

# Regular Polytope Networks

Federico Pernici, Matteo Bruni, Claudio Bacchi and Alberto Del Bimbo, *Member, IEEE*,

**Abstract**—Neural networks are widely used as a model for classification in a large variety of tasks. Typically, a learnable transformation (i.e. the classifier) is placed at the end of such models returning a value for each class used for classification. This transformation plays an important role in determining how the generated features change during the learning process. In this work, we argue that this transformation not only can be fixed (i.e. set as non-trainable) with no loss of accuracy and with a reduction in memory usage, but it can also be used to learn stationary and maximally separated embeddings. We show that the stationarity of the embedding and its maximal separated representation can be theoretically justified by setting the weights of the fixed classifier to values taken from the coordinate vertices of the three regular polytopes available in  $\mathbb{R}^d$ , namely: the  $d$ -Simplex, the  $d$ -Cube and the  $d$ -Orthoplex. These regular polytopes have the maximal amount of symmetry that can be exploited to generate stationary features angularly centered around their corresponding fixed weights. Our approach improves and broadens the concept of a fixed classifier, recently proposed in [1], to a larger class of fixed classifier models. Experimental results confirm the theoretical analysis, the generalization capability, the faster convergence and the improved performance of the proposed method. Code will be publicly available.

**Index Terms**—Deep Neural Networks, Fixed classifiers, Internal feature representation.

## I. INTRODUCTION

DEEP Convolutional Neural Networks (DCNNs) have achieved state-of-the-art performance on a variety of tasks [2], [3] and have revolutionized Computer Vision in both classification [4], [5] and representation [6], [7]. In DCNNs, both representation and classification are typically jointly learned in a single network. The classification layer placed at the end of such models transforms the  $d$ -dimension of the network internal feature representation to the  $K$ -dimension of the output class probabilities. Despite the large number of trainable parameters that this layer adds to the model (i.e.  $d \times K$ ), it has been verified that its removal only causes a slight increase in error [8]. Moreover, the most recent architectures tend to avoid the use of fully connected layers [9] [10] [11]. It is also well known that DCNNs can be trained to perform metric learning without the explicit use of a classification layer [12] [13] [14]. In particular, it has been shown that excluding from learning the parameters of the classification layer causes little or no decline in performance while allowing a reduction in the number of trainable parameters [1]. Fixed classifiers also have an important role in the theoretical convergence analysis of training models with batch-norm [15]. Very recently it has been shown that DCNNs with a fixed classifier and batch-norm in each layer establish a principle of equivalence between different learning rate schedules [16].

MICC, Media Integration and Communication Center, University of Florence, Dipartimento di Ingegneria dell'Informazione Firenze, Italy.

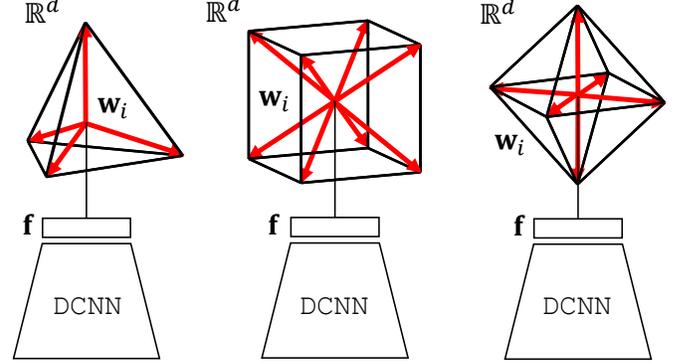


Fig. 1. Regular Polytope Networks (RePoNet). The fixed classifiers derived from the three regular polytopes available in  $\mathbb{R}^d$  with  $d \geq 5$  are shown. From left: the  $d$ -Simplex, the  $d$ -Cube and the  $d$ -Orthoplex fixed classifier. The trainable parameters  $\mathbf{w}_i$  of the classifier are replaced with fixed values taken from the coordinate vertices of a regular polytope (shown in red).

All these works seem to suggest that the final fully connected layer used for classification is somewhat redundant and does not have a primary role in learning and generalization. In this paper we show that a special set of fixed classification layers *has a key role in modeling the internal feature representation* of DCNNs, while ensuring little or no loss in classification accuracy and a significant reduction in memory usage.

In DCNNs the internal feature representation for an input sample is the feature vector  $\mathbf{f}$  generated by the penultimate layer, while the last layer (i.e. the classifier) outputs score values according to the inner product as:

$$z_i = \mathbf{w}_i^\top \cdot \mathbf{f} \quad (1)$$

for each class  $i$ , where  $\mathbf{w}_i$  is the weight vector of the classifier for the class  $i$ . To evaluate the loss, the scores are further normalized into probabilities via the softmax function [17]. Since the values of  $z_i$  can be also expressed as  $z_i = \mathbf{w}_i^\top \cdot \mathbf{f} = \|\mathbf{w}_i\| \|\mathbf{f}\| \cos(\theta)$ , where  $\theta$  is the angle between  $\mathbf{w}_i$  and  $\mathbf{f}$ , the score for the correct label with respect to the other labels is obtained by optimizing the length of the vectors  $\|\mathbf{w}_i\|$ ,  $\|\mathbf{f}\|$  and the angle  $\theta$  they are forming. This simple formulation of the final classifier provides the intuitive explanation of how feature vector directions and weight vector directions align simultaneously with each other at training time so that their average angle is made as small as possible. If the parameters  $\mathbf{w}_i$  of the classifier in Eq. 1 are fixed (i.e. set as non-trainable), *only the feature vector directions* can align toward the classifier weight vector directions and not the opposite. Therefore, weights can be regarded as fixed angular references to which features align.

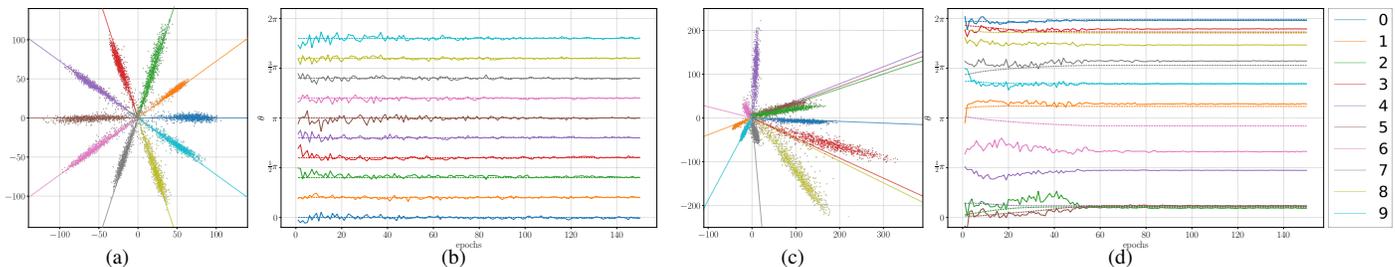


Fig. 2. Feature learning on the MNIST dataset in a 2D embedding space. Fig. (a) and Fig. (c) show the 2D features learned by RePoNet and by a standard trainable classifier respectively. Fig. (b) and Fig. (d) show the training evolution of the classifier weights (dashed) and their corresponding class feature means (solid) respectively. Both are expressed according to their angles. Although the two methods achieve the same classification accuracy, features in the proposed method are both stationary and maximally separated.

According to this, we obtain a precise result on the spatio-temporal statistical properties of the generated features during the learning phase. Supported by the empirical evidence in [1] we show that not only the final classifier of a DCNN can be set as non-trainable with no loss of accuracy and with a significant reduction in memory usage, but that an appropriate set of values assigned to its weights allows learning a maximally separated and strictly stationary embedding while training. That is, the features generated by the Stochastic Gradient Descent (SGD) optimization have constant mean and are angularly centered around their corresponding fixed class weights. Constant known mean implies that features cannot have non-constant trends while learning. Maximally separated features and their stationarity are obtained by setting the classifier weights according to values following a highly symmetrical configuration in the embedding space.

DCNN models with trainable classifiers are typically convergent and therefore, after a sufficient learning time has elapsed, some form of stationarity in the learned features can still be achieved. However, until that time, it is not possible to know where the features will be projected by the learned model in the embedding space. An advantage of the approach proposed in this paper is that it allows to define (and therefore to know in advance) where the features will be projected before starting the learning process.

Our result can be understood by looking at the basic functionality of the final classifier in a DCNN. The main role of a trainable classifier is to dynamically adjust the decision boundaries to learn class feature representations. When the classifier is set as non-trainable this dynamic adjustment capability is no longer available and it is automatically demanded to all the previous layers. Specifically, the work [1] reports empirical evidence that the expressive power of DCNN models is large enough to account for the missing dynamic adjustment capability of the classifier. We provide more systematic empirical evidence confirming and broadening the general validity of DCNNs with fixed classifiers (Sec. V-A).

We show that our approach can be theoretically justified and easily implemented by setting the classifier weights to values taken from the coordinate vertices of a regular polytope in the embedding space. Regular polytopes are the generalization in any number of dimensions of regular polygons and regular polyhedra (i.e. Platonic Solids). Although there are infinite

regular polygons in  $\mathbb{R}^2$  and 5 regular polyhedra in  $\mathbb{R}^3$ , there are only three regular polytopes in  $\mathbb{R}^d$  with  $d \geq 5$ , namely the  $d$ -Simplex, the  $d$ -Cube and the  $d$ -Orthoplex. Having different symmetry, geometry and topology, each regular polytope will reflect its properties into the classifier and the embedding space which it defines. Fig. 1 illustrates the three basic architectures defined by the proposed approach termed Regular Polytope Networks (RePoNet). Fig. 2 provides a first glance at our main result in a 2D embedding space. Specifically, the main evidence from Fig. 2(a) and 2(b) is that the features learned by RePoNet remain aligned with their corresponding fixed weights and maximally exploit the available representation space directly from the beginning of the training phase.

We apply our method to multiple vision datasets showing that it is possible to generate stationary and maximally separated features without affecting the generalization performance of DCNN models and with a significant reduction in GPU memory usage at training time. A preliminary exploration of this work was presented in [18], [19].

## II. RELATED WORK

**Fixed Classifier.** Empirical evidence shows that convolutional neural networks with a fixed classification layer (i.e. not subject to learning) initialized by random numbers does not worsen the performance on the CIFAR-10 dataset [20]. A recent paper [1] explores in more detail the idea of excluding from learning the parameters  $w_i$  in Eq.1. The work shows that a fixed classifier causes little or no reduction in classification performance for common datasets while allowing a significant reduction in trainable parameters, especially when the number of classes is large. Setting the last layer as not trainable also reduces the computational complexity for training as well as the communication cost in distributed learning. The paper in question sets the classifier with the coordinate vertices of orthogonal vectors taken from the columns of the Hadamard<sup>1</sup> matrix and does not investigate on the internal feature representation. A major limitation of this method is that, when the number of classes is higher than the dimension of the feature space, it is not possible to have mutually orthogonal columns. As a consequence some of the classes are constrained to lie in a common subspace which causes a reduction in classification

<sup>1</sup>The Hadamard matrix is a square matrix whose entries are either +1 or -1 and whose rows are mutually orthogonal.

performance. In our solution, we improve and generalize this work by finding a novel set of unique directions overcoming the limitations of the Hadamard matrix.

The work [21] trains a neural network according to the triplet loss with a set of fixed vertices on a hyper-sphere (i.e. a sphere lattice). The work aims at learning a function that maps real-valued vectors to a uniform distribution over a  $d$ -dimensional sphere.

As shown in [16], fixed classifiers are also related to BatchNorm [15] and learning rate schedules. BatchNorm parametrizes the weights of a layer to “normalize” its activations (i.e. the features), but typically it is not applied to normalize the outputs of the classifier (i.e. logits). They show that, with BatchNorm layers and a fixed classifier layer, training with  $L2$  regularization is equivalent to training with an exponentially increasing learning rate.

**Softmax Angular Optimization.** As originally described in [22], under softmax loss<sup>2</sup> the class prediction is *largely* determined by the angular similarity since softmax loss can be factorized, as shown in Eq. 1, into an amplitude component and an angular component. Several papers have followed this intuition and proposed to train DCNNs by direct angle optimization [23], [24], [25], [26]. The angle encodes the required discriminative information for class recognition. The wider the angles the better the classes are separated from each other and, accordingly, their representation is more discriminative. The common idea of these works is to constrain the features and/or the classifier to be unit normalized. The works [27], [28] and [26] normalize both features and the classifier weights thus obtaining an exact optimization of the angle in Eq. 1. With weight normalization only, label prediction is largely determined by the angular similarity [22] [24]. This is not only because Eq. 1 can be factorized into amplitude and angular component, but also because decision boundaries between adjacent classes are determined by their angular bisectors.

Differently from weight normalization, feature normalization cannot directly perform angle optimization but encourages intra-class compactness of learned features [23]. Specifically, [23] also proposes adding a multiplicative scale parameter after feature normalization based on the property that increasing the norm of samples can decrease the softmax loss [26], [29]. Although with a different goal than learning discriminative features, the work [1], in addition to fixing the classifier, normalizes both the weights and the features and applies the multiplicative scale parameter.

In agreement with [30], [1], [23] and [26] we found that applying the feature normalization and the multiplicative scale parameter makes optimization hard with general datasets, having a significant dependence on image quality. According to this, we follow the work [24] that normalizes the classifier weights. Normalizing classifier weights typically also includes setting the classifier biases to zero. As discussed in [26] and in [29] this encourages well-separated features to have bigger magnitudes. This avoids features collapsing into the origin,

making angles between the weights and features a reliable metric for classification.

As conjectured in [26], if all classes are well-separated, weight normalization will roughly correspond to computing the mean of features in each class. The maximal and fixed separation proposed in this paper further strengthens the conjecture, producing features more centered around their fixed weights as the training process progresses.

Another close related work to ours is [31] in which separability of learned features is improved by injecting a single dynamic virtual negative class into the original softmax. A virtual class is a class that is active in the classifier but has no data available from which to learn. Injecting the virtual class enlarges the inter-class margin and compresses intra-class distribution by strengthening the decision boundary constraint. In our case, we can profitably exploit virtual classes when the number of classes of the fixed classifier does not match the number of vertices of a regular polytope.

While all the above works impose large angular distances between the classes, they provide solutions to enforce such constraint in a local manner without considering global inter-class separability and intra-class compactness. For this purpose, very recently the works [32], [33] and [34] add a regularization loss to specifically force the classifier weights to be far from each other in a global manner. These works draw inspiration from a well-known problem in physics – the Thomson problem [35], where given  $K$  charges confined to the surface of a sphere, one seeks to find an arrangement of the charges which minimizes the total electrostatic energy. Electrostatic force repels charges each other inversely proportional to their mutual distance. In [32], [33] and [34] global equiangular features are obtained by adding to the standard categorical cross-entropy loss a further loss inspired by the Thomson problem.

In our research, we follow a similar principle for global separability. We consider that minimal energies are often concomitant with special geometric configurations of charges that recall the geometry of Platonic Solids in high dimensional spaces [36]. We have reported a few preliminary and qualitative results of using Regular Polytope Networks for compact feature learning in [18], where we have demonstrated that the angular parameter of the margin loss of [7] can be analytically determined to maximize feature compactness.

### III. MAIN CONTRIBUTIONS:

Our technical contributions can be summarized as follows:

- 1) We generalize the concept of fixed classifiers and show they can generate stationary and maximally separated features at training time with no loss of performance and in many cases with slightly improved performance.
- 2) We performed extensive evaluations across a range of datasets and modern CNN architectures reaching state-of-the-art performance. We observed faster speed of convergence and a significant reduction in model parameters.
- 3) We further provide a formal characterization of the class decision boundaries according to the dual relationship

<sup>2</sup>The combination of cross-entropy loss and the softmax function at the last fully connected layer.

between regular polytopes and statistically verify the validity of our method on random permutations of the labels.

#### IV. REGULAR POLYTOPES AND MAXIMALLY SEPARATED STATIONARY EMBEDDINGS

We are basically concerned with the following question: *How should the non-trainable weights of the classifier be distributed in the embedding space such that they generate stationary and maximally separated features?*

Let  $\mathbb{X} = \{(x_i, y_i)\}_{i=1}^N$  be the training set containing  $N$  samples, where  $x_i$  is the raw input to the DCNN and  $y_i \in \{1, 2, \dots, K\}$  is the label of the class that supervises the output of the DCNN. Then, the cross entropy loss can be written as:

$$\mathcal{L} = -\frac{1}{N} \sum_{i=1}^N \log \left( \frac{\exp(\mathbf{w}_{y_i}^\top \mathbf{f}_i + \mathbf{b}_{y_i})}{\sum_{j=1}^K \exp(\mathbf{w}_j^\top \mathbf{f}_i + \mathbf{b}_j)} \right), \quad (2)$$

where  $\mathbf{W} = \{\mathbf{w}_j\}_{j=1}^K$  are the classifier weight vectors for the  $K$  classes. Following the discussion in [24] we normalize the weights and zero the biases ( $\hat{\mathbf{w}}_j = \frac{\mathbf{w}_j}{\|\mathbf{w}_j\|}$ ,  $\mathbf{b}_j = 0$ ) to directly optimize angles, enabling the network to learn angularly distributed features. Angles therefore encode the required discriminative information for class recognition and the wider they are, the better the classes are represented. As a consequence, the representation in this case is maximally separated when features are distributed at *equal angles* maximizing the available space.

If we further consider the feature vector parametrized by its unit vector as  $\mathbf{f}_i = \kappa_i \hat{\mathbf{f}}_i$  where  $\kappa_i = \|\mathbf{f}_i\|$  and  $\hat{\mathbf{f}}_i = \frac{\mathbf{f}_i}{\|\mathbf{f}_i\|}$ , then Eq.2 can be rewritten as:

$$\mathcal{L} = -\frac{1}{N} \sum_{i=1}^N \log \left( \frac{\exp(\kappa_i \hat{\mathbf{w}}_{y_i}^\top \hat{\mathbf{f}}_i)}{\sum_{j=1}^K \exp(\kappa_i \hat{\mathbf{w}}_j^\top \hat{\mathbf{f}}_i)} \right) \quad (3)$$

The equation above can be interpreted as if  $N$  realizations from a set of  $K$  von Mises-Fisher distributions with different concentration parameters  $\kappa_i$  are passed through the softmax function. The probability density function of the von Mises-Fisher distribution for the random  $d$ -dimensional unit vector  $\hat{\mathbf{f}}$  is given by:  $P(\hat{\mathbf{f}}; \hat{\mathbf{w}}, \kappa) \propto \exp(\kappa \hat{\mathbf{w}}^\top \hat{\mathbf{f}})$  where  $\kappa \geq 0$ . Under this parameterization  $\hat{\mathbf{w}}$  is the mean direction on the hypersphere and  $\kappa$  is the concentration parameter. The greater the value of  $\kappa$  the higher the concentration of the distribution around the mean direction  $\hat{\mathbf{w}}$ . The distribution is unimodal for  $\kappa > 0$  and is uniform on the sphere for  $\kappa = 0$ .

As with this formulation each weight vector is the mean direction of its associated features on the hypersphere, equiangular features maximizing the available space can be obtained by arranging accordingly their corresponding weight vectors around the origin. This problem is equivalent to distributing points uniformly on the sphere and is a well-known geometric problem, called Tammes problem [37] which is a generalization of the physic problem firstly addressed by Thomson [35]. In 2D the problem is that of placing  $K$  points on a circle so that they are as far as possible from each other. In this case the optimal solution is that of placing the points at the vertices of a regular  $K$ -sided polygon. The 3D analogous of regular

TABLE I  
NUMBER OF REGULAR POLYTOPES AS DIMENSION  $d$  INCREASES.

Dimension $d$	1	2	3	4	$\geq 5$
Number of Regular Polytopes	1	$\infty$	5	6	3

polygons are Platonic Solids. However, the five Platonic solids are not always the unique solutions of the Thomson problem. In fact, only the tetrahedron, octahedron and the icosahedron are the unique solutions for  $K = 4, 6$  and  $12$  respectively. For  $K = 8$ : the cube is not optimal in the sense of the Thomson problem. This means that the energy stabilizes at a minimum in configurations that are not symmetric from a geometric point of view. The unique solution in this case is provided by the vertices of an irregular polytope [38].

The non geometric symmetry between the locations causes the global charge to be different from zero. Therefore in general, when the number of charges is arbitrary, their position on the sphere cannot reach a configuration for which the global charge vanishes to zero. A similar argument holds in higher dimensions for the so called generalized Thomson problem [36]. According to this, we argue that, *the geometric limit to obtain a zero global charge in the generalized Thomson problem is equivalent to the impossibility to learn maximally separated features for an arbitrary number of classes.*

However, since classification is not constrained in a specific dimension as in the case of charges, our approach addresses this issue by *selecting the appropriate dimension of the embedding space so as to have access to symmetrical fixed classifiers directly from regular polytopes.* In dimensions five and higher, there are only three ways to do that (See Tab. I) and they involve the symmetry properties of the three well known regular polytopes available in any high dimensional spaces [39]. These three special classes exist in every dimension and are: the  $d$ -Simplex, the  $d$ -Cube and the  $d$ -Orthoplex. In the next paragraphs the three fixed classifiers derived from them are presented.

**The  $d$ -Simplex Fixed Classifier.** In geometry, a simplex is a generalization of the notion of a triangle or tetrahedron to arbitrary dimensions. Specifically, a  $d$ -Simplex is a  $d$ -dimensional polytope which is the convex hull of its  $d + 1$  vertices. A regular  $d$ -Simplex may be constructed from a regular  $(d - 1)$ -Simplex by connecting a new vertex to all original vertices by the common edge length. According to this, the weights for this classifier can be computed as:

$$\mathbf{W}_S = \left\{ e_1, e_2, \dots, e_{d-1}, \alpha \sum_{i=1}^{d-1} e_i \right\}$$

where  $\alpha = \frac{1-\sqrt{d+1}}{d}$  and  $e_i$  with  $i \in \{1, 2, \dots, d - 1\}$  denotes the standard basis in  $\mathbb{R}^{d-1}$ . The final weights will be shifted about the centroid and normalized. The  $d$ -Simplex fixed classifier defined in an embedding space of dimension  $d$  can accommodate a number of classes equal to its number of vertices:

$$K = d + 1. \quad (4)$$

This classifier has the largest number of classes that can be embedded in  $\mathbb{R}^d$  such that their corresponding class features are equidistant from each other. It can be shown (see appendix) that the angle subtended between *any* pair of weights is equal to:

$$\theta_{\mathbf{w}_i, \mathbf{w}_j} = \arccos\left(-\frac{1}{d}\right) \quad \forall i, j \in \{1, 2, \dots, K\} : i \neq j. \quad (5)$$

**The  $d$ -Orthoplex Fixed Classifier.** This classifier is derived from the  $d$ -Orthoplex (or Cross-Polytope) regular polytope that is defined by the convex hull of points, two on each Cartesian axis of an Euclidean space, that are equidistant from the origin. The weights for this classifier can therefore be defined as:

$$\mathbf{W}_O = \{\pm e_1, \pm e_2, \dots, \pm e_d\}.$$

Since it has  $2d$  vertices, the derived fixed classifier can accommodate in its embedding space of dimension  $d$  a number of distinct classes equal to:

$$K = 2d. \quad (6)$$

Each vertex is adjacent to other  $d - 1$  vertices and the angle between adjacent vertices is

$$\theta_{\mathbf{w}_i, \mathbf{w}_j} = \frac{\pi}{2} \quad \forall i, j \in \{1, 2, \dots, K\} : j \in C(i) \quad (7)$$

Where each  $j \in C(i)$  is an adjacent vertex and  $C$  is the set of adjacent vertices defined as  $C(i) = \{j : (i, j) \in E\}$ .  $E$  is the set of edges of the graph  $G = (\mathbf{W}_O, E)$ . The  $d$ -Orthoplex is the dual polytope of the  $d$ -Cube and vice versa (i.e. the normals of the  $d$ -Orthoplex faces correspond to the the directions of the vertices of the  $d$ -Cube).

**The  $d$ -Cube Fixed Classifier.** The  $d$ -Cube (or Hypercube) is the regular polytope formed by taking two congruent parallel hypercubes of dimension  $(d - 1)$  and joining pairs of vertices, so that the distance between them is 1. A  $d$ -Cube of dimension 0 is one point. The fixed classifier derived from the  $d$ -Cube is constructed by creating a vertex for each binary number in a string of  $d$  bits. Each vertex is a  $d$ -dimensional boolean vector with binary coordinates  $-1$  or  $1$ . Weights are finally obtained from the normalized vertices:

$$\mathbf{W}_C = \left\{ \mathbf{w} \in \mathbb{R}^d : \left[ -\frac{1}{\sqrt{d}}, \frac{1}{\sqrt{d}} \right]^d \right\}.$$

The  $d$ -Cube can accommodate a number of distinct classes equal to:

$$K = 2^d. \quad (8)$$

The vertices are connected by an edge whenever the Hamming distance of their binary numbers is one therefore forming a  $d$ -connected graph. It can be shown (see appendix) that the angle between a vertex with its adjacent (i.e. connected) vertices is:

$$\theta_{\mathbf{w}_i, \mathbf{w}_j} = \arccos\left(\frac{d-2}{d}\right), \forall i, j \in \{1, \dots, K\} : j \in C(i) \quad (9)$$

where  $C(i)$  is the set of vertices adjacent to vertex  $i$ .

Fig. 3 shows the angle between a weight and its adjacent weights computed from Eqs. 5, 7 and 9 as the dimension of the

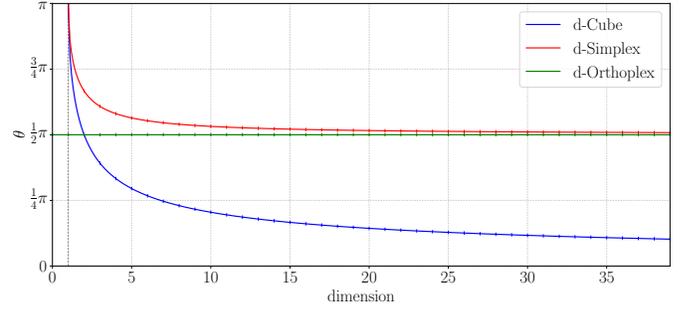


Fig. 3. The angular space defined by RePoNet classifiers. Curves represent the angle between a weight and its adjacent weights as the dimension of the embedding space increases. The angle between class features follows the same trend.

embedding space increases. Having the largest angle between the weights, the  $d$ -Simplex fixed classifier achieves the best inter-class separability. However, as the embedding space dimension increases, its angle tends towards  $\pi/2$ . Therefore, the larger the dimension is, the more similar it becomes to the  $d$ -Orthoplex classifier. The main difference between the two classifiers is in their neighbor connectivity. The different connectivity of the three regular polytope classifiers has a direct influence on the evaluation of the loss. In the case of the  $d$ -Simplex classifier, all the summed terms in the loss of Eq. 3 have always comparable magnitudes in a mini batch.

The  $d$ -Cube classifier has the most compact feature embedding and the angle between each weight and its  $d$  neighbors decreases as the dimension increases. Due to this, it is the hardest to optimize.

#### A. Implementation

Given a classification problem with  $K$  classes, the three RePoNet fixed classifiers can be simply instantiated by defining a non-trainable fully connected layer of dimension  $d$ , where  $d$  is computed from Eqs. 4, 6 and 8 as summarized in Tab II.

TABLE II  
FEATURE DIMENSION  $d$  AS A FUNCTION OF THE NUMBER OF CLASSES  $K$ .

RePoNet	$d$ -Simplex	$d$ -Cube	$d$ -Orthoplex
Layer dim.	$d = K - 1$	$d = \lceil \log_2(K) \rceil$	$d = \lceil \frac{K}{2} \rceil$

In order to accommodate different CNN architectures having different convolutional activation output size (e.g., from the 2048 size of the ResNet50 to the feature size of 10 of the fixed  $d$ -Cube classifier with the 1000 classes of ImageNet), a middle “junction” linear layer (without ReLu) is required.

#### B. Exceeding Vertices as Virtual Negative Classes

Except for the  $d$ -Simplex that allows to assign all its vertices for a given number of classes  $K$ , for both the  $d$ -Cube and the  $d$ -Orthoplex classifiers some of the vertices may be in excess for a given number of classes. As implied by Eq. 6, in the case of the  $d$ -Orthoplex one vertex remains unassigned when the number of classes  $K$  is odd. In the case of the  $d$ -Cube classifier, due to the exponential dependency in Eq. 8,

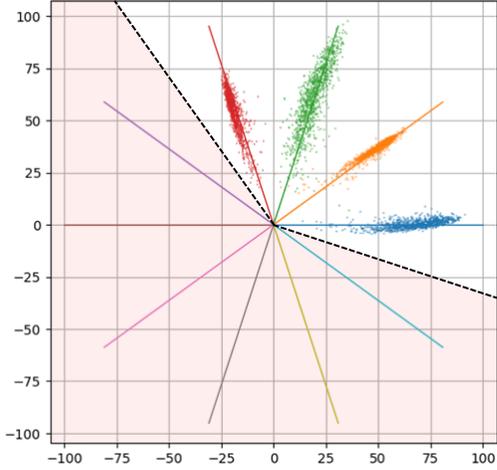


Fig. 4. Learning with unassigned classes in a 2D embedding space. The features of the first four digits of the MNIST dataset are learned using a 10-sided regular polygon in which six of the classes are virtual. Unassigned classes (colored lines inside the shaded region) force a large angular margin region (shaded region) from which features are pushed out.

a large number of vertices may remain not assigned. For example, assuming  $K = 100$  the  $d$ -Cube fixed classifier has 128 vertices (see Tab. II) and 28 of them are not assigned to any class. As shown in [31], unassigned classes act as virtual negative classes forcing a margin around the unassigned weights without affecting the correctness of softmax based cross entropy optimization. Virtual negative classes do not change substantially the objective function of Eq. 3 that can be rewritten as:

$$\mathcal{L} = -\frac{1}{N} \sum_{i=1}^N \log \left( \frac{\exp(\kappa_i \hat{\mathbf{w}}_i^\top \hat{\mathbf{f}}_i)}{\sum_{j=1}^K \exp(\kappa_j \hat{\mathbf{w}}_j^\top \hat{\mathbf{f}}_i) + \sum_{j=K+1}^{K_V} \exp(\kappa_j \hat{\mathbf{w}}_j^\top \hat{\mathbf{f}}_i)} \right) \quad (10)$$

where  $K_V$  is the number of virtual classes (i.e. the exceeding polytope vertices). Fig. 4 illustrates an example similar to Fig. 2(a) in which a 10-sided polygon fixed classifier is learned to classify the first four digits of the MNIST dataset (0, 1, 2 and 3). The remaining six “empty slots” of the classifier are not assigned to any class data and therefore the classifier acts as a virtual negative classifier forcing a large margin (the shaded region) around the virtual class weights (colored lines). This result generalizes the proposed method to any arbitrary number of classes.

### C. Fixed Classifier Decision Boundaries

In binary-classification, the posterior probabilities obtained by softmax in Eq.3 are:

$$p_1 = \frac{\exp(\kappa \hat{\mathbf{w}}_1^\top \hat{\mathbf{f}})}{\exp(\kappa \hat{\mathbf{w}}_1^\top \hat{\mathbf{f}}) + \exp(\kappa \hat{\mathbf{w}}_2^\top \hat{\mathbf{f}})} \quad (11)$$

$$p_2 = \frac{\exp(\kappa \hat{\mathbf{w}}_2^\top \hat{\mathbf{f}})}{\exp(\kappa \hat{\mathbf{w}}_1^\top \hat{\mathbf{f}}) + \exp(\kappa \hat{\mathbf{w}}_2^\top \hat{\mathbf{f}})} \quad (12)$$

where  $\mathbf{f}$  is the learned feature vector and  $\mathbf{w}_1$   $\mathbf{w}_2$  are the fixed classifier weights. The predicted label will be assigned to the class 1 if  $p_1 > p_2$  and to the class 2 if  $p_1 < p_2$ . By comparing the two probabilities  $p_1$  and  $p_2$ ,  $\kappa \hat{\mathbf{w}}_1^\top \hat{\mathbf{f}} + \kappa \hat{\mathbf{w}}_2^\top \hat{\mathbf{f}}$

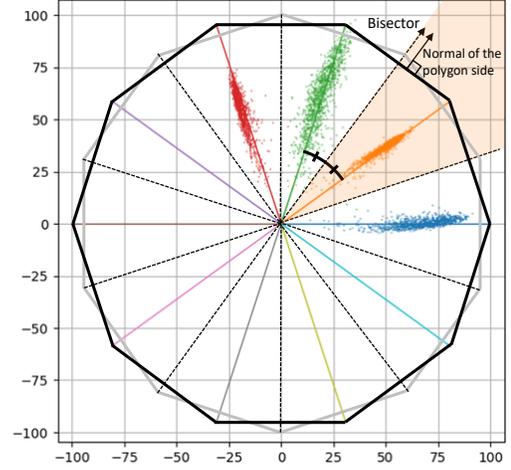


Fig. 5. The intuition behind the decision boundaries in RePoNet (10-sided polygon). The bisector directions (dotted lines), represent the class decision boundaries. They have the same direction of the normal of the corresponding polygon side (only one shown for clarity). The decision boundaries form a regular polygon that is related with the classifier 10-sided polygon according to duality. For clarity, only one class region is highlighted (shaded region).

determines the classification result. The decision boundary is therefore  $\kappa \hat{\mathbf{w}}_1^\top \hat{\mathbf{f}} + \kappa \hat{\mathbf{w}}_2^\top \hat{\mathbf{f}} = 0$ . Due to weight normalization the posterior probabilities result in  $p_1 = \kappa \|\hat{\mathbf{f}}\| \cos(\theta_1)$  and  $p_2 = \kappa \|\hat{\mathbf{f}}\| \cos(\theta_2)$  and since  $p_1$  and  $p_2$  share the same feature  $\hat{\mathbf{f}}$  the equation  $\cos(\theta_1) - \cos(\theta_2) = 0$  is verified at the *angular bisector* between  $\mathbf{w}_1$  and  $\mathbf{w}_2$ . Although the above analysis is built on binary-class case, it can be generalized to the multi-class case [24].

In RePoNet angular bisectors define class decision boundaries that follow a symmetry similar to that of the regular polytope defining the classifier. Specifically, the class decision boundaries and the weights of the classifier are related by the duality relationship that holds between regular polytopes. More practically:

- the set of decision boundaries of the  $d$ -Simplex classifier is shaped as a  $d$ -Simplex;
- the set of the decision boundaries of the  $d$ -Cube classifier is shaped as a  $d$ -Orthoplex;
- the set of the decision boundaries of the  $d$ -Orthoplex classifier is shaped as a  $d$ -Cube.

Class decision boundaries are still defined as a regular polytope and the features located within such boundaries are therefore maximally separated. The basic intuition behind this result can be better appreciated in 2D exploiting the well known result that all the regular polygons are self-dual [39]. That is, the normal of each side of a regular polygon is parallel to the direction from the origin towards the vertex of its dual polygon. Fig. 5 shows the example introduced in Fig. 4 in which decision boundaries are highlighted with dotted lines according to the dual regular polygon. Fig. 6 illustrates the duality relationship between the weights of the three fixed classifiers proposed and their decision boundaries in the 3D embedding space.

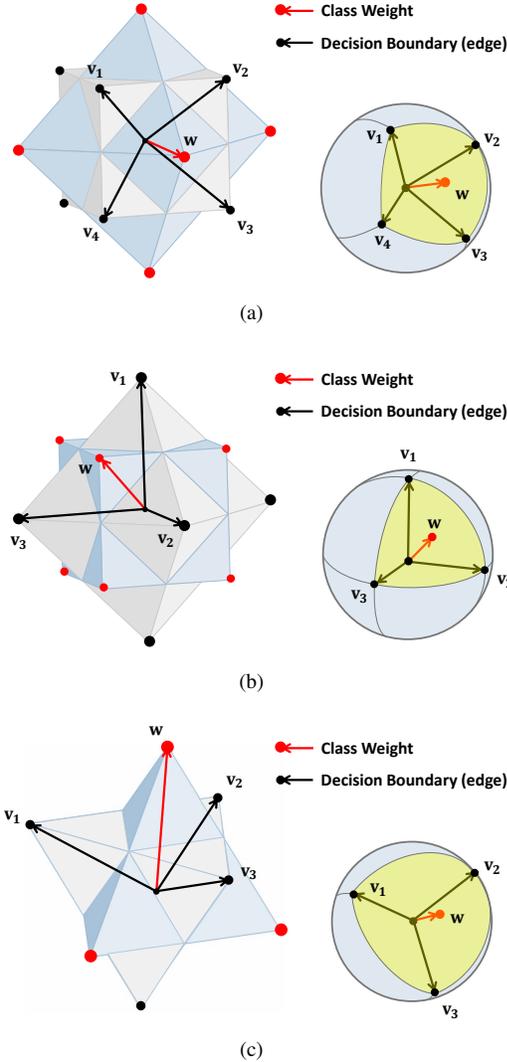


Fig. 6. The RePoNet fixed classifiers decision boundaries in a 3D embedding space: (a):  $d$ -Orthoplex classifier; (b):  $d$ -Cube classifier; (c):  $d$ -Simplex classifier. On the *left*: the regular polytope classifier (light blue); its dual polytope (grey); a classifier weight  $w$  (red) and its edge decision boundaries  $v_1, v_2, \dots$  (black). On the *right*: the same entities on the unit sphere. The yellow region indicates where class features are located (only one class weight and the corresponding edge decision boundaries are shown for clarity). The characterization extends to arbitrary dimensions.

## V. EXPERIMENTAL RESULTS

We evaluate the correctness (Sec. V-A) and the no loss of performance of our approach with respect to standard baselines using trainable and fixed classifiers across a range of datasets and architectures (Sec. V-B). All the experiments are conducted with the well known MNIST, FashionMNIST [40], EMNIST [41], CIFAR-10, CIFAR-100 [42] and ImageNet (ILSVRC2012) [43] datasets. We chose several common CNN architectures (i.e. LeNet, VGG, ResNet, DenseNet), as well as more recent ones (i.e. SeResNeXt50 [44], SkResNeXt50 [45] and EfficientNet [46]) that have shown to improve performance while maintaining or, in some cases, reducing computational complexity and model size.

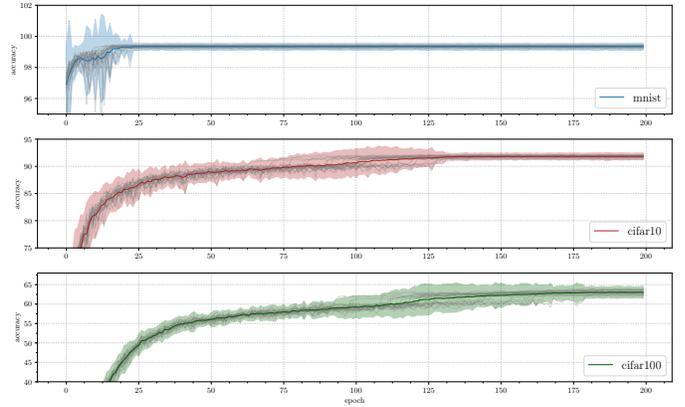


Fig. 7. Class permutation verification. Average accuracy curves and confidence interval computed from the MNIST, CIFAR-10 and CIFAR-100 datasets (from top to bottom, respectively) under different random permutations of the ground truth labels position.

### A. Hard Permutations Verification

Since fixed classifiers cannot rely on an adjustable set of subspaces for class feature representation, we verified if some permutations are harder than others for our proposed method. The presence of such hard permutations would preclude the general applicability of our method. The standard trainable classifier does not suffer from this problem, when features cannot be well separated trainable classifiers can rearrange the feature subspace directions so that the previous convolutional layers can better disentangle the non-linear interactions between complex data patterns. Instead, fixed classifiers demand this capability to all the previous layers.

According to this, we generate random permutations of the ground truth label positions<sup>3</sup> and a new model is learned for each permuted dataset. Fig. 7 shows the mean and the 95% confidence interval computed from the accuracy curves of the learned models. To provide further insight into this analysis, 20 out of 500 accuracy curves computed for each dataset are also shown. Specifically, the evaluation is performed on three different datasets with an increasing level of complexity (i.e. MNIST, CIFAR-10 and CIFAR-100). All the models are trained for 200 epochs to make sure that the models trained with CIFAR-100 achieve convergence.

In order to address the most severe possible outcomes that may happen, for this experiment we used the  $d$ -Cube fixed classifier. Being the hardest to optimize, this experiment can be regarded as a worst case analysis scenario for our method. As shown in the same figure, the performance is substantially insensitive to both permutations and datasets. The average reduction in performance at the end of the training process is negligible and the confidence intervals reflect the complexity of the datasets. Although the space of permutations cannot be exhaustively evaluated even for a small number of classes, we have achieved proper convergence for the whole set of 1500 learned models. The experiment took 5 days on a Nvidia DGX-1.

<sup>3</sup>This is equivalent to randomly permuting the classifier weight set  $\mathbf{W} = \{\mathbf{w}_j\}_{j=1}^K$ .

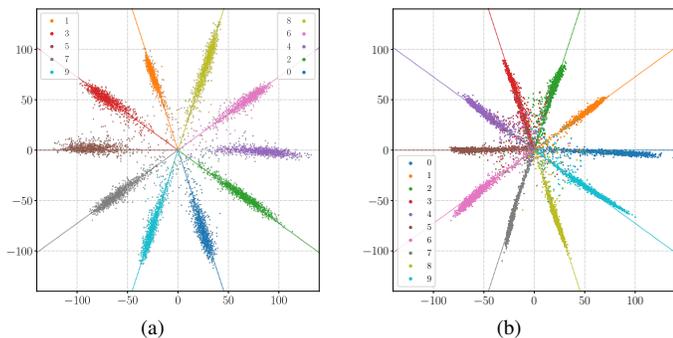


Fig. 8. The distribution of features learned using a 10-sided regular polygon. (a): A special permutation of classes is shown in which the MNIST even and odd digits are placed in the positive and negative half-space of the abscissa respectively. (b): The features learned using the CIFAR-10 dataset.

On the basis of this evidence, we can conclude that fixing the classifier (therefore not having access to a set of adjustable subspaces for class feature representation) does not affect the expressive power of neural networks. This experiment also provides a novel and more systematic empirical evidence of the general applicability and correctness of fixed classifiers with respect to [1] where only one permutation was tested.

We finally report qualitative results of a learned permuted dataset. Fig. 8(a) shows features learned in a  $k$ -sided polygon ( $2d$  embedding space) on the MNIST dataset. In particular the model is learned with a special permutation of the labels (manually selected) that places even and odd digits features respectively on the positive and negative half space of the abscissa. Fig. 8(b) shows the features of on CIFAR-10 learned with a similar 10-sided-polygon. It can be noticed that features are distributed following the same polygonal pattern shown in Fig. 8(a).

### B. Generalization and Performance Evaluation

Having verified that the order position of the class labels does not adversely affect the proposed method, in this section we evaluate the classification performance of RePoNet on the following datasets: MNIST, EMNIST, FashionMNIST, CIFAR-10, CIFAR-100 and ImageNet. The RePoNet method is compared with CNN baselines with learned classifiers and the fixed classifier method reported in [1], that has been implemented for different architectures and different dimensions of the embedding space. Except for the final fixed classifier all the compared methods have exactly the same architecture and training settings as the one that RePoNet uses.

1) **MNIST and CIFAR**: We trained the so called LeNet++ architecture [47] on all the MNIST family datasets. The network is a modification of the LeNet [48] to a deeper and wider network including parametric rectifier linear units (pReLU) [49]. For the evaluation on the CIFAR-10 and CIFAR-100 datasets, we further trained VGG [50] with depth 13 and 19, ResNet50 [11], SeNet [51] and DenseNet169 [52]. Popular network architectures for ImageNet require modifications to adapt to the CIFAR 32x32 input size. According to this, our

experiments follow the publicly available implementations<sup>4</sup>. We compared all the variants of our approach for each architecture including trainable classifiers with different dimensions of the feature space. The mini batch size is 256 for both the MNIST family datasets and the CIFAR-10/100 datasets. For the CIFAR datasets, we compared both a “vanilla” learning setup with no hyperparameters tuning based on the Adam optimizer (learning rate 0.0005) for the VGG architectures and a learning setup based on SGD with a specific learning rate schedule (starting from 0.1 and decreasing by a factor of 10 after 150 and 250 epochs) for ResNet50, SEnet18 and DenseNet169 architectures. As hyperparameters tuning is an integral part of Deep Learning we provided two opposite learning setup.

Test-set accuracy for this experiment is reported in Tab. III, IV and V for MNISTs, CIFAR-10 and CIFAR-100, respectively. In addition to the well-known MNIST and FashionMnist, we included EMNIST dataset having 47 classes including lower/upper case letters and digits. This allows to quantify with a specific dataset and architecture, as in CIFAR-10 and CIFAR-100, the classification accuracy with a higher number of classes. Each entry in the tables report the test-set accuracy. The subscript indicates the specific feature space dimension  $d$  used. The results reveal and confirm that the RePoNet method achieves comparable classification accuracy of other trainable classifier models. This evidence is in agreement on all the combinations of datasets, architectures, number of classes and feature space dimensions considered. All the RePoNet variants exhibit similar behavior even in hard combinations such as the CIFAR-100 dataset in a low dimensional feature space. For example, the RePoNet  $d$ -Cube fixed classifier implemented with the VGG19 architecture achieves an accuracy of 65.32% in a  $d = 7$  dimensional feature space. A fully trainable classifier in a feature space of dimension  $d = 512$  (i.e. two orders of magnitude larger), achieves a moderate improvement of about 3% (68.47%). On the other hand, with a significantly lower feature dimension of  $d = 50$ , RePoNet  $d$ -Orthoplex improves the accuracy to 69.76%. All the RePoNet variants exhibit similar behavior also in the case of more sophisticated architectures trained with SGD scheduled learning rates to match state-of-the-art performance. RePoNet classifiers are both agnostic to architectures and training setup and are able to improve accuracy similar to trainable classifiers.

Results also show that the Hadamard fixed classifier [1] does not succeed to learn when the number of classes is larger than the number of unique weight directions in the embedding space (i.e.  $d < K$ ). As expected, this effect is present for simple datasets as the MNIST digits dataset, however as reported in [1] Section 4.2 (Possible Caveats) as the number of classes  $K$  increases the effect is less pronounced.

When  $d \approx K$  or  $d > K$ , classification performance is similar. However, as shown in Fig. 9(a) RePoNet converges faster than [1], and with the same speed as the trainable baselines. Our conjecture is that with our symmetrical fixed classifiers, each term in the loss function tends to have the same magnitude

<sup>4</sup><https://github.com/bearpaw/pytorch-classification> and <https://github.com/kuangliu/pytorch-cifar>

TABLE III  
REPORTED ACCURACY (%) OF THE REPONET METHOD ON MNIST, EMNIST, FASHIONMNIST DATASETS ON DIFFERENT COMBINATIONS OF ARCHITECTURES AND RELATIVE LEARNED CLASSIFIER BASELINES.

ARCHITECTURE	MNIST	EMNIST	FASHIONMNIST
	( $K = 10$ )	( $K = 47$ )	( $K = 10$ )
	LENET++		
RePoNet $K$ -sided-polygon	99.24 <sub><math>d=2</math></sub>	72.81 <sub><math>d=2</math></sub>	92.48 <sub><math>d=2</math></sub>
Hadamard fixed classifier [1]	21.14 <sub><math>d=2</math></sub>	4.12 <sub><math>d=2</math></sub>	19.89 <sub><math>d=2</math></sub>
Learned Classifier	99.21 <sub><math>d=2</math></sub>	73.08 <sub><math>d=2</math></sub>	92.79 <sub><math>d=2</math></sub>
RePoNet $d$ -Cube	99.58 <sub><math>d=4</math></sub>	88.12 <sub><math>d=6</math></sub>	94.01 <sub><math>d=4</math></sub>
Hadamard fixed classifier [1]	41.99 <sub><math>d=4</math></sub>	15.12 <sub><math>d=6</math></sub>	37.16 <sub><math>d=4</math></sub>
Learned Classifier	99.41 <sub><math>d=4</math></sub>	86.96 <sub><math>d=6</math></sub>	93.94 <sub><math>d=4</math></sub>
RePoNet $d$ -Orthoplex	99.66 <sub><math>d=5</math></sub>	88.19 <sub><math>d=24</math></sub>	94.84 <sub><math>d=5</math></sub>
Hadamard fixed classifier [1]	79.34 <sub><math>d=5</math></sub>	60.34 <sub><math>d=24</math></sub>	74.22 <sub><math>d=5</math></sub>
Learned Classifier	99.07 <sub><math>d=5</math></sub>	87.66 <sub><math>d=24</math></sub>	94.21 <sub><math>d=5</math></sub>
RePoNet $d$ -Simplex	99.71 <sub><math>d=9</math></sub>	88.89 <sub><math>d=46</math></sub>	94.29 <sub><math>d=9</math></sub>
Hadamard fixed classifier [1]	99.12 <sub><math>d=9</math></sub>	88.48 <sub><math>d=46</math></sub>	94.30 <sub><math>d=9</math></sub>
Learned classifier	99.41 <sub><math>d=9</math></sub>	88.33 <sub><math>d=46</math></sub>	94.41 <sub><math>d=9</math></sub>
Hadamard fixed classifier [1]	99.54 <sub><math>d=512</math></sub>	88.35 <sub><math>d=512</math></sub>	94.14 <sub><math>d=512</math></sub>
Learned classifier	99.29 <sub><math>d=512</math></sub>	88.87 <sub><math>d=512</math></sub>	94.28 <sub><math>d=512</math></sub>

TABLE IV  
REPORTED ACCURACY (%) OF THE REPONET METHOD ON THE CIFAR-10 DATASET ON DIFFERENT COMBINATIONS OF ARCHITECTURES AND RELATIVE BASELINES.

ARCHITECTURE	CIFAR-10 ( $K = 10$ )				
	VGG13	VGG19	RESNET50	SENET18	DENSENET169
	Training hyperparameters	ADAM	ADAM	SGD	SGD
RePoNet $K$ -sided-polygon	90.79 <sub><math>d=2</math></sub>	91.54 <sub><math>d=2</math></sub>	92.78 <sub><math>d=2</math></sub>	92.63 <sub><math>d=2</math></sub>	92.74 <sub><math>d=2</math></sub>
Hadamard fixed classifier [1]	19.45 <sub><math>d=2</math></sub>	19.19 <sub><math>d=2</math></sub>	19.69 <sub><math>d=2</math></sub>	19.77 <sub><math>d=2</math></sub>	19.76 <sub><math>d=2</math></sub>
Learned classifier	90.41 <sub><math>d=2</math></sub>	91.17 <sub><math>d=2</math></sub>	93.15 <sub><math>d=2</math></sub>	93.25 <sub><math>d=2</math></sub>	92.89 <sub><math>d=2</math></sub>
RePoNet $d$ -Cube	92.26 <sub><math>d=4</math></sub>	92.58 <sub><math>d=4</math></sub>	94.86 <sub><math>d=4</math></sub>	94.96 <sub><math>d=4</math></sub>	93.94 <sub><math>d=4</math></sub>
Hadamard fixed classifier [1]	37.19 <sub><math>d=4</math></sub>	36.95 <sub><math>d=4</math></sub>	37.89 <sub><math>d=4</math></sub>	38.05 <sub><math>d=4</math></sub>	38.12 <sub><math>d=4</math></sub>
Learned classifier	92.14 <sub><math>d=4</math></sub>	92.21 <sub><math>d=4</math></sub>	95.03 <sub><math>d=4</math></sub>	94.95 <sub><math>d=4</math></sub>	94.97 <sub><math>d=4</math></sub>
RePoNet $d$ -Orthoplex	92.51 <sub><math>d=5</math></sub>	92.47 <sub><math>d=5</math></sub>	95.25 <sub><math>d=5</math></sub>	95.05 <sub><math>d=5</math></sub>	95.16 <sub><math>d=5</math></sub>
Hadamard fixed classifier [1]	73.77 <sub><math>d=5</math></sub>	72.46 <sub><math>d=5</math></sub>	75.99 <sub><math>d=5</math></sub>	75.95 <sub><math>d=5</math></sub>	75.73 <sub><math>d=5</math></sub>
Learned classifier	92.28 <sub><math>d=5</math></sub>	92.21 <sub><math>d=5</math></sub>	95.18 <sub><math>d=5</math></sub>	95.08 <sub><math>d=5</math></sub>	95.41 <sub><math>d=5</math></sub>
RePoNet $d$ -Simplex	92.71 <sub><math>d=9</math></sub>	92.59 <sub><math>d=9</math></sub>	95.66 <sub><math>d=9</math></sub>	95.36 <sub><math>d=9</math></sub>	95.32 <sub><math>d=9</math></sub>
Hadamard fixed classifier [1]	92.03 <sub><math>d=9</math></sub>	92.37 <sub><math>d=9</math></sub>	95.53 <sub><math>d=9</math></sub>	95.25 <sub><math>d=9</math></sub>	94.92 <sub><math>d=9</math></sub>
Learned classifier	91.89 <sub><math>d=9</math></sub>	92.60 <sub><math>d=9</math></sub>	95.08 <sub><math>d=9</math></sub>	95.20 <sub><math>d=9</math></sub>	95.32 <sub><math>d=9</math></sub>
Hadamard fixed classifier [1]	90.11 <sub><math>d=512</math></sub>	88.32 <sub><math>d=512</math></sub>	95.36 <sub><math>d=512</math></sub>	95.49 <sub><math>d=512</math></sub>	95.68 <sub><math>d=512</math></sub>
Learned classifier	92.34 <sub><math>d=512</math></sub>	92.42 <sub><math>d=512</math></sub>	95.53 <sub><math>d=512</math></sub>	95.26 <sub><math>d=512</math></sub>	95.68 <sub><math>d=512</math></sub>

centered around the mean of the distribution (i.e. the von Mises-Fisher distribution is similar to the Normal distribution) and therefore the average computed in the loss is a good estimator. Instead, in the Hadamard classifier the terms may have different magnitudes and “important” errors in the loss may not be taken into account correctly by simple averaging.

2) *ImageNet*: Finally, we evaluated our method on the 1000 object category classification problem defined by the ImageNet dataset. This dataset consists of a 1.2M image training set and a 100k image test set. We compared all the variants of our approach on different combinations of architectures and their relative trainable classifiers. The comparison also

includes the Hadamard classifier.

Experiments have been conducted in two different configurations of the training hyperparameters. *First*, we performed experiments using the Adam optimizer and simple augmentation based on random cropping and horizontal flipping on well-established networks such as ResNet50 [11] and DenseNet169 [52]. The learning rate is automatically adjusted when a plateau in model performance is detected. We trained for 250 epochs with batch size 64 with an initial learning rate of 0.0005. With this configuration, we aim to evaluate our method without performing any specific hyperparameter optimization or exploiting large computational resources. *Second*, we eval-

TABLE V  
 REPORTED ACCURACY (%) OF THE REPONET METHOD ON THE CIFAR-100 DATASET ON DIFFERENT COMBINATIONS OF ARCHITECTURES AND RELATIVE BASELINES.

ARCHITECTURE	CIFAR-100 ( $K = 100$ )				
	VGG13	VGG19	RESNET50	SENET18	DENSENET169
Training hyperparameters	ADAM	ADAM	SGD	SGD	SGD
RePoNet $K$ -sided-polygon	36.22 <sub><math>d=2</math></sub>	37.65 <sub><math>d=2</math></sub>	33.39 <sub><math>d=2</math></sub>	35.26 <sub><math>d=2</math></sub>	30.04 <sub><math>d=2</math></sub>
Hadamard fixed classifier [1]	1.75 <sub><math>d=2</math></sub>	1.75 <sub><math>d=2</math></sub>	1.61 <sub><math>d=2</math></sub>	1.80 <sub><math>d=2</math></sub>	1.64 <sub><math>d=2</math></sub>
Learned classifier	37.56 <sub><math>d=2</math></sub>	35.83 <sub><math>d=2</math></sub>	33.30 <sub><math>d=2</math></sub>	40.57 <sub><math>d=2</math></sub>	32.87 <sub><math>d=2</math></sub>
RePoNet $d$ -Cube	64.35 <sub><math>d=7</math></sub>	65.32 <sub><math>d=7</math></sub>	67.27 <sub><math>d=7</math></sub>	69.38 <sub><math>d=7</math></sub>	68.99 <sub><math>d=7</math></sub>
Hadamard fixed classifier [1]	5.96 <sub><math>d=7</math></sub>	5.52 <sub><math>d=7</math></sub>	5.91 <sub><math>d=7</math></sub>	6.27 <sub><math>d=7</math></sub>	6.08 <sub><math>d=7</math></sub>
Learned classifier	64.11 <sub><math>d=7</math></sub>	65.29 <sub><math>d=7</math></sub>	74.96 <sub><math>d=7</math></sub>	75.29 <sub><math>d=7</math></sub>	75.51 <sub><math>d=7</math></sub>
RePoNet $d$ -Orthoplex	68.78 <sub><math>d=50</math></sub>	69.76 <sub><math>d=50</math></sub>	78.23 <sub><math>d=50</math></sub>	77.24 <sub><math>d=50</math></sub>	79.41 <sub><math>d=50</math></sub>
Hadamard fixed classifier [1]	43.88 <sub><math>d=50</math></sub>	43.89 <sub><math>d=50</math></sub>	50.33 <sub><math>d=50</math></sub>	49.56 <sub><math>d=50</math></sub>	50.65 <sub><math>d=50</math></sub>
Learned classifier	68.13 <sub><math>d=50</math></sub>	68.41 <sub><math>d=50</math></sub>	78.22 <sub><math>d=50</math></sub>	77.15 <sub><math>d=50</math></sub>	78.83 <sub><math>d=50</math></sub>
RePoNet $d$ -Simplex	68.61 <sub><math>d=99</math></sub>	68.69 <sub><math>d=99</math></sub>	79.02 <sub><math>d=99</math></sub>	78.20 <sub><math>d=99</math></sub>	80.01 <sub><math>d=99</math></sub>
Hadamard fixed classifier [1]	67.23 <sub><math>d=99</math></sub>	67.18 <sub><math>d=99</math></sub>	78.82 <sub><math>d=99</math></sub>	77.21 <sub><math>d=99</math></sub>	79.41 <sub><math>d=99</math></sub>
Learned classifier	68.15 <sub><math>d=99</math></sub>	68.87 <sub><math>d=99</math></sub>	78.58 <sub><math>d=99</math></sub>	77.42 <sub><math>d=99</math></sub>	79.05 <sub><math>d=99</math></sub>
Hadamard fixed classifier [1]	63.16 <sub><math>d=512</math></sub>	64.46 <sub><math>d=512</math></sub>	78.78 <sub><math>d=512</math></sub>	77.94 <sub><math>d=512</math></sub>	79.44 <sub><math>d=512</math></sub>
Learned classifier	68.56 <sub><math>d=512</math></sub>	68.47 <sub><math>d=512</math></sub>	77.96 <sub><math>d=512</math></sub>	77.63 <sub><math>d=512</math></sub>	79.63 <sub><math>d=512</math></sub>

uated our method with more sophisticated CNN architectures, namely SKresNeXt, SEresNeXt, and EfficientNet with related training hyperparameters. With this configuration, the aim is to evaluate whether our method can reach state-of-the-art performance. The SKresNeXt and SEresNeXt architectures integrate the SE and SK blocks, [51] and [53] respectively, with the ResNeXt architecture [54]. The benefit of these variants is to maintain computational complexity and model size similar to the SEnet and SKnet architectures while further improving performance. The third architecture, EfficientNet [46], achieves state-of-the-art performance using significantly fewer parameters than other state-of-the-art models. As these architectures typically require a large effort to tune the training hyperparameters, we trained our method on top of these models following the settings reported in the original papers. Specifically, we train EfficientNet-B2 following [46]: RMSProp optimizer with decay 0.9 and momentum 0.9; batch norm momentum 0.99; initial learning rate 0.256 that decays by 0.97 every 2.4 epochs; weight decay  $1e-5$ . Analogously, SKresNeXt50 and SEresNeXt50 are trained following the ResNeXt50 [54]: SGD optimizer, weight decay 0.0001; momentum 0.9; initial learning rate of 0.1, divided by 10 for three times using a specific schedule reported in the paper. For all the three models we used automated data augmentation techniques from [55] (RandAugment) with distortion magnitude 7. SeResNeXt50 and SkResNeXt50 were trained for 250 epochs with 192 batch size. EfficientNet-B2 was trained for 450 epochs with 120 batch size. Our evaluation is based on the pytorch-image-models<sup>5</sup> repository.

Tab. VI summarizes our results. As can be clearly noticed, except for the  $d$ -Cube there is no substantial difference between the performance of our fixed classifiers and the

learned classifiers. This holds also in the case of the learned classifiers in their original architecture implementation (shown in the bottom line of the Table). The table also shows that RePoNet accuracy is comparable with the Hadamard fixed classifier [1]. As in the cases of CIFAR-10 and CIFAR-100, also with ImageNet the accuracy of the  $d$ -Cube is lower than the corresponding learned classifiers. We argue this is mainly due to the difficulty of performing optimization in the  $d = 10$  dimensional space due to the fact that the angle between each class weight vector and its  $d$  adjacent weight vectors approaches to zero as the dimension increases (Fig. 3). However, the  $d$ -Cube classifier shows the largest relative improvement as the representational power of the architecture increases (left to right). For example, the accuracy of the  $d$ -Cube-EFFICIENTNET-B2 fixed classifier is 12.18 percentage points larger than the  $d$ -Cube-RESNET50 (i.e.  $75.62 - 63.44 = 12.18$ ). This relative performance improvement is substantially higher than that of the corresponding learned classifier (i.e.  $77.42 - 68.82 = 8.6$ ). This result is quantitatively consistent with the underlying assumption of this paper and provides further support on the fact that the adjustable capability of the final classifier can be successfully demanded to previous layers. The other two RePoNet variants substantially achieve the same accuracy of the learned classifiers, irrespective whether they have similar ( $d = \{10, 500, 999\}$ ) or higher feature space dimension ( $d = \{2048, 1669, 1408\}$ ) as in their original architecture implementations. They do not show sensible relative performance improvement with increasing representational power of the network. More importantly, both the  $d$ -Simplex and the  $d$ -Orthoplex classifiers reach state-of-the-art accuracy (around 80%) when combined with competitive architectures. This confirms the validity and the absence of a loss of generalization of our method.

<sup>5</sup><https://github.com/rwightman/pytorch-image-models>

TABLE VI  
REPORTED ACCURACY (%) OF THE REPO NET METHOD ON THE IMAGENET DATASET ON DIFFERENT COMBINATIONS OF ARCHITECTURES AND RELATIVE CLASSIFIER BASELINES.

ARCHITECTURE	RESNET50	DENSENET169	SERESNEXT50	SKRESNEXT50	EFFICIENTNET-B2
Training hyperparameters	ADAM	ADAM	SGD+RandAug	SGD+RandAug	RMSPROP+RandAug
RePoNet $d$ -Cube	63.44 <sub><math>d=10</math></sub>	63.63 <sub><math>d=10</math></sub>	73.58 <sub><math>d=10</math></sub>	74.80 <sub><math>d=10</math></sub>	75.62 <sub><math>d=10</math></sub>
Learned classifier	68.82 <sub><math>d=10</math></sub>	68.03 <sub><math>d=10</math></sub>	76.66 <sub><math>d=10</math></sub>	77.49 <sub><math>d=10</math></sub>	77.42 <sub><math>d=10</math></sub>
RePoNet $d$ -Orthoplex	73.71 <sub><math>d=500</math></sub>	74.20 <sub><math>d=500</math></sub>	79.95 <sub><math>d=500</math></sub>	79.66 <sub><math>d=500</math></sub>	80.07 <sub><math>d=500</math></sub>
Learned classifier	73.67 <sub><math>d=500</math></sub>	73.70 <sub><math>d=500</math></sub>	77.60 <sub><math>d=500</math></sub>	80.18 <sub><math>d=500</math></sub>	79.27 <sub><math>d=500</math></sub>
RePoNet $d$ -Simplex	74.13 <sub><math>d=999</math></sub>	74.03 <sub><math>d=999</math></sub>	80.25 <sub><math>d=999</math></sub>	80.17 <sub><math>d=999</math></sub>	80.61 <sub><math>d=999</math></sub>
Learned classifier	73.96 <sub><math>d=999</math></sub>	73.37 <sub><math>d=999</math></sub>	77.99 <sub><math>d=999</math></sub>	80.08 <sub><math>d=999</math></sub>	79.36 <sub><math>d=999</math></sub>
Hadamard fixed classifier [1]	74.07 <sub><math>d=2048</math></sub>	73.95 <sub><math>d=1669</math></sub>	80.25 <sub><math>d=2048</math></sub>	80.19 <sub><math>d=2048</math></sub>	79.74 <sub><math>d=1408</math></sub>
Learned classifier	74.11 <sub><math>d=2048</math></sub>	74.01 <sub><math>d=1669</math></sub>	79.95 <sub><math>d=2048</math></sub>	80.09 <sub><math>d=2048</math></sub>	80.57 <sub><math>d=1408</math></sub>

Finally, Tab. VII shows the total number of parameters for each network in comparison with their original learned classifiers (i.e. bottom line in Tab. VI). The  $d$ -Orthoplex-EFFICIENTNET-B2 fixed classifier saves 7.74% of the network parameters while achieving the same accuracy (around 80%). It is worth to notice that the  $d$ -Cube-EFFICIENTNET-B2 with 7.7M of parameters (15.31% savings) achieves similar accuracy of a vanilla ResNet50 baseline (i.e. around 75% accuracy) having 25.5M of parameters.

### C. Training Time

The time it takes to train a neural network model to address a classification problem is typically considered as the product of the training time per epoch and the number of epochs which need to be performed to reach the desired level of accuracy [56]. Although in our case the training time per epoch is lower (the weights of the fixed classifier do not require back-propagation), it has a negligible effect due to the number of epochs required to reach a reasonable desired level of accuracy. In Fig. 9 and Fig. 10 we report the classification accuracy over the epochs for the two different configurations of the training hyperparameters we evaluated.

Specifically, Fig. 9(a)(top) and Fig. 9(a)(bottom) show the training error and the classification accuracy, respectively, over the epochs. The curves are obtained on the CIFAR100 dataset, using the VGG19 architecture and training is performed according to the Adam stochastic optimization. Fig. 9(b) shows the accuracy curves of the proposed three fixed classifiers

and the best performing learned classifier (i.e.  $d = 999$ ). The curves are obtained on the ImageNet dataset with the DenseNet169 architecture and learned according to the Adam optimizer. As can be noticed,  $d$ -Simplex and  $d$ -Orthoplex classifiers have lower or equal time to reach any desired level of accuracy than the learned and Hadamard fixed classifiers. The  $d$ -Cube classifier is the slowest and does not reach a comparable final performance. This is due to the different feature dimension ( $d = 10$ ) and topology. However, when compared with a learned classifier with same feature dimension (as discussed in the next paragraph) the training time is similar.

Fig. 10(a) and Fig. 10(b) show the time to reach accuracy using SGD+RandAug on SeResNeXt50 for the  $d$ -Simplex and  $d$ -Orthoplex (in red) fixed classifiers and the learned classifiers (in blue), respectively. As evidenced in the figure, the learned classifiers require about 150 epochs to obtain the accuracy that the  $d$ -Simplex and  $d$ -Orthoplex achieve in 50 and 90 epochs, respectively. Although this gain reduces as training progresses towards the end, our method achieves consistently better results. Fig. 10(c) shows that the  $d$ -Cube classifier requires similar training time with slightly lower final accuracy.

The general behavior of the curves shown in Fig. 9 and Fig. 10 is consistent across combinations of datasets, architectures, classifiers and training strategies. The training time to reach the same accuracy is shorter or equal for our method and the time reduction follows the complexity of the embeddings defined by each regular polytope fixed classifier.

Overall, we have demonstrated that Regular Polytope Networks provide a novel, effective and easy approach to fixed classifiers that achieves comparable state-of-the-art performance with the standard trainable classifiers. They provide faster speed of convergence and a significant reduction in model parameters. To facilitate replication of our experiments, the code will be made publicly available.

## VI. DISCUSSION: POTENTIAL AND CHALLENGES

Our finding may have implications in those Deep Neural Network learning contexts in which a classifier must be robust

TABLE VII  
THE NUMBER OF PARAMETERS OF EACH NETWORK AND THE PERCENTAGE (%) OF SAVED PARAMETERS ON THE IMAGENET DATASET ( $d$  INDICATES THE FEATURE DIMENSION).

ARCHITECTURE	PARAM#	SAVED PARAMS ( % )		
		$d$ -CUBE	$d$ -ORTHOPEX	$d$ -SIMPLEX
DenseNet169	14.15M	11.65 <sub><math>d=10</math></sub>	5.89 <sub><math>d=500</math></sub>	0.02 <sub><math>d=999</math></sub>
ResNet50	25.56M	7.94 <sub><math>d=10</math></sub>	4.01 <sub><math>d=500</math></sub>	0.01 <sub><math>d=999</math></sub>
SeResNeXt50	27.56M	7.36 <sub><math>d=10</math></sub>	3.72 <sub><math>d=500</math></sub>	0.01 <sub><math>d=999</math></sub>
SkResNeXt50	27.50M	7.38 <sub><math>d=10</math></sub>	3.73 <sub><math>d=500</math></sub>	0.01 <sub><math>d=999</math></sub>
EfficientNet-B2	9.11M	15.31 <sub><math>d=10</math></sub>	7.74 <sub><math>d=500</math></sub>	0.03 <sub><math>d=999</math></sub>

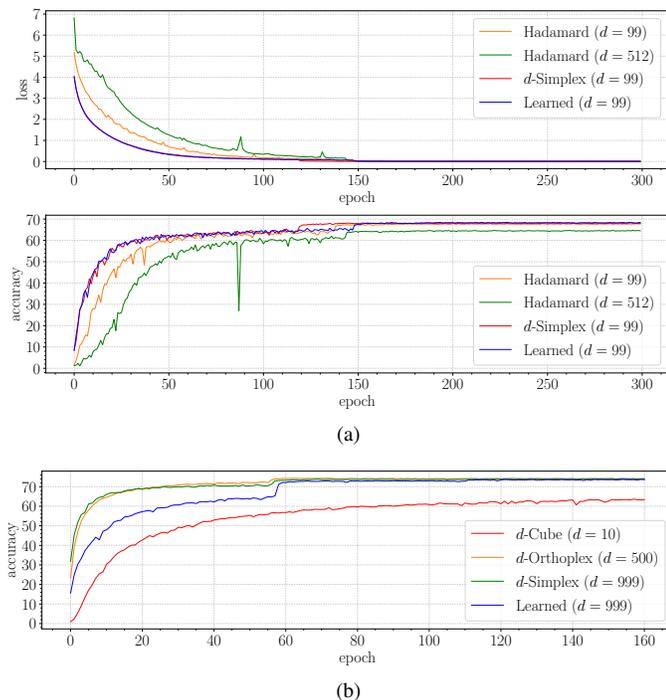


Fig. 9. Speed of convergence comparison (ADAM). (a): Training error curves (top) and test accuracy curves (bottom) using the CIFAR-100 dataset with the VGG19 architecture. (b): ImageNet learning speed using DenseNet169. As evidenced from the figures, the proposed method has faster convergence.

against changes of the feature representation while learning. This is the case of incremental learning settings, especially when features are stored in memory banks while learning [57], [58], [59]. Despite recent advances, methods inspired by memory-augmented deep neural networks are still limited when it comes to incremental learning. The method [60] simplifies the original fully differentiable end-to-end idea. Except for the nearest-neighbor query to the memory bank, their approach is fully differentiable, can be trained end-to-end and operates in an incremental manner (i.e. without the need of resetting during training). However, the features stored in the memory bank remain fixed (i.e. they are not undergoing learning) and only the memory bank is learned. Our approach may have a promising potential for learning both the feature and the memory without considering their joint learning. The intuition is that every time the internal feature representation changes the memory bank must be relearned from scratch. Our method can mitigate the need of feature relearning by keeping the compatibility of features between learning steps thanks to their stationarity. Concurrent to this work, [61] addresses a similar problem in terms of feature “back-compatibility” and exploits a pre-trained fixed classifier to avoid re-indexing a memory bank containing the gallery features of a retrieval system that has been updated.

This basic idea can be in principle applied to the many computer vision tasks that have benefited from memory based learning. Among them we mention [62], [63], [64] for cumulative learning of face appearance models from video stream, [65], [66], [67], [68], [69] for object detection, [70] for video object segmentation and [71] for visual object tracking. The

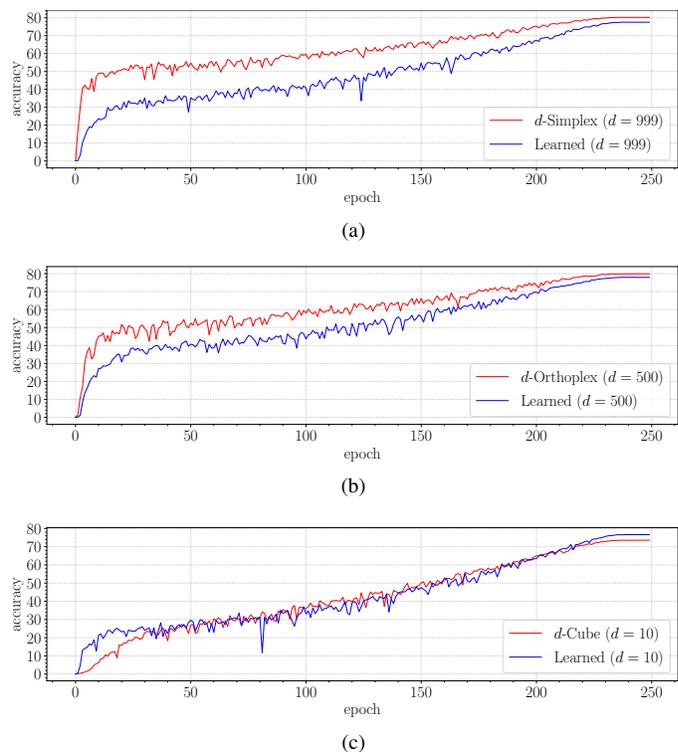


Fig. 10. Speed of convergence comparison (SGD+RANDAUG). Test accuracy curves over the epochs on the ImageNet test set for the SeResNeXt50 architecture using the proposed fixed classifiers (red) and the standard trainable baselines (blue). (a): The  $d$ -Simplex classifier, (b): the  $d$ -Orthoplex classifier and (c): the  $d$ -Cube classifier. The time to reach the same accuracy is shorter or equal for our method.

works [65], [66], [67], [68], [69] accumulate context from pre-computed feature banks (with fixed pre-trained feature extractors i.e. not undergoing learning). The feature banks extend the time horizon of their network up to 60 second in [68] or to one month in [65] and achieve strong results on spatiotemporal localization. The works [62], [63], [64] accumulate extracted face features in a memory bank to preserve all the past knowledge without forgetting and at the same time handle the non-stationarity of the data stream. At a high level, all these approaches can be framed as a non-parametric estimation method (like nearest neighbors) sitting on top of a high-powered parametric function (Faster R-CNN in the case of object detection [65], a face feature extractor in [62] and [63], a SiamFC feature extractor [72] for object tracking in [71]). These methods use a fixed representation that is not incrementally learned as it would require re-encoding all the images in the memory bank. Avoiding re-encoding images can be advantageous in applications where images cannot be stored for privacy reasons (i.e. face recognition, applications in medical imaging, etc.). Clearly, also Multi-Object Tracking [73], [74] can benefit from memory based learning.

## VII. CONCLUSION

We have shown that a special set of fixed classifiers based on regular polytopes generates stationary features by maximally exploiting the available representation space. The proposed method is simple to implement and theoretically correct.

Experimental results confirm both the theoretical analysis and the generalization capability of the approach across a range of datasets, baselines and architectures. Our RePoNet solution improves and generalizes the concept of a fixed classifier, recently proposed in [1], to a larger class of fixed classifier models exploiting the inherent symmetry of regular polytopes in the feature space.

Our findings may have implications in all of those Deep Neural Network learning contexts in which a classifier must be robust against changes of the feature representation while learning as in incremental and continual learning settings.

## APPENDIX

### Computing the Angle Between Adjacent Classifier Weights

The angle between a vertex and its adjacent vertices in a regular polytope can be computed following the same mathematical formulation used to compute its dihedral angle. The dihedral angle of a regular  $d$ -Simplex is the acute angle formed by a pair of intersecting faces. In the case  $d = 2$  the dihedral angle is the angle at the vertex of an equilateral triangle, while in the case  $d = 3$  is the angle formed by the faces of the regular tetrahedron.

Because the dual polytope of a regular  $d$ -Simplex is also a regular  $d$ -Simplex, the angle  $\theta$  between pairs of vertices can be expressed as [75]:

$$\theta = \pi - \delta, \quad (13)$$

where  $\delta$  is the dihedral angle. Since the dihedral angle of a regular  $d$ -Simplex is known to be [75][39]:

$$\delta = \arccos\left(\frac{1}{d}\right), \quad (14)$$

substituting Eq. 14 into Eq. 13 we obtain:

$$\theta = \pi - \arccos\left(\frac{1}{d}\right)$$

which simplifies to:

$$\theta = \arccos\left(-\frac{1}{d}\right).$$

The Eq. above provides the value between any pair of vectors in a  $d$ -Simplex.

The calculation for the the  $d$ -Cube follows a similar argument. The dihedral angle of a regular  $d$ -Orthoplex is known to be  $\arccos((2-d)/d)$ . Since the  $d$ -Cube is the dual of the  $d$ -Orthoplex the angle defined by a vertex of a  $d$ -Cube and its adjacent vertices is:

$$\theta = \arccos\left(\frac{d-2}{d}\right). \quad (15)$$

## ACKNOWLEDGMENT

This work was partially supported by the European Commission under European Horizon 2020 Programme, grant number 951911 - AI4Media. This research was also partially supported by NVIDIA Corporation with the donation of Titan Xp GPUs, Leonardo Finmeccanica S.p.A., Italy. We thank Francesco Calabrò (Leonardo), Giuseppe Fiameni (Nvidia) and Marco Rorro (CINECA) for their support.

## REFERENCES

- [1] E. Hoffer, I. Hubara, and D. Soudry, "Fix your classifier: the marginal value of training the last weight layer," in *International Conference on Learning Representations (ICLR)*, 2018.
- [2] Y. LeCun, Y. Bengio, and G. Hinton, "Deep learning," *nature*, vol. 521, no. 7553, pp. 436–444, 2015.
- [3] J. Schmidhuber, "Deep learning in neural networks: An overview," *Neural networks*, vol. 61, pp. 85–117, 2015.
- [4] B. Zoph, V. Vasudevan, J. Shlens, and Q. V. Le, "Learning transferable architectures for scalable image recognition," in *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2018.
- [5] H. Touvron, A. Vedaldi, M. Douze, and H. Jégou, "Fixing the train-test resolution discrepancy," in *Advances in Neural Information Processing Systems (NeurIPS)*, 2019.
- [6] Q. Cao, L. Shen, W. Xie, O. M. Parkhi, and A. Zisserman, "Vggface2: A dataset for recognising faces across pose and age," in *13th IEEE International Conference on Automatic Face & Gesture Recognition, FG 2018, Xi'an, China, May 15-19, 2018*, 2018, pp. 67–74.
- [7] J. Deng, J. Guo, N. Xue, and S. Zafeiriou, "Arcface: Additive angular margin loss for deep face recognition," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 4690–4699.
- [8] M. D. Zeiler and R. Fergus, "Visualizing and understanding convolutional networks," in *European conference on computer vision*. Springer, 2014, pp. 818–833.
- [9] M. Lin, Q. Chen, and S. Yan, "Network in network," *International Conference on Learning Representations (ICLR)*, 2014.
- [10] 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.
- [11] 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.
- [12] S. Chopra, R. Hadsell, and Y. LeCun, "Learning a similarity metric discriminatively, with application to face verification," in *2005 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR 2005), 20-26 June 2005, San Diego, CA, USA*, 2005, pp. 539–546.
- [13] F. Schroff, D. Kalenichenko, and J. Philbin, "Facenet: A unified embedding for face recognition and clustering," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2015, pp. 815–823.
- [14] E. Hoffer and N. Ailon, "Deep metric learning using triplet network," in *International Workshop on Similarity-Based Pattern Recognition*. Springer, 2015, pp. 84–92.
- [15] S. Ioffe and C. Szegedy, "Batch normalization: Accelerating deep network training by reducing internal covariate shift," in *Proceedings of the 32nd International Conference on Machine Learning*, ser. Proceedings of Machine Learning Research, F. Bach and D. Blei, Eds., vol. 37. Lille, France: PMLR, 07–09 Jul 2015, pp. 448–456.
- [16] Z. Li and S. Arora, "An exponential learning rate schedule for deep learning," in *International Conference on Learning Representations*, 2019.
- [17] I. Goodfellow, Y. Bengio, A. Courville, and Y. Bengio, *Deep learning*. MIT Press, 2016, vol. 1.
- [18] F. Pernici, M. Bruni, C. Baecchi, and A. Del Bimbo, "Maximally compact and separated features with regular polytope networks," in *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR) Workshops*, June 2019.
- [19] F. Pernici, M. Bruni, C. Baecchi, F. Turchini, and A. D. Bimbo, "Class-incremental learning with pre-allocated fixed classifiers," in *25th International Conference on Pattern Recognition, ICPR 2020, Milan, Italy, January 10-15, 2021*. IEEE Computer Society, 2020.
- [20] M. Hardt and T. Ma, "Identity matters in deep learning," *International Conference on Learning Representations (ICLR)*, 2017.
- [21] A. Sablayrolles, M. Douze, C. Schmid, and H. Jégou, "Spreading vectors for similarity search," *International Conference on Learning Representations (ICLR)*, 2019.
- [22] W. Liu, Y. Wen, Z. Yu, and M. Yang, "Large-margin softmax loss for convolutional neural networks," in *ICML*, vol. 2, no. 3, 2016, p. 7.
- [23] R. Ranjan, C. D. Castillo, and R. Chellappa, "L2-constrained softmax loss for discriminative face verification," *arXiv preprint arXiv:1703.09507*, 2017.
- [24] W. Liu, Y. Wen, Z. Yu, M. Li, B. Raj, and L. Song, "Sphereface: Deep hypersphere embedding for face recognition," in *CVPR*, 2017.

- [25] W. Liu, Z. Liu, Z. Yu, B. Dai, R. Lin, Y. Wang, J. M. Rehg, and L. Song, "Decoupled networks," in *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018.
- [26] F. Wang, X. Xiang, J. Cheng, and A. L. Yuille, "Normface: 1 2 hypersphere embedding for face verification," in *Proceedings of the 2017 ACM on Multimedia Conference*. ACM, 2017, pp. 1041–1049.
- [27] Y. Liu, H. Li, and X. Wang, "Learning deep features via congenerous cosine loss for person recognition," *arXiv preprint: 1702.06890*, 2017.
- [28] M. Hasnat, J. Bohné, J. Milgram, S. Gentric, L. Chen *et al.*, "von mises-fisher mixture model-based deep learning: Application to face verification," *arXiv preprint arXiv:1706.04264*, 2017.
- [29] Y. Yuan, K. Yang, and C. Zhang, "Feature incay for representation regularization," *arXiv preprint arXiv:1705.10284*, 2017.
- [30] F. Wang, J. Cheng, W. Liu, and H. Liu, "Additive margin softmax for face verification," *IEEE Signal Processing Letters*, vol. 25, no. 7, pp. 926–930, 2018.
- [31] B. Chen, W. Deng, and H. Shen, "Virtual class enhanced discriminative embedding learning," in *Advances in Neural Information Processing Systems 31*, S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett, Eds. Curran Associates, Inc., 2018, pp. 1942–1952.
- [32] W. Liu, R. Lin, Z. Liu, L. Liu, Z. Yu, B. Dai, and L. Song, "Learning towards minimum hyperspherical energy," *NIPS*, 2018.
- [33] K. Zhao, J. Xu, and M.-M. Cheng, "Regularface: Deep face recognition via exclusive regularization," in *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2019.
- [34] Y. Duan, J. Lu, and J. Zhou, "Uniformface: Learning deep equidistributed representation for face recognition," in *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2019.
- [35] J. J. Thomson, "XXIV. On the structure of the atom: an investigation of the stability and periods of oscillation of a number of corpuscles arranged at equal intervals around the circumference of a circle; with application of the results to the theory of atomic structure," *The London, Edinburgh, and Dublin Philosophical Magazine and Journal of Science*, vol. 7, no. 39, pp. 237–265, 1904.
- [36] J. Batle, A. Bagdasaryan, M. Abdel-Aty, and S. Abdalla, "Generalized thomson problem in arbitrary dimensions and non-euclidean geometries," *Physica A: Statistical Mechanics and its Applications*, vol. 451, pp. 237–250, 2016.
- [37] P. M. L. Tammes, "On the origin of number and arrangement of the places of exit on the surface of pollen-grains," *Recueil des travaux botaniques néerlandais*, vol. 27, no. 1, pp. 1–84, 1930.
- [38] B. Bagchi, "How to stay away from each other in a spherical universe," *Resonance*, vol. 2, no. 9, pp. 18–26, 1997.
- [39] H. S. M. Coxeter, *Regular Polytopes*, ser. Macmillan mathematics paperbacks. Macmillan, 1963.
- [40] H. Xiao, K. Rasul, and R. Vollgraf, "Fashion-mnist: a novel image dataset for benchmarking machine learning algorithms," *arXiv preprint arXiv:1708.07747*, 2017.
- [41] G. Cohen, S. Afshar, J. Tapson, and A. van Schaik, "EMNIST: extending MNIST to handwritten letters," in *2017 International Joint Conference on Neural Networks, IJCNN 2017, Anchorage, AK, USA, May 14-19, 2017*, 2017, pp. 2921–2926.
- [42] A. Krizhevsky, "Learning multiple layers of features from tiny images," 2009.
- [43] 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.
- [44] J. Hu, L. Shen, and G. Sun, "Squeeze-and-excitation networks," in *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2018.
- [45] X. Li, W. Wang, X. Hu, and J. Yang, "Selective kernel networks," in *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2019.
- [46] M. Tan and Q. V. Le, "Efficientnet: Rethinking model scaling for convolutional neural networks," in *Proceedings of the 36th International Conference on Machine Learning, ICML 2019, 9-15 June 2019, Long Beach, California, USA*, ser. Proceedings of Machine Learning Research, K. Chaudhuri and R. Salakhutdinov, Eds., vol. 97. PMLR, 2019, pp. 6105–6114.
- [47] Y. Wen, K. Zhang, Z. Li, and Y. Qiao, "A discriminative feature learning approach for deep face recognition," in *European Conference on Computer Vision*. Springer, 2016, pp. 499–515.
- [48] Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner, "Gradient-based learning applied to document recognition," *Proceedings of the IEEE*, vol. 86, no. 11, pp. 2278–2324, 1998.
- [49] K. He, X. Zhang, S. Ren, and J. Sun, "Delving deep into rectifiers: Surpassing human-level performance on imagenet classification," in *Proceedings of the IEEE international conference on computer vision*, 2015, pp. 1026–1034.
- [50] K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image recognition," *ICLR*, 2015.
- [51] J. Hu, L. Shen, and G. Sun, "Squeeze-and-excitation networks," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 7132–7141.
- [52] G. Huang, Z. Liu, L. van der Maaten, and K. Q. Weinberger, "Densely connected convolutional networks," in *2017 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2017, Honolulu, HI, USA, July 21-26, 2017*, 2017, pp. 2261–2269.
- [53] X. Li, W. Wang, X. Hu, and J. Yang, "Selective kernel networks," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2019, pp. 510–519.
- [54] S. Xie, R. Girshick, P. Dollár, Z. Tu, and K. He, "Aggregated residual transformations for deep neural networks," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 1492–1500.
- [55] E. D. Cubuk, B. Zoph, J. Shlens, and Q. V. Le, "Randaugment: Practical automated data augmentation with a reduced search space," in *The IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) Workshops*, June 2020.
- [56] D. Justus, J. Brennan, S. Bonner, and A. S. McGough, "Predicting the computational cost of deep learning models," in *IEEE International Conference on Big Data, Big Data 2018, Seattle, WA, USA, December 10-13, 2018*, N. Abe, H. Liu, C. Pu, X. Hu, N. K. Ahmed, M. Qiao, Y. Song, D. Kossmann, B. Liu, K. Lee, J. Tang, J. He, and J. S. Saltz, Eds. IEEE, 2018, pp. 3873–3882.
- [57] A. Graves, G. Wayne, and I. Danihelka, "Neural turing machines," *CoRR*, vol. abs/1410.5401, 2014.
- [58] J. Weston, S. Chopra, and A. Bordes, "Memory networks," in *3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings*, Y. Bengio and Y. LeCun, Eds., 2015.
- [59] A. Graves, G. Wayne, M. Reynolds, T. Harley, I. Danihelka, A. Grabska-Barwińska, S. G. Colmenarejo, E. Grefenstette, T. Ramalho, J. Agapiou *et al.*, "Hybrid computing using a neural network with dynamic external memory," *Nature*, vol. 538, no. 7626, pp. 471–476, 2016.
- [60] L. Kaiser, O. Nachum, A. Roy, and S. Bengio, "Learning to remember rare events," in *5th International Conference on Learning Representations, ICLR 2017, Toulon, France, April 24-26, 2017, Conference Track Proceedings*. OpenReview.net, 2017.
- [61] Y. Shen, Y. Xiong, W. Xia, and S. Soatto, "Towards backward-compatible representation learning," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2020.
- [62] F. Pernici, M. Bruni, and A. Del Bimbo, "Self-supervised on-line cumulative learning from video streams," *Computer Vision and Image Understanding*, vol. 197-198, p. 102983, 2020.
- [63] F. Pernici, F. Bartoli, M. Bruni, and A. Del Bimbo, "Memory based on-line learning of deep representations from video streams," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2018.
- [64] F. Pernici and A. Del Bimbo, "Unsupervised incremental learning of deep descriptors from video streams," in *2017 IEEE International Conference on Multimedia & Expo Workshops (ICMEW)*. IEEE, 2017, pp. 477–482.
- [65] S. Beery, G. Wu, V. Rathod, R. Votel, and J. Huang, "Context rnn: Long term temporal context for per-camera object detection," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2020.
- [66] H. Deng, Y. Hua, T. Song, Z. Zhang, Z. Xue, R. Ma, N. M. Robertson, and H. Guan, "Object guided external memory network for video object detection," in *2019 IEEE/CVF International Conference on Computer Vision, ICCV 2019, Seoul, Korea (South), October 27 - November 2, 2019*. IEEE, 2019, pp. 6677–6686.
- [67] M. Shvets, W. Liu, and A. C. Berg, "Leveraging long-range temporal relationships between proposals for video object detection," in *2019 IEEE/CVF International Conference on Computer Vision, ICCV 2019, Seoul, Korea (South), October 27 - November 2, 2019*. IEEE, 2019, pp. 9755–9763.
- [68] C. Wu, C. Feichtenhofer, H. Fan, K. He, P. Krähenbühl, and R. B. Girshick, "Long-term feature banks for detailed video understanding," in *IEEE Conference on Computer Vision and Pattern Recognition, CVPR*

2019, Long Beach, CA, USA, June 16-20, 2019. Computer Vision Foundation / IEEE, 2019, pp. 284–293.

- [69] H. Wu, Y. Chen, N. Wang, and Z. Zhang, “Sequence level semantics aggregation for video object detection,” in *2019 IEEE/CVF International Conference on Computer Vision, ICCV 2019, Seoul, Korea (South), October 27 - November 2, 2019*. IEEE, 2019, pp. 9216–9224.
- [70] S. W. Oh, J. Lee, N. Xu, and S. J. Kim, “Video object segmentation using space-time memory networks,” in *2019 IEEE/CVF International Conference on Computer Vision, ICCV 2019, Seoul, Korea (South), October 27 - November 2, 2019*. IEEE, 2019, pp. 9225–9234.
- [71] T. Yang and A. B. Chan, “Learning dynamic memory networks for object tracking,” in *Computer Vision - ECCV 2018 - 15th European Conference, Munich, Germany, September 8-14, 2018, Proceedings, Part IX*, ser. Lecture Notes in Computer Science, V. Ferrari, M. Hebert, C. Sminchisescu, and Y. Weiss, Eds., vol. 11213. Springer, 2018, pp. 153–169.
- [72] L. Bertinetto, J. Valmadre, J. F. Henriques, A. Vedaldi, and P. H. S. Torr, “Fully-convolutional siamese networks for object tracking,” in *Computer Vision - ECCV 2016 Workshops - Amsterdam, The Netherlands, October 8-10 and 15-16, 2016, Proceedings, Part II*, ser. Lecture Notes in Computer Science, G. Hua and H. Jégou, Eds., vol. 9914, 2016, pp. 850–865.
- [73] G. Ciaparrone, F. Luque Sánchez, S. Tabik, L. Troiano, R. Tagliaferri, and F. Herrera, “Deep learning in video multi-object tracking: A survey,” *Neurocomputing*, vol. 381, pp. 61 – 88, 2020.
- [74] P. Salvagnini, F. Pernici, M. Cristani, G. Lisanti, I. Masi, A. Del Bimbo, and V. Murino, “Information theoretic sensor management for multi-target tracking with a single pan-tilt-zoom camera,” in *IEEE Winter Conference on Applications of Computer Vision*. IEEE, 2014, pp. 893–900.
- [75] H. R. Parks and D. C. Wills, “An elementary calculation of the dihedral angle of the regular  $n$ -simplex,” *The American mathematical monthly*, vol. 109, no. 8, pp. 756–758, 2002.



**Federico Pernici** Federico Pernici received the laurea degree in Information Engineering in 2002, the post-laurea degree in Internet Engineering in 2003 and the Ph.D. in Information and Telecommunication Engineering in 2005 from the University of Firenze, Italy. Since 2002 he has been a research assistant at MICC Media Integration and Communication Center, assistant professor and adjunct professor at the University of Firenze. His scientific interests are computer vision and machine learning with a focus on different aspects of visual tracking, incremental learning and representation learning. Presently, he is Associate Editor of *Machine Vision and Applications* journal.



**Matteo Bruni** Matteo Bruni received the M.S. degree (cum laude) in Computer Engineering from the University of Firenze, Italy, in 2016. Presently, he is a Ph.D. student at the University of Firenze at MICC, Media Integration and Communication Center, University of Firenze. His research interests include pattern recognition and computer vision with specific focus on feature embedding, face recognition and incremental learning.



**Claudio Baccchi** Claudio Baccchi received his Master's Degree in computer engineering from University of Florence in 2013 and the Ph.D in 2017. Currently he is working at the Visual Information and Media Lab at Media Integration and Communication Centre, University of Florence. His current researches in the computer vision field cover sentiment and polarity classification in web videos, face and emotion recognition.



**Alberto Del Bimbo** Prof. Del Bimbo is Full Professor at the University of Firenze, Italy and the Director of MICC Media Integration and Communication Center. He is the author of over 350 scientific publications in computer vision and multimedia and principal investigator of technology transfer projects with industry and governments. He was the Program Chair of ICPR 2012, ICPR 2016 and ACM Multimedia 2008, and the General Chair of IEEE ICMCS 1999, ACM Multimedia 2010, ICMR 2011 and ECCV 2012. He is the General Chair of the forthcoming ICPR 2020. He is the Editor in Chief of *ACM TOMM Transactions on Multimedia Computing Communications and Applications* and Associate Editor of *Multimedia Tools and Applications* and *Pattern Analysis and Applications* journals. He was Associate Editor of *IEEE Transactions on Pattern Analysis and Machine Intelligence*, *IEEE Transactions on Multimedia* and *Pattern Recognition* and also served as the Guest Editor of many Special Issues in highly ranked journals. Prof. Del Bimbo is IAPR Fellow and ACM Distinguished Scientist and is the recipient of the 2016 ACM SIGMM Award for *Outstanding Technical Contributions to Multimedia Computing Communications and Applications*.