). The proposed algorithm works per class; therefore, all

181 the calculations are done for each class separately.

$$P \leftarrow 1; \quad \mu \leftarrow x_i; \tag{3}$$

182 where μ denotes the global mean of data samples of the given class. P
 183 is the total number of the identified prototypes from the observed data
 184 samples (images).

Each class C is initialized by the first data sample of that class:

$$\begin{aligned} C_1 &\leftarrow x_1; \quad p_1 \leftarrow x_1; \\ Support_1 &\leftarrow 1; \quad r_1 \leftarrow r^*; \quad \hat{I}_1 \leftarrow I_1 \end{aligned} \tag{4}$$

185 where, p_1 is the vector of features that describe the prototype \hat{I} of the C_1 ; \hat{I}
 186 is the identified prototype; $Support_1$ is the corresponding support (number
 187 of members) associated with this prototype; r_1 is the corresponding radius
 188 of the area of influence of C_1 .

189 In this paper, we use $r^* = \sqrt{2 - 2\cos(30^\circ)}$ same as (Angelov and Gu,
 190 2019); the rationale is that two vectors for which the angle between them
 191 is less than $\pi/6$ or 30° are pointing in close/similar directions d . That
 192 is, we consider that two feature vectors can be considered to be similar if
 193 the angle between them is smaller than 30 degrees. Note that r^* is data
 194 derived, not a problem- or user- specific parameter. In fact, it can be
 195 defined without *prior* knowledge of the specific problem or data through
 196 the following equation (5).

$$d(x_i, p_i) = \left\| \frac{x_i}{\|x_i\|} - \frac{p_i}{\|p_i\|} \right\|. \tag{5}$$

197 2. Density layer:

198 The density layer defines the mutual proximity of the images in the data
 199 space defined by the features from the previous layer. The data density,
 200 if use Euclidean form of distance, has a Cauchy form (15) (Angelov and
 201 Gu, 2019):

$$D(x_i) = \frac{1}{1 + \frac{\|x_i - \mu_N\|^2}{\|\sigma\|_N^2}}, \quad (6)$$

202 where D is the density, μ is the global mean, and σ is the variance. The
 203 reason it is Cauchy is not arbitrary (Angelov and Gu, 2019). It can be
 204 demonstrated theoretically that if Euclidean or Mahalanobis type of dis-
 205 tances in the feature space are considered, the data density reduces to
 206 Cauchy type as referred in equation (15). Density can also be updated
 207 online (Angelov, 2012):

$$D(x_i) = \frac{1}{1 + \|x_i - \mu_i\|^2 + \sum_i -\|\mu_i\|^2}. \quad (7)$$

208 where μ_i and the scalar product, \sum_i can be updated recursively as follows:

$$\mu_i = \frac{i-1}{i} \mu_{i-1} + \frac{1}{i} x_i, \quad (8)$$

$$\sum_i = \frac{i-1}{i} \sum_{i-1} + \frac{1}{i} \|x_i\|^2 \quad \sum_1 = \|x_1\|^2. \quad (9)$$

209 Data samples (images) that are closer to the global mean have higher
 210 density values. Therefore, the value of the data density indicates how
 211 strongly a particular data sample is influenced by other data samples in
 212 the data space due to their mutual proximity.

213 3. Typicality layer:

Typicality is an empirically derived form of probability distribution func-
 tion (pdf). *Typicality* τ is given by the equation (10). The value of τ even
 at the point $x = p_i$ is much less than 1; the integral of $\int_{-\infty}^{\infty} \tau dx = 1$
 (Angelov and Gu, 2019).

$$\tau(x_i) = \frac{\sum_{i=1}^c \text{Support}_i D(x_i)}{\sum_{i=1}^c \text{Support}_i \int_{-\infty}^{\infty} D(x_i) dx} \quad (10)$$

214 4. Prototypes layer:

215 The prototypes identification layer is the core of the proposed xDNN clas-
 216 sifier. This layer is responsible to provide the clearly explainable model.

217 The xDNN classifier is free from *prior* assumptions about the data dis-
 218 tribution type, as well as the random or deterministic nature of the data.
 219 **In contrast, it empirically extracts the distribution from the data sam-**
 220 **ples (images) bottom up (Angelov and Gu, 2019).** The prototypes are
 221 independent from each other. Therefore, one can change the structure
 222 by adding a new prototype without influencing the other already existing
 223 prototypes. In other words, the proposed xDNN is highly parallelizable
 224 and suitable for evolving form of application where new prototypes may
 225 be added (if the data pattern requires this). The proposed xDNN method
 226 is trained per class forming a set of prototypes per class. Therefore, all the
 227 calculations are done for each class separately. Prototypes are the local
 228 peaks of the data density (and *typicality*) identified in the previous layers/
 229 stages of the algorithm from the images of the corresponding class based
 230 on their feature vectors. The prototypes can be used to form linguistic
 231 logical *IF...THEN* rules of the following form:

232 R_c : IF $(I \sim \hat{I}_P)$ THEN (*class c*)

233 where \sim stands for similarity, it also can be seen as a fuzzy degree of
 234 membership; p is the identified prototype; P is the number of identified
 235 prototypes; c is the class $c = 1, 2, \dots, C$, I denotes an image.

236 One rule per prototype can be formed. All rules per class can be combined
 237 together using logical OR, also known as disjunction or S-norm:

238 R_c : IF $(I \sim \hat{I}_1)$ OR $(I \sim \hat{I}_2)$ OR ... OR $(I \sim \hat{I}_P)$ THEN (*class c*)

239 Figure 4 illustrates the area of influence of the identified prototypes. These
 240 areas around the identified prototypes are called *data clouds* (Angelov and
 241 Gu, 2019). Thus, each prototype defines a *data cloud*.

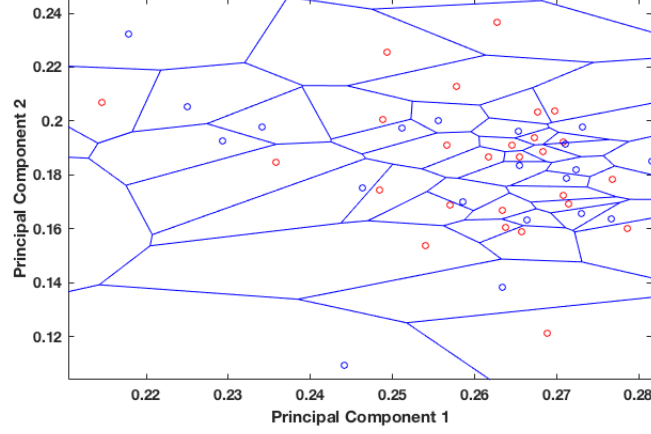


Figure 4: Identified prototypes – Voronoi Tesselation.

242 We call all data points associated with a prototype *data clouds*, because
 243 their shape is not regular (e.g., hyper-spherical, hyper-ellipsoidal, etc.)
 244 and the prototype is not necessarily the statistical and geometric mean ,
 245 but actual image (Angelov and Gu, 2019). The algorithm absorbs the new
 246 new data samples one by one by assigning them to the nearest (in the feature
 247 space) prototype:

$$j^* = \underset{j=1,2,\dots,P}{\operatorname{argmin}} (||x_i - p_j||^2) \quad (11)$$

In case, the following condition (Angelov and Gu, 2019) is met:

$$\begin{aligned} &IF (D(x_i) \geq \max_{j=1,2,\dots,P} D(p_j)) \\ &OR (D(x_i) \leq \min_{j=1,2,\dots,P} D(p_j)) \end{aligned} \quad (12)$$

THEN (add a new data cloud ($P \leftarrow P + 1$))

It means that x_i is out of the influence area of p_j . Therefore, the vector of features x_i becomes a new prototype of a new *data cloud* with meta-

parameters initialized by equation (13). Add a new *data cloud*:

$$\begin{aligned}
P &\leftarrow P + 1; & C_P &\leftarrow x_i; p_P \leftarrow I_i; & Support_P &\leftarrow 1; \\
r_P &\leftarrow r_o; \hat{I}_P &\leftarrow I_i;
\end{aligned} \tag{13}$$

248 Otherwise, *data cloud* parameters are updated online by equation (14). It
249 has to be stressed that all calculations per *data cloud* are performed on
250 the basis of data points associated with a certain *data cloud* only (i. e.
251 locally, not globally, on the basis of all data points).

$$\begin{aligned}
C_{j^*} &\leftarrow C_{j^*} + 1; \\
p_{j^*} &\leftarrow \frac{Support_{j^*}}{Support_{j^*} + 1} p_{j^*} + \frac{Support_{j^*}}{Support_{j^*} + 1} x_i; \\
Support_{j^*} &\leftarrow Support_{j^*} + 1; \\
r_{j^*}^2 &\leftarrow \frac{r_{j^*}^2 + (1 - \|p_{j^*}\|^2)}{2}.
\end{aligned} \tag{14}$$

252 The xDNN learning procedure can be summarized by the following algo-
253 rithm.

254 **xDNN: Learning Procedure**

- 255 1: Read the first feature vector sample x_i representing the image I_i of
256 the class c ;
- 257 2: Set $i \leftarrow 1; n \leftarrow 1; P_1 \leftarrow 1; p_1 \leftarrow x_i; \mu \leftarrow x_1; Support \leftarrow 1; r_1 \leftarrow$
258 $r_o; \hat{I}_1 \leftarrow I_1$;
- 259 3: **FOR** $i = 2, \dots$
- 260 4: Read x_i ;
- 261 5: Calculate $D(x_i)$ and $D(p_j)$ ($j = 1, 2, \dots, P$) according to equation
262 (9);
- 263 6: **IF** equation (12) holds
- 264 7: Create rule according to equation (13);
- 265 8: **ELSE**
- 266 9: Search for p_j according to equation (11);
- 267 10: Update rule according to equation (14);

268 11: **END**

269 12: **END**

270 5. *MegaClouds* layer:

271 In the *MegaClouds* layer the *clouds* formed by the prototypes in the pre-
272 vious layer are merged if the neighbouring prototypes have the same class
273 label. In other words, they are merged if they belong to the same class.
274 *MegaClouds* are used to facilitate the human interpretability. Figure 5
275 illustrates the formation of the *MegaClouds*.

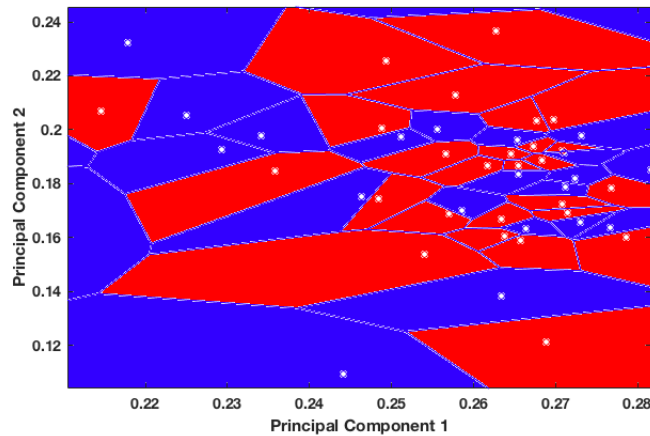


Figure 5: *MegaClouds* – Voronoi Tesselation.

276 Rules in the *MegaClouds* layer have the following format:

277 R_c : IF ($x \sim MC_1$) OR ($x \sim MC_2$) OR ... OR ($x \sim MC_{mc}$) THEN (*class*
278 *c*)

279 where *MC* are the *MegaClouds*, or the areas formed from the merging of
280 the *clouds*, and *mc* is the number of identified *MegaClouds*. Multimodal
281 *typicality*, τ , can also be used to illustrate the *MegaClouds* as illustrated
282 by Figure 6.

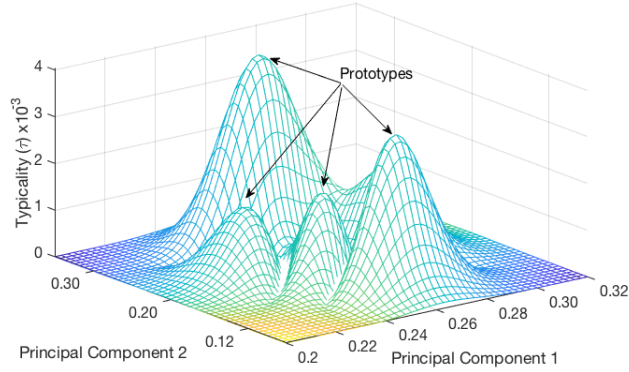


Figure 6: *Typicality* for the iRoads dataset.

283 3.2. Architecture and Validation of the proposed xDNN

284 Architecture for the validation process of the proposed xDNN method is
 illustrated by Figure 7.

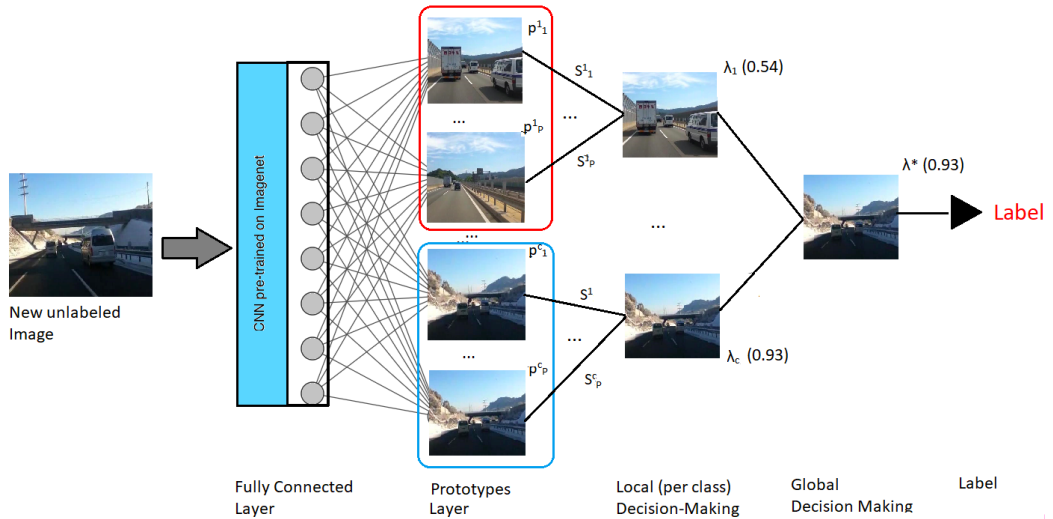


Figure 7: Architecture for the validation process of the proposed xDNN.

285

286 The validation process of xDNN is composed of the following layers:

- 287 1. Features layer;
 288 2. Similarity layer (density);
 289 3. Local decision-making.
 290 4. Global decision-making.

291 Which is detailed described as following:

292 **1. Features layer:**

293 Similarly to the features layer described in the training process.

294 **2. Prototypes layer:**

295 In this layer the degrees of similarity to the nearest prototypes (per class)
 296 are extracted for each unlabeled (new/validation) data sample/image I_i
 297 defined as follows:

$$S(x, p_i) = \frac{1}{1 + \frac{\|x - p_i\|^2}{\|\sigma\|_N^2}}, \quad (15)$$

298 where S denotes the similarity degree.

299 **3. Local (per class) decision-making layer:**

300 Local (per class) decision-making is calculated based on the ‘winner-takes-
 301 all’ principle and can be obtained by:

$$\lambda_c = \max_{j=1,2,\dots,P} (S_j), \quad (16)$$

302 **4. Global decision-making layer:** The global decision-making layer is
 303 in charge of forming the decision by assigning labels to the validation
 304 images based on the degree of similarity of the prototypes obtained by the
 305 prototype identification layer as illustrated by Figure 7 and determining
 306 the winning class.

$$\lambda_c^* = \max_{c=1,2,\dots,C} (\lambda_c), \quad (17)$$

307 In order to determine the overall degree of satisfaction, the maximum of
 308 the local, per class winners is applied.

309 The label is obtained by the following equation (18):

$$label = \underset{c=1,2,\dots,C}{\operatorname{argmax}} (\lambda_c^*), \quad (18)$$

310 4. Experimental Data

311 We validated our proposed approach, xDNN using several complex, well-
312 known image classification benchmark datasets (iRoads, Caltech-256, Caltech-
313 101) as well as we propose our own dataset for SARS-CoV-2 identification.

314 4.1. iRoads dataset

315 The iROADS dataset (Rezaei and Terauchi, 2013) was considered in the
316 analysis first. The dataset contains 4,656 image frames recorded from moving
317 vehicles on a diverse set of road scenes, recorded in day, night, under various
318 weather and lighting conditions, as described below:

- 319 • Daylight - 903 images
- 320 • Night - 1050 images
- 321 • Rainy day - 1049 images
- 322 • Rainy night - 431 images
- 323 • Snowy - 569 images
- 324 • Sun strokes - 307 images
- 325 • Tunnel - 347 images

326 4.2. Caltech-256

327 Caltech-256 has 30,607 images divided into 257 object categories (one of
328 which is the background) (Griffin et al., 2007).

329 *4.3. Caltech-101*

330 Caltech-101 is divided into 102 object categories (one of which is the back-
331 ground) (Fei-Fei et al., 2004).

332 *4.4. COVID-CT dataset*

333 COVID-CT dataset contains 275 computed tomography scans positive for
334 COVID-19 (Zhao et al., 2020).

335 *4.5. Performance Evaluation*

336 We used the following metrics for classification evaluation:

$$ACC(\%) = \frac{TP + TN}{TP + FP + TN + FN} \times 100, \quad (19)$$

337 Precision:

$$Precision(\%) = \frac{TP}{TP + FP} \times 100, \quad (20)$$

338 Recall:

$$Recall(\%) = \frac{TP}{TP + FN} \times 100, \quad (21)$$

339 *F1 Score:*

$$F1\ Score(\%) = 2 \times \frac{Precision \times Recall}{Precision + Recall} \times 100, \quad (22)$$

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

341 The area under the curve, AUC , is defined through the TP rate and FN
342 rate.

343 All the experiments were conducted with MATLAB 2018a using a personal
344 computer with a 1.8 GHz Intel Core i5 processor, 8-GB RAM, and MacOS
345 operating system. The classification experiments were executed using 10-fold
346 cross validation under the same ratio of training-to-testing (90% to 10%) sample
347 sets.

348 5. Results and Analysis

349 Computational simulations were performed to assess the accuracy of the
350 proposed explainable deep learning method, xDNN against other state-of-the-
351 art approaches.

352 5.1. *iRoads Dataset*

353 Table 1 shows that the proposed xDNN method provides the best result
354 in terms of classification accuracy as well as time/complexity and simplicity
355 of the model structure (number of parameters/prototypes). The number of
356 model parameters for xDNN (and DRB) is, strictly speaking, zero, because the
357 2 parameters (mean, μ and standard deviation, σ) per prototype (*data cloud*)
358 are derived from the data and are not algorithmic parameters or user-defined
359 parameters. For kNN method one can argue that the number of parameters
360 is the number of data samples, N . The proposed explainable DNN surpasses
361 in terms of accuracy the state-of-the-art VGG-16 algorithm which is a well-
362 established convolutional deep neural network. Moreover, the proposed xDNN
363 has at its top layer a set of a very small number of *MegaClouds* (27 or, on average,
364 4 *MegaClouds* per class) which makes it very easy to explain and visualize. For
365 comparison, our earlier version of deep rule-based models, called DRB (Angelov
366 and Gu, 2018) also produced a high accuracy and was trained a bit faster,
367 but ended up with 521 prototypes (on average 75 prototypes per class) (Soares
368 et al., 2019). With xDNN we do generate meaningful *IF...THEN* rules as well
369 as generate an analytical description of the *typicality* which is the empirically
370 derived pdf in a closed form which lends itself for further analysis and processing.

Table 1: Performance Comparasion: iRoads Dataset

Method	Accuracy	Time(s)	# Parameters
xDNN	99.59%	4.32	27
VGG-16 (He et al., 2016)	99.51 %	836.28	Not reported
DRB (Angelov and Gu, 2019)	99.02%	2.95	521
SVM (Suykens and Vandewalle, 1999)	94.17%	5.67	Not reported
KNN (Bishop, 2006)	93.49%	4.43	4656
Naive Bayes (Bishop, 2006)	88.35%	5.31	Not reported

371 *MegaClouds* generated by the proposed xDNN model can be visualized in
 372 terms of rules as illustrated by the Fig. 10.

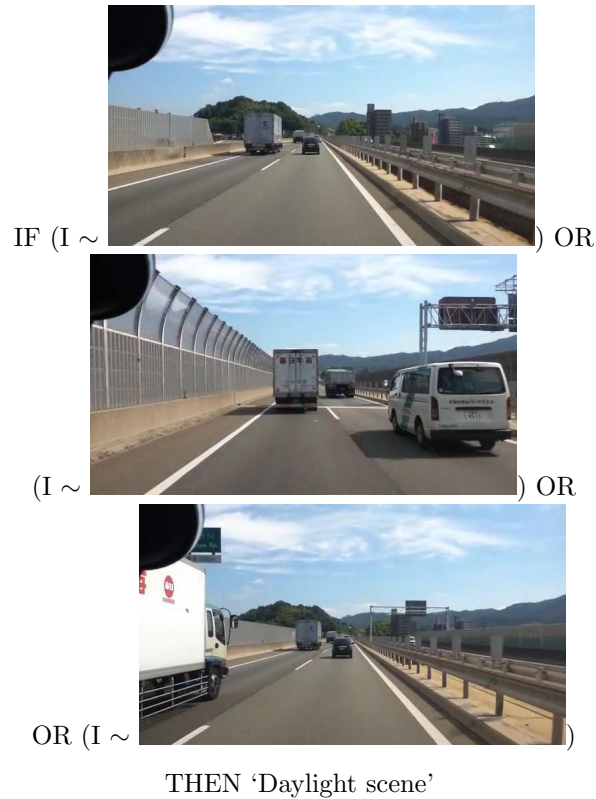


Figure 8: xDNN rule generated for the 'Daylight scene'.

373 Voronoi tessellation can also be used to visualize the resulting *MegaClouds*
 374 as illustrated by Figure 9.

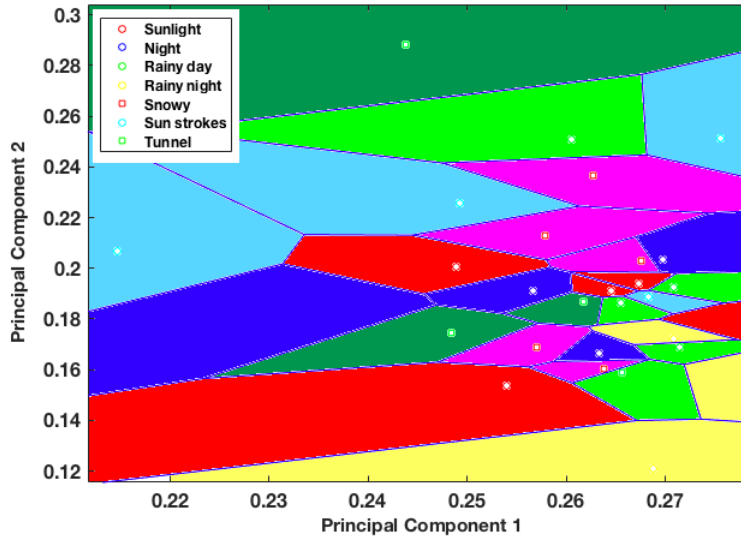


Figure 9: *MegaClouds* for the iRoads dataset.

375 5.2. Caltech-256 and Caltech-101 Dataset

376 Results for Caltech-256 are presented in Table 2.

Table 2: Performance Comparasion: Caltech-256 Dataset

Method	Accuracy
xDNN	75.41%
MSVM (Cao et al., 2019)	70.18%
VGG-16 (He et al., 2016)	73.2%
VGG-19 (He et al., 2016)	70.62 %
ResNet-101 (Simonyan and Zisserman, 2014)	75.14 %
GoogLeNet (Szegedy et al., 2015)	72.42 %
Softmax(7) (Zeiler and Fergus, 2014)	74.2%

377 Results presented in Table 2 demonstrate that the proposed xDNN approach
 378 can obtain highly accurate results compared to state-of-the-art approaches for
 379 this complex problem, it is important to highlight that we just compared the
 380 proposed approach with DNNs that do not use any trick for image augmentation.
 381 The proposed approach offers explainable models which can be visualized in
 382 terms of *IF...THEN* rules. xDNN produced on average 3 *MegaClouds* per
 383 class (a total of 721) which are clearly explainable. Rules have the following
 384 format:

$$\text{IF } (x \sim \text{CD}) \text{ OR } (x \sim \text{CD}) \text{ OR } (x \sim \text{CD}) \\ \text{THEN 'CD'}$$


385 We also tested the proposed xDNN approach on the Caltech-101 dataset.
 386 Results for the Caltech-101 dataset demonstrated on Table 3 showed that the
 387 proposed approach could surpass other state-of-the-art approaches in terms of
 388 accuracy.

Table 3: Performance Comparison: Caltech-101 Dataset

Method	Accuracy
xDNN	94.31%
SPP-net (He et al., 2015)	91.44%
ResNet-50 (He et al., 2016)	90.39%
CNN S TUNE-CLS (Chatfield et al., 2014)	88.35%
(Zeiler and Fergus, 2014)	86.5%
VGG-16 (He et al., 2016)	90.32%
KNN (Bishop, 2006)	85.65%
DT (Quinlan, 1986)	54.42%

389 We compared the proposed xDNN approach with the best published single-
 390 label classifiers methods and achieved better result. There are couple of alter-

391 native methods that report higher results on Caltech problems, but they use
 392 additional information such as the context (Leng et al., 2019) or multiple labels
 393 (Qian et al., 2019) processes in order to enhance the classification performance,
 394 include extra features (labels and descriptions) and this makes the underlying
 395 problem different even if the name is still the same (Caltech-101 or Caltech-
 396 256). We believe that the comparison has to be in the same playing field using
 397 the same amount of information and therefore, we do not report these meth-
 398 ods. Apart from them, to the best of our knowledge, there is no better result
 399 achieved on Caltech data sets.

400 5.3. COVID CT-scan dataset

401 In this section we report the results obtained by the proposed xDNN clas-
 402 sification approach when applied to the COVID CT-scan dataset (Zhao et al.,
 403 2020). Results presented in Table 4 compare the proposed algorithm with other
 404 state-of-the-art approaches, including traditional "black-box" Deep Neural Net-
 405 work, Support Vector Machines, etc.

Table 4: Performance Comparison: COVID CT-scan Dataset

Metric Method	Accuracy	Precision	Recall	F1 Score	AUC
xDNN	88.6%	89.7%	88.6%	89.2%	88.6%
Baseline (Zhao et al., 2020)	84.7%	97.0%	76.2%	85.3%	82.4%
SVM (Suykens and Vandewalle, 1999)	80.5%	84.4%	83.5%	84%	79.7%
KNN (Bishop, 2006)	83.9%	90.4%	82.4%	86.2%	84.3%
AdaBoost (Hastie et al., 2009)	83.9%	87.7%	83.5%	85.5%	84%
Naive Bayes (Bishop, 2006)	70.5%	77%	73.6%	75.3%	69.6%

406 The proposed xDNN classifier provided better results in terms of accuracy,
 407 recall, $F1$ score, and AUC. Moreover, the proposed approach also provided
 408 highly interpretable results that may be helpful for specialists (in this case, med-

409 ical doctors). The proposed classifier identified 30 prototypes for non-COVID
 410 and 33 prototypes for COVID patients. Rules generated by the identified pro-
 411 totypes for COVID and non-COVID patients are illustrated by Figures 10 and
 412 11 respectively. The baseline approach Zhao et al. (2020) is a Deep Neural
 413 Network approach which is ‘black box’ (offers no interpretability).

414 Using the proposed method we extracted from the data linguistic *IF...THEN*
 415 rules which involve actual images of both cases (COVID-19 and non-COVID)
 416 as illustrated in Figures 10 and 11. Such transparent rules can be used in the
 417 decision-making process for early diagnostics for COVID-19 infection. Rapid
 418 detection with high sensitivity of viral infection may allow better control of the
 419 viral spread. Early diagnosis of COVID-19 is crucial for the disease treatment
 420 and control.

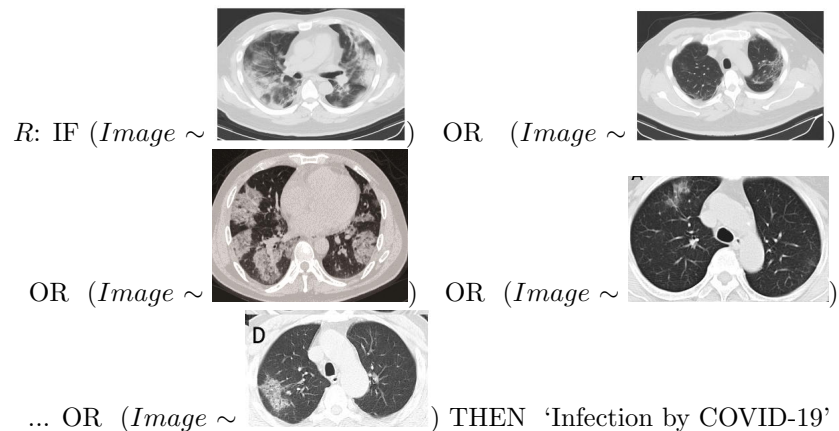


Figure 10: Final rule given by the proposed xDNN classifier for the COVID-19 identification. Differently from ‘black box’ approaches as deep neural networks, the proposed approach provides highly interpretable rules which can be used by human experts for the early evaluation of patients suspected of SARS-Cov-2 infection.

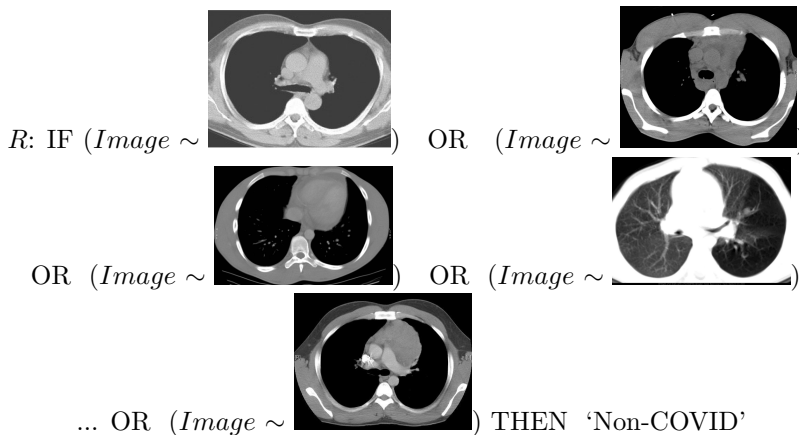


Figure 11: Non-Covid final rule given by the proposed eXplainable Deep Learning classifier.

421 Figure 12 illustrates the evolving nature of the proposed approach. xDNN
422 is able to continuously learn as new data is presented to it. Therefore, no full
423 re-training is required due to its life-long learning architecture. On the contrary,
424 the baseline approach Zhao et al. (2020) is based on a Deep Neural Network
425 that requires full re-training for any new data sample, which can be very costly
426 in terms of time, computational complexity and requirements for hardware and
427 computer experts. xDNN continuously learns as new training data arrives to the
428 system. It can be observed that with 478 training data samples the proposed
429 approach could obtain better results in terms of accuracy (84.56%) than the
430 baseline approach (84.0%) with 537 training data samples Zhao et al. (2020).
431 The baseline approach is a Deep Neural Network that needs a large number of
432 training data to obtain a high performance in terms of classification accuracy
433 and once trained can not be further improved unless fully re-trained. In contrast,
434 the proposed approach can obtain higher performance using less training data
435 due to its prototype-based nature.

436 Experiments have demonstrated that the proposed xDNN approach is able
437 to produce highly accurate results surpassing state-of-the-art methods for differ-
438 ent challenging datasets. Moreover, xDNN presents highly interpretable results
439 that can be presented in the form of *IF...THEN* logical rules, Voronoi tessella-

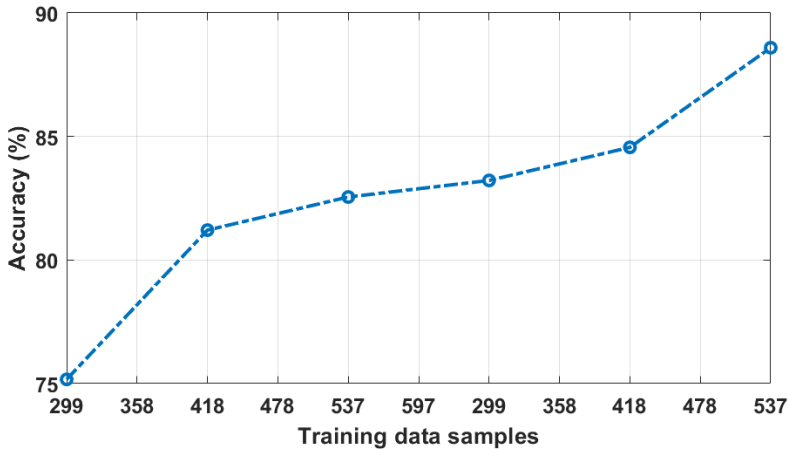


Figure 12: The figure illustrates the evolving nature of the proposed xDNN approach

440 tions, and/or *typicality* (empirically derived form of pdf) in a closed analytical
 441 form allowing further analysis. Because of its recursive, non-iterative and non-
 442 parametric form it allows computationally very efficient implementations to be
 443 realized.

444 6. Conclusion

445 In this paper we propose a new method, explainable deep neural network
 446 (xDNN), that is directly addressing the bottlenecks of the traditional deep learn-
 447 ing approaches and offers an explainable internal architecture that can outper-
 448 form the existing methods. The proposed xDNN approach requires very little
 449 computational resources (no need for GPUs) and short training times (in the
 450 order of seconds). The proposed approach, xDNN is prototype-based. Pro-
 451 totypes are actual training data samples (images), which have local peaks of
 452 the empirical data distribution called *typicality* as well as of the data density.
 453 This generative model is identified in a closed form and equates to the pdf but
 454 is derived automatically and entirely from the training data with no user- or
 455 problem-specific thresholds, parameters or intervention. The proposed xDNN
 456 offers a new deep learning architecture that combines reasoning and learning in

457 a synergy. It is non-iterative and non-parametric, which explains its efficiency
458 in terms of time and computational resources. From the user perspective, the
459 proposed approach is clearly understandable to human users. Results for some
460 well-known benchmark data sets such as iRoads, Caltech-256, Caltech-101, and
461 COVID CT-scan show that xDNN outperforms the other methods including
462 state-of-the-art deep learning approaches in terms of accuracy, time to train
463 and offers an explainable classifier. Future research will concentrate on the
464 development of a tree-based architecture, synthetic data generation, and local
465 optimization in order to improve the proposed deep explainable approach.

466 **References**

- 467 P. Angelov. *Autonomous learning systems: from data streams to knowledge in*
468 *real-time*. John Wiley & Sons, 2012.
- 469 P. P. Angelov and X. Gu. Deep rule-based classifier with human-level perfor-
470 mance and characteristics. *Information Sciences*, 463:196–213, 2018.
- 471 P. P. Angelov and X. Gu. *Empirical approach to machine learning*. Springer,
472 2019.
- 473 P. P. Angelov, X. Gu, and J. C. Príncipe. A generalized methodology for data
474 analysis. *IEEE transactions on cybernetics*, 48(10):2981–2993, 2017.
- 475 M. Biehl, B. Hammer, and T. Villmann. Distance measures for prototype based
476 classification. In *International Workshop on Brain-Inspired Computing*, pages
477 100–116. Springer, 2013.
- 478 M. Biehl, B. Hammer, and T. Villmann. Prototype-based models in machine
479 learning. *Wiley Interdisciplinary Reviews: Cognitive Science*, 7(2):92–111,
480 2016.
- 481 C. M. Bishop. *Pattern recognition and machine learning*. springer, 2006.

- 482 J. Cao, M. Wang, Y. Li, and Q. Zhang. Improved support vector machine clas-
483 sification algorithm based on adaptive feature weight updating in the hadoop
484 cluster environment. *PloS one*, 14(4), 2019.
- 485 K. Chatfield, K. Simonyan, A. Vedaldi, and A. Zisserman. Return of the
486 devil in the details: Delving deep into convolutional nets. *arXiv preprint*
487 *arXiv:1405.3531*, 2014.
- 488 F. Doshi-Velez and B. Kim. Towards a rigorous science of interpretable machine
489 learning. *arXiv preprint arXiv:1702.08608*, 2017.
- 490 L. Fei-Fei, R. Fergus, and P. Perona. Learning generative visual models from few
491 training examples: An incremental bayesian approach tested on 101 object
492 categories. In *2004 conference on computer vision and pattern recognition*
493 *workshop*, pages 178–178. IEEE, 2004.
- 494 I. Goodfellow, Y. Bengio, and A. Courville. *Deep learning*. MIT press, 2016.
- 495 G. Griffin, A. Holub, and P. Perona. Caltech-256 object category dataset. 2007.
- 496 T. Hastie, S. Rosset, J. Zhu, and H. Zou. Multi-class adaboost. *Statistics and*
497 *its Interface*, 2(3):349–360, 2009.
- 498 K. He, X. Zhang, S. Ren, and J. Sun. Spatial pyramid pooling in deep convolu-
499 tional networks for visual recognition. *IEEE transactions on pattern analysis*
500 *and machine intelligence*, 37(9):1904–1916, 2015.
- 501 K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recogni-
502 tion. In *Proceedings of the IEEE conference on computer vision and pattern*
503 *recognition*, pages 770–778, 2016.
- 504 J. Hu, J. Lu, and Y.-P. Tan. Deep transfer metric learning. In *Proceedings of the*
505 *IEEE conference on computer vision and pattern recognition*, pages 325–333,
506 2015.

- 507 A. Krizhevsky, I. Sutskever, and G. E. Hinton. Imagenet classification with deep
508 convolutional neural networks. In *Advances in neural information processing*
509 *systems*, pages 1097–1105, 2012.
- 510 Y. LeCun, Y. Bengio, and G. Hinton. Deep learning. *nature*, 521(7553):436–444,
511 2015.
- 512 J. Leng, Y. Liu, and S. Chen. Context-aware attention network for image
513 recognition. *Neural Computing and Applications*, 31(12):9295–9305, 2019.
- 514 O. Li, H. Liu, C. Chen, and C. Rudin. Deep learning for case-based reasoning
515 through prototypes: A neural network that explains its predictions. In *Thirty-*
516 *Second AAAI Conference on Artificial Intelligence*, 2018.
- 517 C. Liu, G. Bellec, B. Vogginger, D. Kappel, J. Partzsch, F. Neumärker,
518 S. Höppner, W. Maass, S. B. Furber, R. Legenstein, et al. Memory-efficient
519 deep learning on a spinnaker 2 prototype. *Frontiers in neuroscience*, 12:840,
520 2018.
- 521 K. Mizuno, Y. Terachi, K. Takagi, S. Izumi, H. Kawaguchi, and M. Yoshimoto.
522 Architectural study of hog feature extraction processor for real-time object
523 detection. In *2012 IEEE Workshop on Signal Processing Systems*, pages 197–
524 202. IEEE, 2012.
- 525 D. Nebel, M. Kaden, A. Villmann, and T. Villmann. Types of (dis-) similarities
526 and adaptive mixtures thereof for improved classification learning. *Neuro-*
527 *computing*, 268:42–54, 2017.
- 528 O. K. Oyedotun and A. Khashman. Prototype-incorporated emotional neural
529 network. *IEEE transactions on neural networks and learning systems*, 29(8):
530 3560–3572, 2017.
- 531 G. Qian, L. Zhang, and Y. Wang. Single-label and multi-label concepter classi-
532 fiers in pre-trained neural networks. *Neural Computing and Applications*, 31
533 (10):6179–6188, 2019.

- 534 J. R. Quinlan. Induction of decision trees. *Machine learning*, 1(1):81–106, 1986.
- 535 S. Ren, K. He, R. Girshick, X. Zhang, and J. Sun. Object detection networks
536 on convolutional feature maps. *IEEE transactions on pattern analysis and*
537 *machine intelligence*, 39(7):1476–1481, 2016.
- 538 M. Rezaei and M. Terauchi. Vehicle detection based on multi-feature clues
539 and Dempster-Shafer fusion theory. In *Pacific-Rim Symposium on Image and*
540 *Video Technology*, pages 60–72. Springer, 2013.
- 541 C. Rudin. Stop explaining black box machine learning models for high stakes
542 decisions and use interpretable models instead. *Nature Machine Intelligence*,
543 1(5):206–215, 2019.
- 544 S. Saralajew, L. Holdijk, M. Rees, and T. Villmann. Prototype-based neu-
545 ral network layers: incorporating vector quantization. *arXiv preprint*
546 *arXiv:1812.01214*, 2018.
- 547 J. Schmidhuber. Deep learning in neural networks: An overview. *Neural net-*
548 *works*, 61:85–117, 2015.
- 549 T. J. Sejnowski. *The deep learning revolution*. MIT Press, 2018.
- 550 K. Simonyan and A. Zisserman. Very deep convolutional networks for large-scale
551 image recognition. *arXiv preprint arXiv:1409.1556*, 2014.
- 552 E. Soares and P. Angelov. Novelty detection and learning from extremely weak
553 supervision. *arXiv preprint arXiv:1911.00616*, 2019.
- 554 E. Soares, P. Angelov, B. Costa, and M. Castro. Actively semi-supervised deep
555 rule-based classifier applied to adverse driving scenarios. In *2019 International*
556 *Joint Conference on Neural Networks (IJCNN)*, pages 1–8. IEEE, 2019.
- 557 B. Solmaz, S. M. Assari, and M. Shah. Classifying web videos using a global
558 video descriptor. *Machine vision and applications*, 24(7):1473–1485, 2013.

- 559 J. A. Suykens and J. Vandewalle. Least squares support vector machine classi-
560 fiers. *Neural processing letters*, 9(3):293–300, 1999.
- 561 C. Szegedy, W. Liu, Y. Jia, P. Sermanet, S. Reed, D. Anguelov, D. Erhan,
562 V. Vanhoucke, and A. Rabinovich. Going deeper with convolutions. In *Pro-
563 ceedings of the IEEE conference on computer vision and pattern recognition*,
564 pages 1–9, 2015.
- 565 C. Szegedy, S. Ioffe, V. Vanhoucke, and A. A. Alemi. Inception-v4, inception-
566 resnet and the impact of residual connections on learning. In *Thirty-First
567 AAAI Conference on Artificial Intelligence*, 2017.
- 568 W. Xiong, L. Wu, F. Alleva, J. Droppo, X. Huang, and A. Stolcke. The microsoft
569 2017 conversational speech recognition system. In *2018 IEEE international
570 conference on acoustics, speech and signal processing (ICASSP)*, pages 5934–
571 5938. IEEE, 2018.
- 572 M. D. Zeiler and R. Fergus. Visualizing and understanding convolutional net-
573 works. In *European conference on computer vision*, pages 818–833. Springer,
574 2014.
- 575 J. Zhao, Y. Zhang, X. He, and P. Xie. Covid-ct-dataset: a ct scan dataset about
576 covid-19. *arXiv preprint arXiv:2003.13865*, 2020.
- 577 F. Zhuang, X. Cheng, P. Luo, S. J. Pan, and Q. He. Supervised representa-
578 tion learning: Transfer learning with deep autoencoders. In *Twenty-Fourth
579 International Joint Conference on Artificial Intelligence*, 2015.