

Deep Learning on 3D Neural Fields

Pierluigi Zama Ramirez* , Luca De Luigi* , Daniele Sirocchi*

Adriano Cardace, Riccardo Spezialetti, Francesco Ballerini, Samuele Salti, Luigi Di Stefano

University of Bologna

{pierluigi.zama, luca.deluigi4, adriano.cardace2, riccardo.spezialetti}@unibo.it

Abstract—In recent years, Neural Fields (*NFs*) have emerged as an effective tool for encoding diverse continuous signals such as images, videos, audio, and 3D shapes. When applied to 3D data, *NFs* offer a solution to the fragmentation and limitations associated with prevalent discrete representations. However, given that *NFs* are essentially neural networks, it remains unclear whether and how they can be seamlessly integrated into deep learning pipelines for solving downstream tasks. This paper addresses this research problem and introduces *nf2vec*, a framework capable of generating a compact latent representation for an input *NF* in a single inference pass. We demonstrate that *nf2vec* effectively embeds 3D objects represented by the input *NFs* and showcase how the resulting embeddings can be employed in deep learning pipelines to successfully address various tasks, all while processing exclusively *NFs*. We test this framework on several *NFs* used to represent 3D surfaces, such as unsigned/signed distance and occupancy fields. Moreover, we demonstrate the effectiveness of our approach with more complex *NFs* that encompass both geometry and appearance of 3D objects such as neural radiance fields.

Index Terms—Neural Fields, INR, Implicit Neural Representations, Representation Learning, Deep Learning on Neural Fields, Signed Distance Function, SDF, Unsigned Distance Function, UDF, Occupancy Field, OF, Neural Radiance Field, NeRF, 3D classification, 3D generation, 3D segmentation, 3D completion, 3D reconstruction, NeRF generation, NeRF classification.

I. INTRODUCTION

Computer vision has always been concerned with understanding the 3D world around us. One of the main challenges when dealing with 3D data is the representation strategy, which was addressed over the years by introducing various discrete representations, including voxel grids, point clouds, and meshes. Each representation has its advantages and disadvantages, especially when it comes to processing it through deep learning, leading to the development of a plethora of ad-hoc algorithms [1]–[3] for each coexisting representation. Hence, no standard way to store and process 3D data has yet emerged.

Recently, a new representation has been proposed, called Neural Fields [4] (*NFs*). They are continuous functions defined at all spatial coordinates, parameterized by a neural network such as a Multi-Layer Perceptron (MLP). In the context of 3D world representation, various types of *NFs* have been explored. Some of the most common *NFs* utilize the Signed/Unsigned Distance Field (*SDF/UDF*) [5]–[8] and the Occupancy Field (*OF*) [9], [10] to represent the 3D surfaces or volumes of the objects in the scene. Alternatively, strategies seeking to capture both geometries and appearances often leverage the Radiance Field (*RF*), as shown in the pioneering approach NeRF [11].

Representing a 3D scene by encoding it with a continuous function parameterized as an MLP separates the memory cost of the representation from the spatial resolution. In other words, starting from the same fixed number of parameters, it is possible to reconstruct a surface with arbitrarily fine resolution or to render an image with arbitrarily high quality. Furthermore, the identical neural network architecture can be applied to learn various field functions, offering the possibility of a unified framework for representing 3D objects.

Owing to their efficacy and potential benefits, 3D *NFs* are garnering growing interest from the scientific community, as evidenced by the frequent publication of novel and impressive results [8], [12]–[14]. This leads us to speculate that, in the near future, *NFs* could establish themselves as a standard way to store and communicate 3D data. It is conceivable that repositories hosting digital twins of 3D objects, exclusively realized as MLPs, might become widely accessible.

The above scenario prompts an intriguing research question: can 3D *NFs* be directly processed using deep learning pipelines for solving downstream tasks, as it is commonly done with discrete representations such as point clouds or images? For instance, is it feasible to classify an object by directly processing the corresponding NeRF without rendering any image from it?

Since *NFs* are neural networks, there is no straightforward way to process them. A recent work in the field, Functua [15], fits the whole dataset with a shared network conditioned on a different embedding for each data. In this formulation, a solution could be to use such embeddings as the input for downstream tasks. Nevertheless, representing an entire dataset through a shared network poses a formidable learning challenge, as the network encounters difficulties in accurately fitting all the samples (see Section VII).

On the contrary, recent studies, including SIREN [16] and others [17]–[21], have demonstrated that it is possible to achieve high-quality reconstructions by tailoring an individual network to each input sample. This holds true even when dealing with complex 3D shapes or images. Furthermore, constructing an individual *NF* for each object is more adaptable to real-world deployment, as it does not require the availability of the entire dataset to fit each individual data. The increasing popularity of such methodologies suggests that adopting the practice of fitting an individual network is likely to become commonplace in learning *NFs*.

Therefore, in the former version of this paper [22], we explored conducting downstream tasks using deep learning pipelines on 3D data represented as **individual** *NFs*. Recently, several methods addressing this topic have been published,

*Joint first authorship.

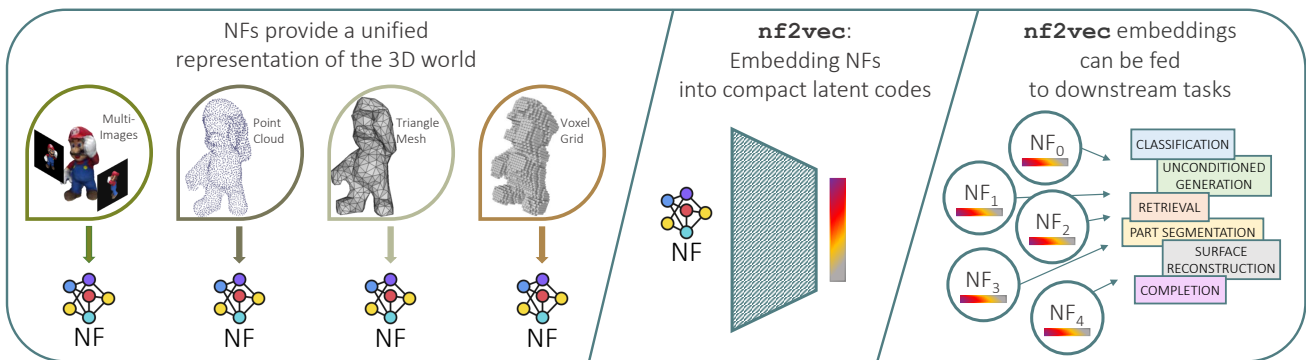


Fig. 1: **Overview of our framework.** **Left:** NFs hold the potential to provide a unified representation of the 3D world. **Center:** Our framework, dubbed $nf2vec$, produces a compact representation for an input NF by looking only at its weights. **Right:** $nf2vec$ embeddings can be used with standard deep-learning machinery to solve various downstream tasks.

such as NFN [23], NFT [23], and DWSNet [24], and all of them process individual NFs , supporting this paradigm.

Using NFs as input or output data is intrinsically non-trivial, as the MLP of a single NF can encompass hundreds of thousands of parameters. However, deep models inherently present a significantly redundant parameterization of the underlying function, as shown in [25], [26]. As a result, we explore whether and how an answer to the research question mentioned earlier might be identified within a representation learning framework. We present an approach that encodes individual NFs into compact and meaningful embeddings, making them suitable for diverse downstream tasks. We name this framework $nf2vec$, shown in Fig. 1.

Our framework has at its core an encoder designed to produce a task-agnostic embedding representing the input NF by processing only the NF weights. These embeddings can seamlessly be used in downstream deep learning pipelines as we validate for various tasks, like classification, retrieval, part segmentation, unconditioned generation, completion, and surface reconstruction. Remarkably, the last two tasks become achievable by learning a straightforward mapping between the embeddings generated using our framework, as embeddings derived from NFs exist in low-dimensional vector spaces, regardless of the underlying implicit function. For instance, we can learn the mapping between NFs of incomplete objects into NFs of normal ones. Then, we can complete shapes by exploiting this mapping, *e.g.*, we can map the NF of an airplane with a missing wing into the NF of a complete airplane. Furthermore, we show that $nf2vec$ learns a smooth latent space, which enables the interpolation of NFs representing previously unseen 3D objects.

This paper builds on our previous work [22], with revisions to the overall framework and thorough experiments on novel scenarios. Specifically, the key differences with [22] are:

- In [22], we focused solely on neural fields representing the surfaces of 3D objects. In this extended version, we also tackle the processing of neural fields capturing objects’ geometry and appearance. Specifically, we extend our framework to perform deep learning tasks on NeRFs by directly processing their MLPs weights.
- The processing of MLPs parametrizing NFs has been

investigated in works published contemporaneously or subsequently to [22]. We extend our literature review by including these recent papers, and we evaluate them to foster progress in this emerging topic and facilitate future comparisons.

Overall, the summary of our work contributions is:

- We propose and investigate the novel research problem of applying deep learning directly on individual NFs representing 3D objects.
- We introduce $nf2vec$, a framework designed to derive a meaningful and compact representation of an input NF solely by processing its weights, without the need to sample the underlying function.
- We demonstrate that a range of tasks, typically tackled with intricate frameworks tailored to specific representations, can be effectively executed using simple deep learning tools on NFs embedded by $nf2vec$, regardless of the signal underlying the NFs .
- We demonstrate the versatility of $nf2vec$ by successfully applying it to neural fields that capture either the geometry alone or the combined information of both geometry and appearance of 3D objects.
- We analyze recent methods for processing NFs in terms of classification accuracy and representation quality. We build the first evaluation benchmark for NF classification.

Additional details, code, and datasets are available at <https://cvlab-unibo.github.io/nf2vec>.

II. RELATED WORK

Neural fields. Recent approaches have shown the ability of MLPs to parameterize fields representing any physical quantity of interest [4]. The works focusing on representing 3D shapes with MLPs rely on fitting functions such as the unsigned distance [6], the signed distance [5], [7], [10], [27]–[29], or the occupancy [9], [30]. Among these approaches, sinusoidal representation networks (SIREN) [16] use periodical activation functions to capture the high-frequency details of the input data. In addition to representing shapes, some of these models have been extended to encode object appearance [11], [27], [31]–[33], or to include temporal information [34]. Among these

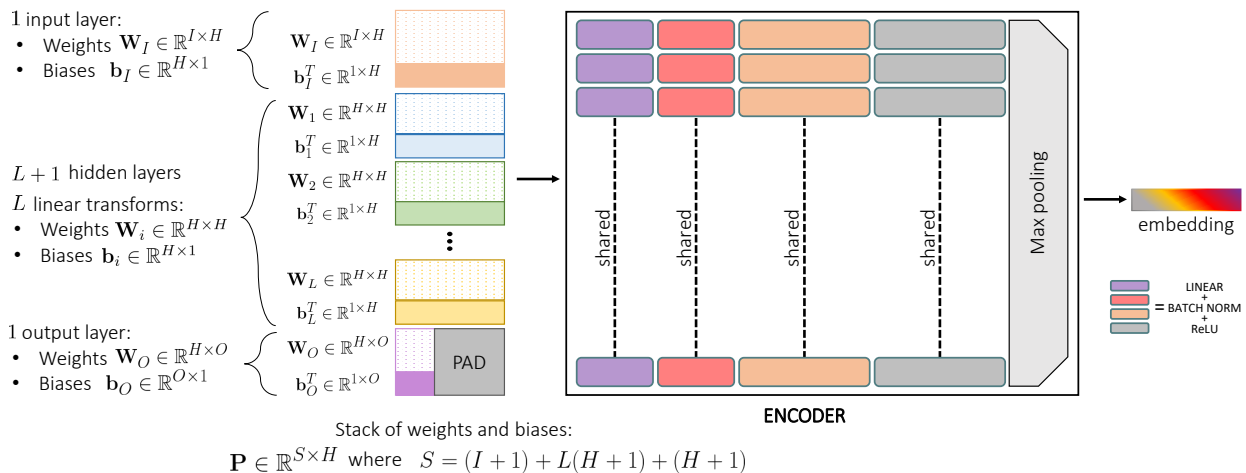


Fig. 2: Encoder architecture.

NF hidden dim.	512	512	512	1024	1024	1024
NF #layers	4	8	16	4	8	16
NF #params	~800K	~2M	~4M	~3M	~7M	~15M
#params our encoder	~3M	~3M	~3M	~3.5M	~3.5M	~3.5M
#params MLP encoder	~800M	~2B	~4B	~3B	~7.5B	~16B

TABLE I: Num. of parameters of our vs an MLP encoder.

recent approaches, modeling the radiance field of a scene [11] has proven to be the critical factor in obtaining excellent scene representations. In our work, we employ *NFs* encoding *SDF*, *UDF*, *OF*, and *RF* as input data for deep learning pipelines.

Deep learning on neural networks. Several works have explored using neural networks to process other neural networks. [35] utilizes a network’s weights as input and forecasts its classification accuracy. Another approach [36] involves learning a network representation through a self-supervised learning strategy applied to the *N*-dimensional weight array. These representations are then employed to predict various characteristics of the input classifier. In contrast, [37]–[39] depict neural networks as computational graphs, subsequently processed by a Graph Neural Network (GNN). This GNN is tasked with predicting optimal parameters, adversarial examples, or branching strategies for verifying neural networks.

These works view neural networks as algorithms, primarily focusing on forecasting properties like accuracy. In contrast, some recent studies handle networks that implicitly represent 3D data, thus tackling various tasks directly from their weights, essentially treating neural networks as input/output data. Functa [15] tackles this scenario by acquiring priors across the entire dataset using a shared network and subsequently encoding each sample into a concise embedding employed for downstream discriminative and generative tasks. We note that in this formulation, each neural field is parametrized by both the shared network and the embedding. It is worth pointing out that, though not originally proposed as a framework to process neural fields, DeepSDF [5] learns dataset priors by optimizing a reconstruction objective through a shared auto-decoder network conditioned on a shape-specific embedding. Thus, the embeddings learned by DeepSDF may be used for neural processing tasks, as done in Functa.

However, shared network frameworks encounter challenges

as they struggle to reconstruct the underlying signal with high fidelity and require an entire dataset to learn the neural field of an object. In response, recent approaches have shifted their focus on processing *NFs* learned on individual data, *e.g.*, a specific object or scene. The first framework doing it was proposed in our previous paper version [22]. This approach leverages representation learning to condense individual neural fields of 3D shapes into embeddings, serving as input for subsequent tasks. Recognizing that MLPs exhibit weight space symmetries [40], where hidden neurons can be permuted across layers without altering the network’s function, recent approaches such as DWSNet [24], NFN [41], and NFT [23] leverage these symmetries as an inductive bias to create innovative architectures tailored for MLPs. DWSNet and NFN design neural layers equivariant to the permutations inherent in MLPs. In contrast, NFT builds on the concept of achieving permutation equivariance by eliminating positional encoding from a Transformer architecture.

III. LEARNING TO REPRESENT *NFs*

This paper explores the possibility and the methodology of directly utilizing *NFs* for downstream tasks. Specifically, can we classify an object implicitly encoded in a *NF*, and if so, how? As outlined in Section I, we condense the redundant information encoded in the weights of *NFs* into latent codes by a representation learning framework. These codes can then be efficiently processed using standard deep-learning pipelines. Our framework, dubbed `nf2vec`, comprises an encoder and a decoder. In the following sections, we first provide some basic knowledge about what a 3D neural field is, then we deepen the reasons behind the architectural choices for both components and describe the representation learning protocol.

3D Neural Fields. A field is a physical quantity defined for all domain coordinates. We focus on fields describing the 3D world, thus operating on \mathbb{R}^3 coordinates $\mathbf{p} = (x, y, z)$. We consider the 3D fields commonly used in computer vision and graphics, such as the *SDF* [5] and *UDF* [6], which map coordinates to the signed and unsigned distance from the closest surface, the *OF* [9], which computes the occupancy probability of each position, and the *RF* [11], that outputs (r, g, b) colors

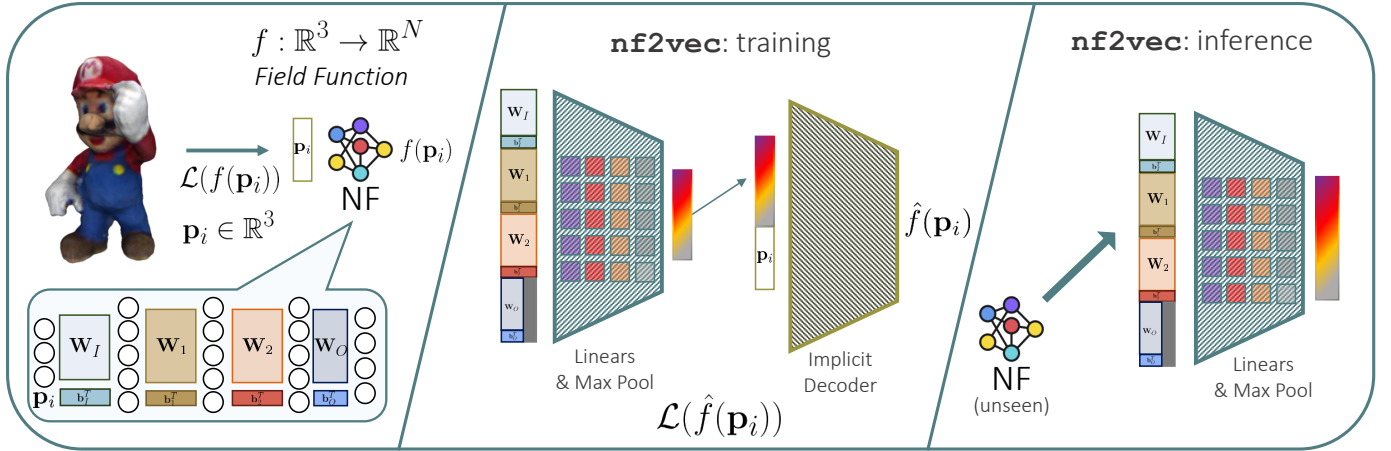


Fig. 3: **Training and inference of `nf2vec`.** **Left:** A Neural Field (NF) represents a 3D object. The NF is composed of an MLP that parametrizes a field function f . **Center:** `nf2vec` encoder is trained together with an implicit decoder. The implicit decoder processes the embedding produced by the encoder to estimate field values \hat{f} . We train the framework similarly to the input NF . **Right:** At inference time, the learned encoder can be used to obtain a compact embedding from unseen NF s.

and density σ for each 3D point. A field can be modeled by a function, f , parameterized by θ . Thus, for any point \mathbf{p} , the field is given by $f(\mathbf{p})$. If parameters θ are the weights of a neural network, f is said to be a Neural Field (NF) [4].

Encoder. The encoder takes as input the weights of a NF and produces a compact embedding that encodes all the relevant information of the input NF . Designing an encoder for NF s poses a challenge in handling weights efficiently to avoid excessive memory usage. While a straightforward solution might involve using an MLP encoder to map flattened weight vectors to desired dimensions, this approach becomes impractical for larger NF s. For instance, given a 4-layer 512-neurons NF , mapping its 800K parameters to a 1024-sized embedding space would require an encoder with roughly 800M parameters, making this approach prohibitive. Thus, we focus on developing an encoder architecture that scales gracefully with the size of the input NF .

Following conventional practice [16]–[20], we consider NF s composed of an MLP with several hidden layers, each with H nodes. The linear transformation between two consecutive hidden layers is parameterized by a matrix of weights $\mathbf{W}_i \in \mathbb{R}^{H \times H}$ and a vector of biases $\mathbf{b}_i \in \mathbb{R}^{H \times 1}$. Thus, stacking \mathbf{W}_i and \mathbf{b}_i^T , the mapping between two consecutive layers can be represented by a single matrix $\mathbf{P}_i \in \mathbb{R}^{(H+1) \times H}$. Additionally, an MLP features also an input layer, parameterized by $\mathbf{W}_I \in \mathbb{R}^{I \times H}$ and a vector of biases $\mathbf{b}_I \in \mathbb{R}^{H \times 1}$, and an output layer, parameterized by $\mathbf{W}_O \in \mathbb{R}^{H \times O}$ and a vector of biases $\mathbf{b}_O \in \mathbb{R}^{O \times 1}$. The input and output layers can be represented by two matrices, $\mathbf{P}_I \in \mathbb{R}^{(I+1) \times H}$ and $\mathbf{P}_O \in \mathbb{R}^{(H+1) \times O}$. Considering that for a NF with $L + 1$ hidden layers there are L linear transformations between them, we can store all the weights of the NF by stacking the input matrix \mathbf{P}_I with all the L matrices \mathbf{P}_i and with the output matrix \mathbf{P}_O . As \mathbf{P}_O has a different number of columns (O), we pad it with zeros before stacking. The final stacked matrix \mathbf{P} has dimension $S \times H$ – where $S = (I + 1) + L(H + 1) + (H + 1)$ – and represents the input for our encoder, as shown in the left part of Fig. 2.

`nf2vec` encoder consists of a series of linear layers with batch normalization and ReLU non-linearity followed by final max pooling. At each stage, the input matrix \mathbf{P} is transformed by one linear layer, processing each row independently. The final max pooling compresses all the rows into a single one, obtaining the desired embedding. An architecture overview is depicted in Fig. 2.

Our proposed architecture scales gracefully to bigger input NF s as supported by the analysis in Table I, that reports the parameters of our encoder *w.r.t.* those of a generic MLP encoder while varying the input NF dimension.

It is worth observing that the randomness involved in fitting an individual NF (weights initialization, data shuffling, etc.) causes the weights in the same position in the NF architecture not to share the same role across NF s. Thus, `nf2vec` encoder would have to deal with input vectors whose elements capture different information across the different data samples, making it impossible to train the framework. However, the use of a shared, pre-computed initialization has been advocated as a good practice when fitting NF s, *e.g.*, to reduce training time by means of meta-learned initialization vectors, as done in MetaSDF [17] and in Functua [15], or to obtain desirable geometric properties [7]. We empirically found that following such a practice, *i.e.*, initializing all NF s with the same random vector, favors the alignment of weights across NF s and enables the convergence of our framework. More analysis can be found in Section VIII.

Decoder. When learning to encode NF s, we are interested in storing the information about the represented object rather than the values of the input weights. Therefore, the adopted decoder predicts the original field values rather than reconstructing the input weights in an auto-encoder fashion. In particular, during training, we adopt an implicit decoder inspired by [5], which takes in input the embeddings produced by the encoder and a spatial coordinate \mathbf{p}_i and decodes the original field values (see Fig. 3 center). We denote as $\hat{f}(\mathbf{p}_i)$ the predicted field value.

Training. We train our encoder and decoder following the

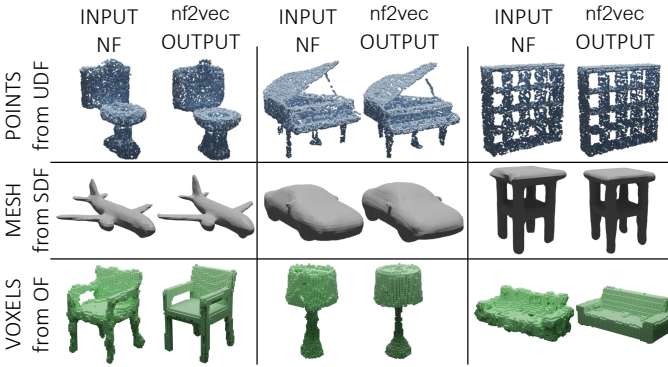


Fig. 4: `nf2vec` reconstructions.

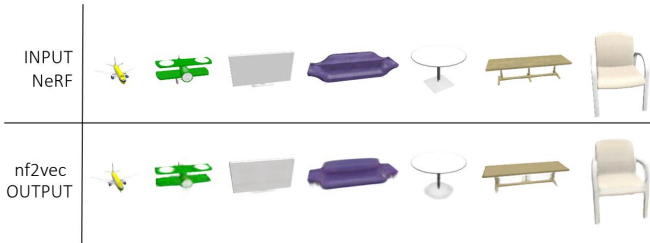


Fig. 5: `nf2vec` reconstructions.

input *NFs* training strategy. For instance, when dealing with *UDF*, *SDF*, and *OF* representing 3D surfaces, we supervise the framework directly using the ground truth field values computed from point clouds, voxel grids, or triangle meshes representing those surfaces. Differently, when processing NeRFs we employ volumetric rendering [11] on the radiance field values predicted by the decoder to obtain the RGB intensities of image pixels, and we supervise the framework directly with a regression loss between predicted and true RGB values.

To better understand the procedure, let us take the example where we aim to learn to represent *UDFs*. We create a set of 3D queries paired with the values of the *UDF* at those locations. The decoder takes in input the embedding produced by the encoder concatenated with the 3D coordinates of a query point and estimates the *UDF* for this location. The whole encoder-decoder is supervised to minimize the discrepancy between the estimated and correct *UDF* values.

Inference. After the overall framework has been trained end to end, the frozen encoder can be used to compute embeddings of unseen *NFs* with a single forward pass (see Fig. 3 right) while the implicit decoder can be used, if needed, to reconstruct the discrete representation given an embedding. *We highlight that no discrete representations are required at inference time.*

The presented `nf2vec` framework shares the same high-level structure of the originally proposed approach `inr2vec` [22]. However, as shown in Section VI, in this extended work, we show that our framework can be applied to more complex neural fields such as NeRFs. Thus, we dub it `nf2vec` to emphasize its generality.

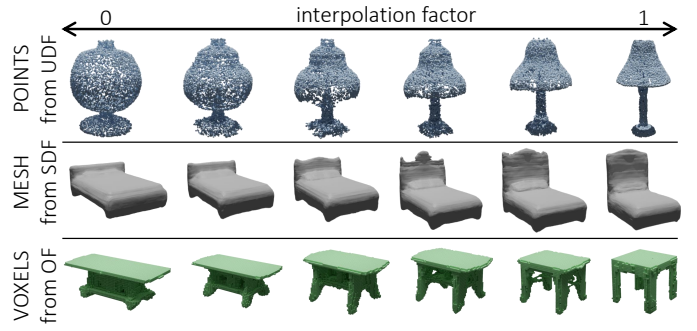


Fig. 6: `nf2vec` latent space interpolation.

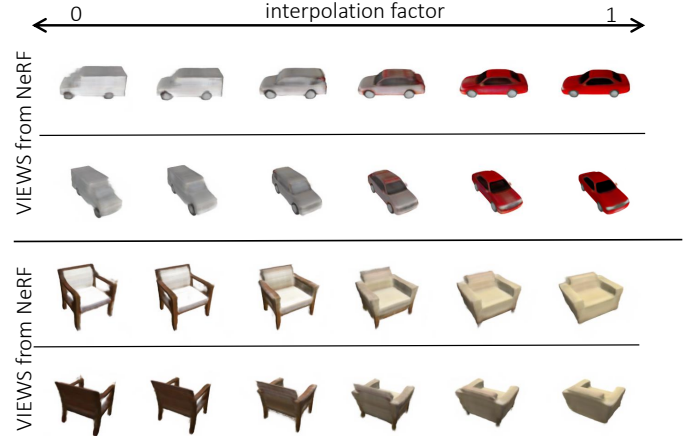


Fig. 7: `nf2vec` latent space interpolation.

IV. LATENT SPACE PROPERTIES

We train `nf2vec` on *NFs* learned from various 3D discrete representations of ShapeNet [42] surfaces. We use either *UDF* learned from point clouds, *SDF* learned from meshes, or *OF* learned from voxel grids. We also train `nf2vec` on NeRFs trained on multi-view renderings of ShapeNet [42] objects. The following sections explore `nf2vec` latent space properties. We investigate the ability to reconstruct the original discrete representation from the compact latent codes extracted with our encoder, the smoothness of the embedding space, and its structure.

Reconstruction. In Fig. 4, we compare 3D shapes reconstructed from *NFs* unseen during training with those reconstructed by the `nf2vec` decoder starting from the latent codes yielded by the encoder. We visualize point clouds with 8192 points, meshes reconstructed by marching cubes [43] from a grid with resolution 128^3 and voxels with resolution 64^3 . Moreover, we conduct the same experiment using as input unseen NeRFs. In Fig. 5, we show the images rendered by `nf2vec` decoder and those rendered from the input NeRF at 224×224 resolution. Though our embedding is dramatically more compact than the original *NF*, the reconstructed discrete data resembles those of the original input *NF*.

Interpolation. In Fig. 6 and Fig. 7, we linearly interpolate between two objects embeddings produced by `nf2vec`. Results highlight that the learned latent spaces enable smooth interpolations between shapes represented as *NFs*. Notably, in Fig. 7, color and shapes change smoothly when interpolating two NeRF embeddings. Moreover, the *interpolated* object has

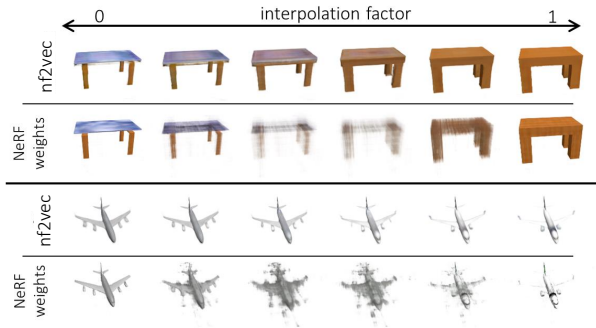


Fig. 8: *nf2vec* latent space interpolation vs NeRF weights interpolation.

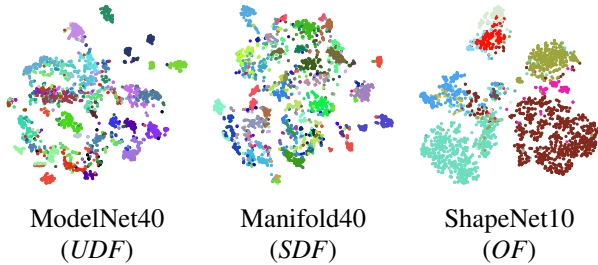


Fig. 9: **t-SNE visualizations of *nf2vec* latent spaces.** We plot the t-SNE components of the embeddings produced by *nf2vec* on the test sets of three datasets, ModelNet40 (left), Manifold40 (center) and Shapenet10 (right). Colors represent the different classes of the datasets.

an underlying 3D consistency, as visible when rendering it from different camera viewpoints. Thus, these results demonstrate the capability to generate new plausible shapes and even realistic *RFs* by means of simple linear interpolation in *nf2vec* latent space.

Additionally, given two input NeRFs, we render images from networks obtained by interpolating their weights. In Fig. 8, we compare these results with those obtained from the interpolation of *nf2vec* embeddings. Notably, our interpolations exhibit superior quality. In particular, renderings obtained by averaging the weights of the two NeRFs (interpolation factor 0.5) appear blurred and lack 3D structure. In contrast, renderings produced by *nf2vec* preserve details and maintain 3D consistency. Our results highlight the efficacy of our representation learning approach in transforming an initially disorganized input weight space into a well-organized latent space.

t-SNE visualization of the latent space. In Fig. 9, we provide the t-SNE visualization of the embeddings produced by *nf2vec* when presented with unseen *NF* of three different datasets: ModelNet40 (*UDF* learned from point clouds), Manifold 40 (*SDF* learned from meshes), and ShapeNet10 (*OF* learned from voxel grids). During the training of our framework, the supervisory signal employed does not impose any particular constraints on the organization of the learned latent space. This lack of constraints was deliberate, as it was not necessary for our primary goal – performing downstream tasks with the generated embeddings. Nevertheless, it is intriguing to note from the t-SNE plots that our algorithm naturally arranges

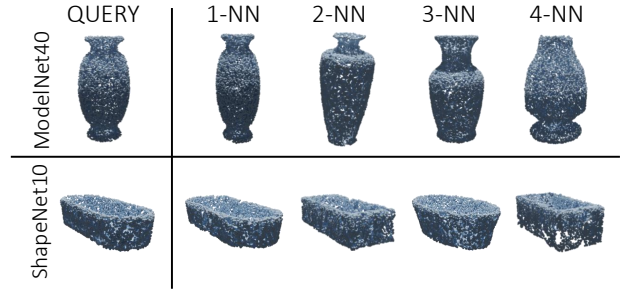


Fig. 10: **Point cloud retrieval qualitative results.** Given the *nf2vec* embedding of a query shape, we show the shapes reconstructed from the closest embeddings (L2 distance).

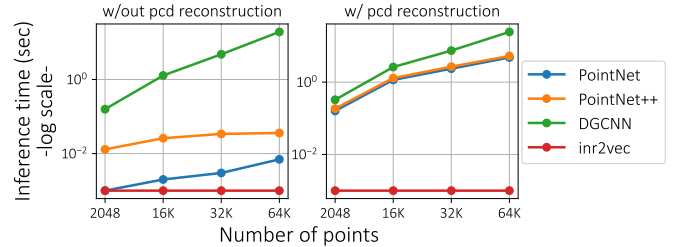


Fig. 11: **Time required to classify *NFs* encoding udf.** We plot the inference time of standard baselines and of our method, both considering the case in which discrete point clouds are available (left) and the one where point clouds must be reconstructed from the input *NFs* (right).

the embeddings in the latent space with a semantic structure, with items of the same category consistently mapped to close positions. This is evident in the colors representing different classes within the datasets under consideration.

V. DEEP LEARNING ON 3D SHAPES

This section shows how several tasks dealing with 3D shapes can be tackled by working only with *nf2vec* embeddings as input and/or output.

General settings. In all the experiments reported in this section, we convert 3D discrete representations into *NFs* featuring 4 hidden layers with 512 nodes each, using the SIREN activation function [16]. We discard the input and output layers of SIREN MLPs when processing them with *nf2vec*. This is based on the observation from the earlier version of this paper [22] that these layers do not provide information beneficial for downstream tasks when using SIREN MLPs as architecture for *NF*. We train *nf2vec* using an encoder composed of four linear layers with respectively 512, 512, 1024, and 1024 features, embeddings with 1024 values, and an implicit decoder with 5 hidden layers with 512 features. In all the experiments, the baselines are trained using standard data augmentation (random scaling and point-wise jittering), while we train both *nf2vec* and the downstream task-specific networks on datasets augmented offline with the same transformations.

Point cloud retrieval. We examine the potential of using *nf2vec* embeddings for representation learning tasks, with 3D retrieval as our benchmark. We follow the procedure introduced in [42], using the Euclidean distance to measure the similarity

Method	Input	ModelNet40			ShapeNet10			ScanNet10		
		mAP@1	mAP@5	mAP@10	mAP@1	mAP@5	mAP@10	mAP@1	mAP@5	mAP@10
PointNet [45]	Point cloud	80.1	91.7	94.4	90.6	96.6	98.1	65.7	86.2	92.6
PointNet++ [1]	Point cloud	85.1	93.9	96.0	92.2	97.5	98.6	71.6	89.3	93.7
DGCNN [2]	Point cloud	83.2	92.7	95.1	91.0	96.7	98.2	66.1	88.0	93.1
<i>nf2vec</i>	<i>NF</i>	81.7	92.6	95.1	90.6	96.7	98.1	65.2	87.5	94.0

TABLE II: Point cloud retrieval quantitative results.

Method	Input	ModelNet40	ShapeNet10	ScanNet10	Manifold40	ShapeNet10
PointNet [45]	Point cloud	88.8	94.3	72.7	-	-
PointNet++ [1]	Point cloud	89.7	94.6	76.4	-	-
DGCNN [2]	Point cloud	89.9	94.3	76.2	-	-
MeshWalker [47]	Mesh	-	-	-	90.0	-
Conv3DNet [48]	Voxels	-	-	-	-	92.1
<i>nf2vec</i>	<i>NF</i>	87.0	93.3	72.1	86.3	93.0

TABLE III: Results on shape classification across representations.

Method	Input	instance mIoU	class mIoU	Class															
				airplane	bag	cap	car	chair	earphone	guitar	knife	lamp	laptop	motor	mug	pistol	rocket	skateboard	table
PointNet [45]	Point cloud	83.1	78.96	81.3	76.9	79.6	71.4	89.4	67.0	91.2	80.5	80.0	95.1	66.3	91.3	80.6	57.8	73.6	81.5
PointNet++ [1]	Point cloud	84.9	82.73	82.2	88.8	84.0	76.0	90.4	80.6	91.8	84.9	84.4	94.9	72.2	94.7	81.3	61.1	74.1	82.3
DGCNN [2]	Point cloud	83.6	80.86	80.7	84.3	82.8	74.8	89.0	81.2	90.1	86.4	84.0	95.4	59.3	92.8	77.8	62.5	71.6	81.1
<i>nf2vec</i>	<i>NF</i>	81.3	76.91	80.2	76.2	70.3	70.1	88.0	65.0	90.6	82.1	77.4	94.4	61.4	92.7	79.0	56.2	68.6	78.5

TABLE IV: Part segmentation quantitative results. We report the IoU for each class, the mean IoU over all the classes (class mIoU) and the mean IoU over all the instances (instance mIoU).

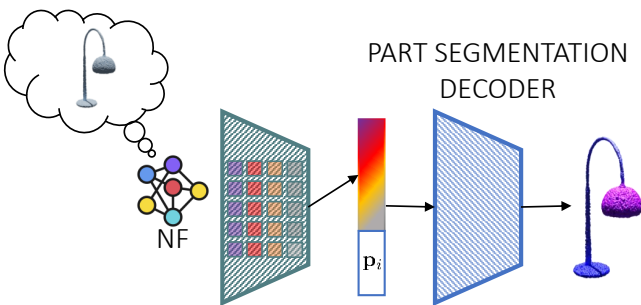


Fig. 12: Point cloud part segmentation. Method.



Fig. 13: Point cloud part segmentation. Qualitatives.

between embeddings of unseen point clouds from the test sets of ModelNet40 [44] and ShapeNet10 (a subset of 10 classes of the popular ShapeNet dataset [42]). For each embedded shape, we select its k -nearest-neighbours and compute a Precision Score comparing the classes of the query and the retrieved shapes, reporting the mean Average Precision for different k (mAP@ k). Beside *nf2vec*, we consider three baselines to embed point clouds, which are obtained by training the PointNet [45], PointNet++ [1] and DGCNN [2] encoders in combination with a fully connected decoder similar to that proposed in [46] to reconstruct the input cloud. The quantitative findings in Table II reveal that *nf2vec* not only matches but sometimes exceeds the performance of other baselines, with an average gap of 1.8 mAP compared to PointNet++. Furthermore, as depicted in Fig. 10, it is evident that the retrieved shapes not only belong to the same class as the query but also exhibit similar coarse structures. These results highlight that the pretext task used to learn *nf2vec* embeddings allows encoding relevant shape information.

Shape classification. We then address the problem of

classifying point clouds, meshes, and voxel grids. We use three datasets for point clouds: ShapeNet10, ModelNet40, and ScanNet10 [49]. When dealing with meshes, we conduct our experiments on the Manifold40 dataset [3]. Finally, we use ShapeNet10 again for voxel grids, quantizing clouds to grids with resolution 64^3 . Despite the different nature of the discrete representations taken into account, *nf2vec* allows us to perform shape classification on *NFs* embeddings, augmented online with E-Stitchup [50], by the very same downstream network architecture, *i.e.*, a simple fully connected classifier consisting of three layers with 1024, 512 and 128 features. We consider as baselines well-known architectures that are optimized to work on the specific input representations of each dataset. For point clouds, we consider PointNet [45], PointNet++ [1] and DGCNN [2]. For meshes, we consider MeshWalker [47], a recent and competitive baseline that processes triangle meshes directly. As for voxel grids, we train a 3D CNN classifier that we implemented following [48] (Conv3DNet from now on). Since only the train and test

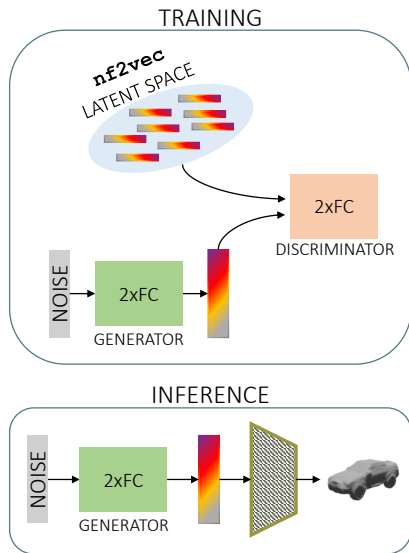


Fig. 14: Learning to generate shapes from $nf2vec$ latent space. Method.

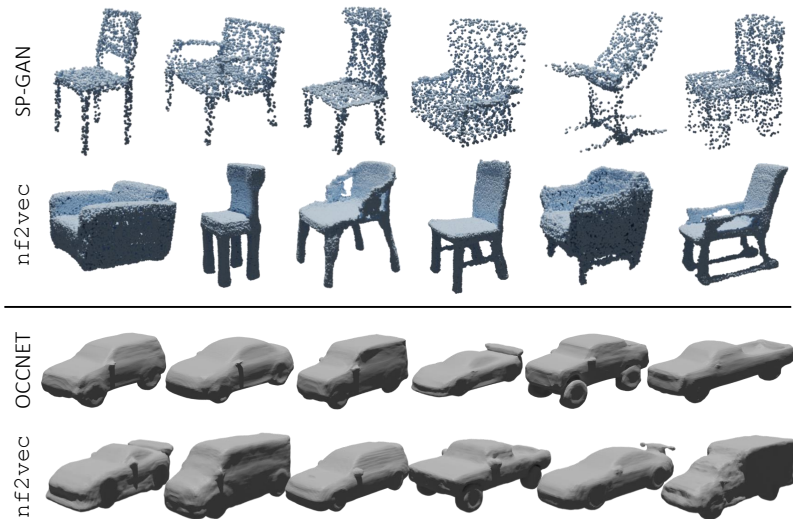


Fig. 15: Learning to generate shapes from $nf2vec$ latent space. Qualitative results.

splits are released for all the datasets, we created validation splits from the training sets in order to follow a proper train/val protocol for both the baselines and our method. As for the test shapes, we evaluated all the baselines on the discrete representations reconstructed from the NFs fitted on the original test sets, as these would be the only data available at test time in a scenario where NFs are used to store and communicate 3D data. The results in Table III show that $nf2vec$ embeddings deliver classification accuracy close to the specialized baselines across all the considered datasets, regardless of the original discrete representation of the shapes in each dataset. Remarkably, our framework allows us to apply the same simple classification architecture to all the considered input modalities, in stark contrast with all the baselines that are highly specialized for each modality, exploit inductive biases specific to each such modality and cannot be deployed on representations different from those they were designed for. Furthermore, while presenting a gap of some accuracy points *w.r.t.* the most recent architectures, like DGCNN and MeshWalker, the simple fully connected classifier that we applied on $nf2vec$ embeddings obtains scores comparable to standard baselines like PointNet and Conv3DNet.

Finally, in Fig. 11 (left), we present the baseline inference times, assuming discrete point clouds are available at test time, and compare it with that of $nf2vec$. Our framework is much faster than competitors. Indeed, by processing directly NFs – where the resolution of the underlying signal is theoretically infinite – $nf2vec$ can classify NFs representing point clouds with different numbers of points with a constant inference time of 0.001 seconds. On the contrary, the examined baselines suffer from the escalating resolution of the input point clouds. While PointNet and PointNet++ maintain a reasonable inference time even with 64K points, DGCNN experiences a significant slowdown as early as 16K points. Additionally, we emphasize that if 3D shapes are stored as NFs , the classification process

using the designated specialized baselines would involve retrieving the original discrete representations through the extensive procedures outlined in [22]. Therefore, in Fig. 11 (right), we present the inference time of standard point cloud classification networks, including the time needed to reconstruct the discrete point cloud from the input NF of the underlying udf at various resolutions. Even at the coarsest resolution (2048 points), all the baselines exhibit an inference time that is one order of magnitude higher than the time required to classify the $nf2vec$ embeddings directly. As the resolution of the reconstructed clouds increases, the inference time of the baselines becomes prohibitively high, while $nf2vec$, not reliant on explicit clouds, maintains a constant inference time of 0.001 seconds.

Point cloud part segmentation. The classification and retrieval tasks explore the potential of utilizing $nf2vec$ embeddings as a global representation of the input shapes. In contrast, in this section, we focus on point cloud part segmentation to examine if $nf2vec$ embeddings also retain local shape properties. Part segmentation aims to predict a semantic (*i.e.*, part) label for each point of a given cloud. We tackle this problem by training a decoder similar to that used to train our framework (see Fig. 12). Such decoder is fed with the $nf2vec$ embedding of the NF representing the input cloud, concatenated with the coordinate of a 3D query, and it is trained to predict the label of the query point. We train it, as well as PointNet, PointNet++, and DGCNN, on the ShapeNet Part Segmentation dataset [51] with point clouds of 2048 points, with the same train/val/test as in the classification task. The outcomes presented in Table IV demonstrate the potential of accomplishing a local discriminative task, such as part segmentation, using the task-agnostic embeddings generated by $nf2vec$. In doing so, the performance achieved is notably close to that of dedicated architectures designed for the specific task. Additionally, in Fig. 13, we show point

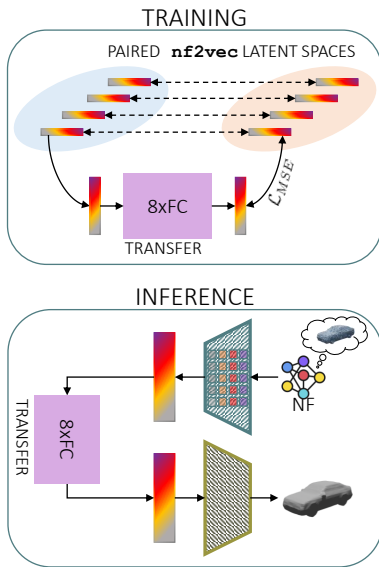


Fig. 16: Learning a mapping between $nf2vec$ latent spaces. Method.



Fig. 17: Learning a mapping between $nf2vec$ latent spaces. Point cloud completion.

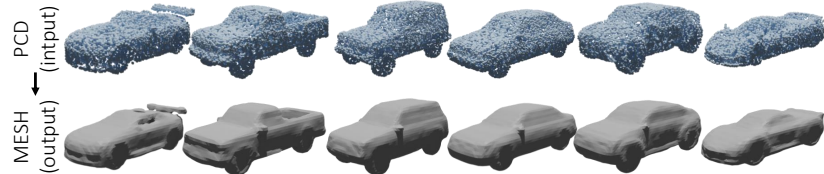


Fig. 18: Learning a mapping between $nf2vec$ latent spaces. Surface reconstruction.

clouds reconstructed at 100K points from the input NFs and segmented with high precision thanks to our formulation based on a semantic decoder conditioned by the $nf2vec$ embedding.

Shape generation. So far, we have validated that NF can be used as input in standard deep learning machinery thanks to $nf2vec$. In this section, we focus on the task of shape generation in an adversarial setting to investigate whether the compact representations produced by our framework can also be adopted as a medium for the output of generative deep learning pipelines. We employ a Latent-GAN [52] to generate embeddings resembling those produced by $nf2vec$ from random noise, as illustrated in Fig. 14. Generated embeddings can be decoded into discrete representations using the implicit decoder from $nf2vec$ training. As our framework is agnostic towards the original discrete representation of shapes used to learn the NFs , we can train Latent-GANs with embeddings representing point clouds or meshes based on the same identical protocol and architecture (two simple fully connected networks as generator and discriminator). For point clouds, we train a Latent-GAN on the *chair* class of ShapeNet10, while we use models of cars provided by [9] when dealing with meshes. In Fig. 15, we report some examples generated with the above procedure, comparing them with SP-GAN [53] for what concerns point clouds and Occupancy Networks [9] (VAE formulation) for meshes. Shapes generated by our Latent-GAN, trained exclusively on $nf2vec$ embeddings, look quite similar to those from the considered baselines regarding diversity and richness of details. Furthermore, by generating embeddings representing NFs , our method allows point cloud sampling at any arbitrary resolution (e.g., 8192 points in Fig. 15). In contrast, SP-GAN necessitates new training for each desired resolution, as the number of generated points must be predetermined during training.

Learning a mapping between $nf2vec$ embedding spaces.

We have shown that $nf2vec$ embeddings can be employed as a proxy of NFs as input to deep learning pipelines and that they can also be obtained as output of generative frameworks. In this section, we advance our exploration by examining the potential of learning a mapping between two distinct latent spaces generated by our framework for two separate datasets of NFs . This involves developing a *transfer* function specifically designed to operate on $nf2vec$ embeddings as both input and output data. Such transfer function can be realized by a simple MLP that maps the input embedding into the output one and is trained by a standard MSE loss (see Fig. 16). As $nf2vec$ generates compact embeddings of the same dimension regardless of the input NF modality, the transfer function described here can be applied seamlessly to a great variety of tasks, usually tackled with ad-hoc frameworks tailored to specific input/output modalities. In particular, We explore two tasks. First, on the dataset presented in [54], we address point cloud completion by learning a mapping from $nf2vec$ embeddings of NFs that represent incomplete clouds to embeddings associated with complete clouds. Then, we tackle the task of surface reconstruction on ShapeNet cars, training the transfer function to map $nf2vec$ embeddings representing point clouds into embeddings that can be decoded into meshes. As we note from Fig. 17 and Fig. 18, in both tasks, the transfer function can learn an effective mapping between $nf2vec$ latent spaces. Indeed, by processing exclusively NFs embedding, we can obtain output shapes that are highly compatible with the input ones whilst preserving the distinctive details, like the pointy wing of the airplane in Fig. 17 or the flap of the first car in Fig. 18.

VI. DEEP LEARNING ON NeRFs

In this section, our focus shifts to processing NFs encoding both geometry and appearance of objects, i.e., NeRFs. The goal

Method	N. views	mAP@1	mAP@5	mAP@10
<i>nf2vec</i>	-	72.38	91.89	95.96
ResNet50 [55] single-view	1	74.65	91.52	95.10
ResNet50 [55] multi-view	9	82.74	91.66	93.79

TABLE V: NeRF retrieval quantitative results.

Method	N. views	Accuracy
<i>nf2vec</i>	-	87.28%
ResNet50 [55] Single-view	1	86.88%
ResNet50 [55] Multi-view	9	93.28%

TABLE VI: Results on NeRF classification.

is to illustrate the efficacy of *nf2vec* in addressing various downstream tasks related to 3D objects implicitly represented by NeRFs.

General settings. In all experiments detailed within this section, we learn NeRFs from images using an MLP comprising three hidden layers with 64 nodes each. We utilize the ReLU activation function between all layers except the final layer, which computes the density and RGB values without any activation function. NeRFs take as input the frequency encoding of the 3D coordinates as in [11]. NeRFs are trained using an L_1 loss between predicted and estimated RGB pixel intensities, weighting background pixels less than foreground pixels (0.8 foreground vs 0.2 background). We use the NeRF formulation without the view direction in input. When training *nf2vec*, we adhere to the same encoder and decoder architectures previously described in Section V. The key distinction lies in the implicit decoder, wherein the dimensions of its layers are doubled. Throughout these experiments, both baseline models and *nf2vec* undergo training with offline data augmentation on the 3D shapes used to generate renderings for training the models. This augmentation involves implementing random deformations and introducing random color changes to each component of the original 3D shapes.

NeRF retrieval. We first investigate the quality of *nf2vec* embeddings with a retrieval task as done in Section V. Given an embedding of an input NeRF, we extract the classes of its k -nearest neighbor within the embedding latent space and compare them to the input NeRF’s class.

We also implement two baseline approaches: the single-view and multi-view baselines. Both strategies rely on a *ResNet50* [55] backbone pre-trained on ImageNet [56]. We extract feature vectors with *ResNet50* from each image. Given a single image for the single-view baseline or 9 images for the multi-view baseline, we find the k -nearest neighbors in the *ResNet50* feature space. We compare the classes of query and retrieved objects and compute the mAP for different k as done in Table II. In the case of the multi-view baseline, we retrieve the nearest neighbors for each of the 9 input images. Then, we select the class with the highest frequency to calculate the mAP. All these experiments are conducted on embeddings of unseen NeRFs from the test set of ShapeNetRender [57]. The quantitative results in Table V indicate that *nf2vec* yields comparable performance to the baselines and, in certain instances, outperforms them. Moreover, Fig. 19 shows that the selected neighbors exhibit similar structures and colors.

NeRF classification. This section investigates the task of predicting the category of an object represented by a NeRF. In this scenario, only NeRFs would be available as input data.

Method	Encoding (ms)	Rendering (ms)	Classification (ms)	Total (ms)
<i>nf2vec</i>	0.75	-	0.28	1.03
Resnet50 [55] Single-View	-	4.63	6.21	10.84
Resnet50 [55] Multi-View	-	41.67	55.89	97.56

TABLE VII: NeRFs Classification Inference Time. All times are in milliseconds. For each method, we report the time needed for each pipeline step and the total time. Encoding: *nf2vec* encoding of NF weights. Rendering: obtaining images from NeRFs. Classification: the time required by the classifier.

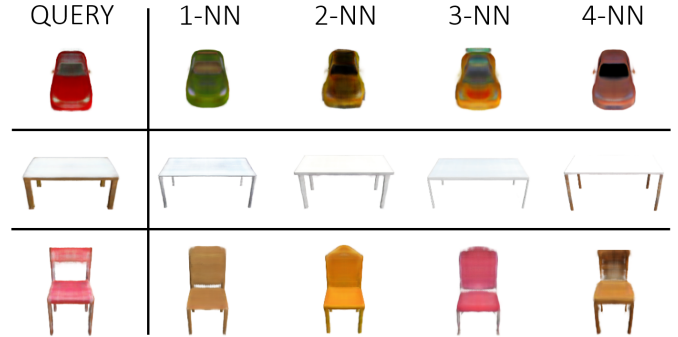


Fig. 19: NeRF retrieval qualitative results. Given the *nf2vec* embedding of a query NeRF, we show the renderings reconstructed from the closest embeddings (L_2 distance).

Our approach processes *nf2vec* embeddings with the same classification architecture as already deployed for shapes, a three-layer MLP with 1024, 512, and 128 neurons each. We highlight that the embeddings are obtained by processing the NeRFs weights without sampling the underlying RF .

As the discrete representations used to learn NeRFs are a set of images depicting the same object, selecting a proper baseline is not straightforward. In our experiment, we choose ResNet50 [55] as the baseline classifier. The network predicts the class for a given input image. Given this architecture, akin to the retrieval experiment, we propose two types of baseline approaches, single-view and multi-view. In the former, we train the network on a single rendering for each NeRF obtained from the same fixed pose, while for the latter, we employ 9 renderings for each NeRF from different viewpoints. At test time, regarding the single-view approach, we test the network on images rendered from unseen NeRFs employing the same training pose. Concerning the multi-view baseline, we render 9 images from the training viewpoints for each unseen NeRF, obtaining 9 distinct predictions per object, that we aggregate by taking the class predicted with the highest frequency.

We report the accuracy results on ShapeNetRender [57] in Table VI. Moreover, we also report the time required to classify NeRFs in Table VII, highlighting the impact of each of the main pipeline steps, such as the *nf2vec* encoding time for our approach and the time to render images for the baseline approaches. We point out that we rely on a fast implementation of NeRFs provided by the NerfAcc [58] library. We note that the classifier directly leveraging *nerf2vec* embeddings achieves a slightly better score than the single-view baseline while being worse than the multi-view one. However, if we analyze the total inference times, our approach is two orders of magnitude faster than the baselines (1.03ms our vs 97.56ms

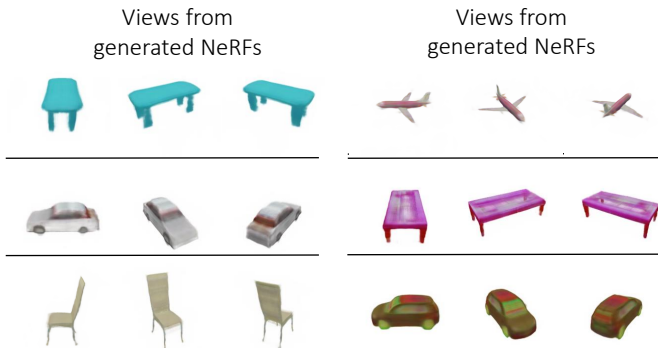


Fig. 20: Learning to generate NeRFs from $nf2vec$ latent space. Qualitative results.

Method	Type	# Params	Mesh Reconstruction	
			CD (mm) ↓	F-score (%) ↑
DeepSDF [5]	Shared MLP + Embedding	2400K	6.6	25.1
Functa [15]	Shared MLP + Modulation	7091K	2.85	21.3
DWS [24]	Individual SIREN MLP	800K	0.26	69.7
$nf2vec$	Individual SIREN MLP	800K	0.26	69.7

TABLE VIII: Properties of input NFs used by recent methods processing neural fields. Mesh reconstruction results on Manifold40 test set.

ResNet50 multi-view). This holds true even excluding the rendering time and looking only at the classification time (third column of Table VII): ResNet50 classifiers processing 224×224 images are remarkably slower than our full pipeline, taking only 1.03ms.

NeRF generation. We experiment here with the task of generating $nf2vec$ embeddings with the same approach depicted in Fig. 14. We trained multiple adversarial networks, one for each class, utilizing the ShapeNetRender dataset [57]. After the embedding generation, we can decode them into discrete representations using the same implicit decoder employed to train the framework. In Fig. 20, we show images rendered from NeRFs generated through this process. The renderings have a good level of realism and diversity. Notably, the 3D consistency of images obtained from different viewpoints is preserved. These results show that generating novel NeRFs in the form of $nf2vec$ embeddings is possible.

Learning a mapping between embedding spaces. We explore here whether it is possible to learn a *transfer* network that maps $nf2vec$ embeddings of $UDFs$ to $nf2vec$ embeddings of NeRFs using the methodology described in Fig. 16. We highlight that it is a non-trivial task as the network has to hallucinate the appearance of the input shape without modifying the underlying 3D structure.

We run this experiment on ShapeNetRender [57], using the rendered images to learn NeRFs and the corresponding 3D models to learn UDF neural fields. Then, we train $nf2vec$ and then the *transfer* network on NFs of the training set. Finally, we test the *transfer* network on unseen test NF , obtaining for each UDF latent code the corresponding mapped $nf2vec$ embedding. We report qualitative results in Fig. 21, showing the point clouds obtained from the input UDF in the first row, and the renderings obtained by feeding the mapped embedding to the implicit decoder of $nf2vec$ in the last three rows. We can appreciate that the mapped NeRF preserves the original



Fig. 21: Learning a mapping between $nf2vec$ latent spaces. UDF to NeRF.

Method	Type	Input	ModelNet40 (UDF)	Manifold40 (SDF)
DeepSDF [5]	Shared	Embedding	41.2	64.9
Functa [15]	Shared	Modulation	87.3	85.9
DWS [24]	Individual	MLP weights	71.6	69.9
$nf2vec$	Individual	MLP weights	87.0	86.3

TABLE IX: Classification accuracy of recent methods processing neural fields.

shape geometry, enabling the rendering of realistic images from various viewpoints. Notably, the renderings exhibit plausible diverse colors associated with different object parts, as can be seen for the glasses and the wheels of the blue car in the second column of Fig. 21.

VII. COMPARISON WITH RECENT APPROACHES

As outlined in Section I, several contemporary works addressing the problem of processing neural fields have been proposed recently. For all methods, the goal is to perform deep learning tasks such as classification using as input data a NF , *i.e.*, data represented with continuous functions. We can divide these methods into two categories. Those relying on a shared network and those focusing on individual NFs . In the former case, referred to here as *Shared*, the NF is defined as a shared network trained on all training samples, plus a distinct vector representing each object. Typically this vector is processed to perform downstream tasks. This is the case of Functa [15] and DeepSDF [5]. In the latter case, denoted as *Individual*, the NF is typically an MLP trained on a single object or scene. In this scenario, the MLP weights are processed directly to perform the downstream tasks. This is the case of our framework, $nf2vec$, as well as NFN [41], NFT [23], and DWS [24].

In this section, we investigate the characteristics of each category of techniques, showing that *Shared* frameworks are problematic, as they cannot reconstruct the underlying signal with high fidelity and need a whole dataset to learn the neural field of an object. Moreover, we build the first benchmark of NF classification by comparing recent approaches in this area.

Representation quality. We first investigate the representation quality of *Shared* approaches compared to *Individual* ones. Specifically, we compare the reconstructions of explicit meshes from SDF neural fields with ground truth meshes on the Manifold40 test set. We report the quantitative comparisons

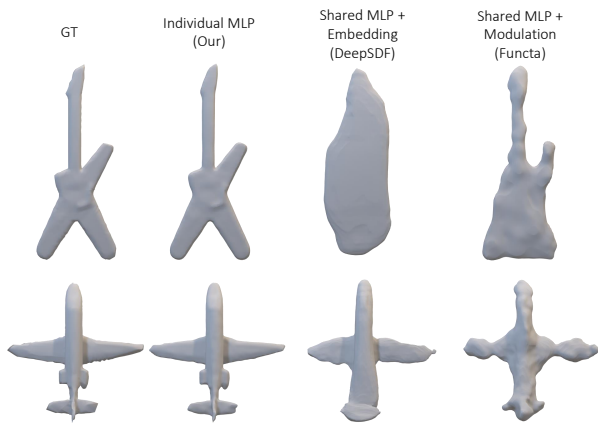


Fig. 22: Reconstruction comparison for Manifold40 meshes obtained from *SDF*.

in Table VIII, using two metrics: the Chamfer Distance (CD) as defined in [46], and the F-Score as defined in [59]. We note that we use the SIREN MLP described in Section V to represent *SDF* with *Individual* frameworks. In the first two rows, we note that *Shared* methods achieve poor reconstruction performance. Indeed, we believe that representing a whole dataset with a shared network is a difficult learning task, and the network struggles to fit accurately the totality of the samples. *Individual* methods instead do not suffer this problem and achieve very good reconstruction performances. Moreover, we believe that the approaches based on *Shared* networks fail to represent unseen samples that are out of the training distribution. Hence, in the foreseen scenario where *NFs* become a standard representation for 3D data hosted in public repositories, leveraging on a single shared network may imply the need to frequently retrain the model upon uploading new samples, which, in turn, would change the embeddings of all the previously stored data. On the contrary, uploading the repository with a new object would not cause any sort of issue with individual *NFs*, where one learns a network for each data point. Finally, we also provide a qualitative perspective of the abovementioned problem in Fig. 22 and Fig. 23. The visualizations confirm the results of Table VIII, with shared network frameworks struggling to represent properly the ground-truth shapes, while individual *NFs* enable high-fidelity reconstructions.

We believe that these results highlight that frameworks based on a single shared network cannot be used as a medium to represent objects as *NFs*, because of their limited representation power when dealing with large and varied datasets and because of their difficulty in representing new shapes not available at training time.

Classification accuracy. We compare recent methods in the *NFs* classification task. The goal is to predict the category of objects represented within the input *NFs* without recreating the discrete signals. Specifically, we test all methods on *UDF* obtained from point clouds of ModelNet40 [44] and on *SDF* learned from meshes of Manifold40 [3]. We compare `nf2vec` with other frameworks designed to process *NFs* realized as Individual MLPs, such as DWSNet [24]. These methods process

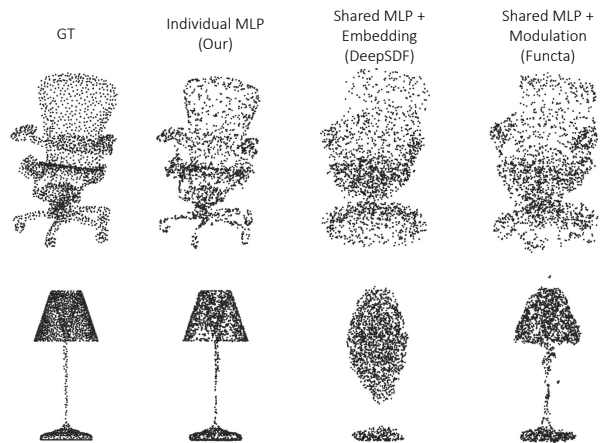


Fig. 23: Reconstruction comparison for ModelNet40 point clouds obtained from *UDF*.

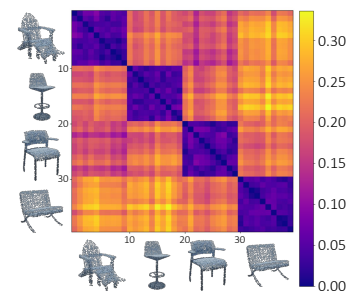


Fig. 24: **L2 distances between `nf2vec` embeddings.** For each shape, we fit 10 *NFs* starting from the same weights initialization (40 *NFs* in total). Then we plot the L2 distances between the embeddings obtained by `nf2vec` for such *NFs*.

the MLP weights of individual *NFs*, implemented as SIREN networks [16]. The MLPs in our benchmark are initialized using the same random seed (see Section VIII for more details). We also tried to run other recent *Individual* approaches, such as NFN [41] and NFT [23], but the training did not converge, probably because the SIREN MLPs used in our experiments are much larger than those used in their paper. Moreover, we compare with *Shared* frameworks where neural fields are realized by a shared network and a small latent vector or modulation, *i.e.*, DeepSDF [5] and Funcsta [15]. Whenever possible, we use the official code released by the authors to run the experiments.

As we can see from results reported in Table IX, Funcsta and `nf2vec` achieve the best results with a large margin over other competitors, with the former slightly more accurate on *UDFs* while the latter on *SDFs*. However, as explained in the previous section, since our method does not rely on shared networks, it does not sacrifice the representation quality of *NFs* and it is more suitable for real-world applications.

VIII. USING THE SAME INITIALIZATION FOR *NFs*

The need to align the multitude of *NFs* that can approximate a given shape is a challenging research problem that has to be dealt with when using *NFs* as input data. We empirically found that fixing the weights initialization to a shared random vector across *NFs* is a viable and simple solution to this problem.

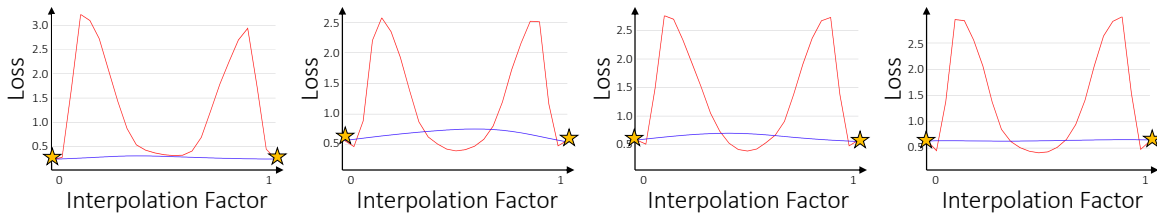


Fig. 25: **Linear mode connectivity study.** Each plot shows the variation of the loss function over the same batch of points when interpolating between two NFs representing the same shape. The red line describes the interpolation between NFs initialized differently, while the blue line shows the same interpolation between NFs initialized from the same random vector. The yellow stars represent the loss value of the boundary NFs .

We report here an experiment to assess if the order of data or other sources of randomness arising while fitting NFs do affect the repeatability of the embeddings computed by `nf2vec`. We fitted 10 NFs on the same discrete shape for 4 different chairs, *i.e.*, 40 NFs in total. Then, we embed all of them with the pretrained `nf2vec` encoder and compute the L2 distance between all pairs of embeddings. The block structure of the resulting distance matrix (see Fig. 24) highlights how, under the assumption of shared initialization, `nf2vec` is repeatable across multiple fittings.

Seeking for a proof with a stronger theoretical foundation, we turn our attention to the recent work *git re-basin* [60], where authors show that the loss landscape of neural networks contains (nearly) a single basin after accounting for all possible permutation symmetries of hidden units. The intuition behind this finding is that, given two neural networks that were trained with equivalent architectures but different random initializations, data orders, and potentially different hyperparameters or datasets, it is possible to find a permutation of the network’s weights such that when linearly interpolating between their weights, all intermediate models enjoy performance similar to them – a phenomenon denoted as *linear mode connectivity*.

Intrigued by this finding, we conducted a study to assess whether initializing NFs with the same random vector, which we found to be key to `nf2vec` convergence, also leads to linear mode connectivity. Thus, given one shape, we fitted it with two different NFs , and then we interpolated linearly their weights, observing at each interpolation step the loss value obtained by the *interpolated NF* on the same batch of points. We repeated the experiment twice for each shape, once initializing the NFs with different random vectors and once with the same random vector.

The results of this experiment are reported for four different shapes in Fig. 25. It is possible to note that, as shown by the blue curves, when interpolating between NFs obtained from the same weights initialization, the loss value at each interpolation step is nearly identical to those of the boundary NFs . On the contrary, the red curves highlight how there is no linear mode connectivity at all between NFs obtained from different weights initializations.

[60] also proposes different algorithms to estimate the permutation needed to obtain linear mode connectivity between two networks. We applied the algorithm proposed in their paper in Section 3.2 (*Matching Weights*) to our NFs and observed the resulting permutations. Remarkably, when applied to NFs

obtained from the same weights initialization, the retrieved permutations are identity matrices, both when the target NFs represent the same shape and when they represent different ones. Instead, the permutations obtained for NFs obtained from different initializations are far from being identity matrices.

All these results favor the hypothesis that our technique of initializing NFs with the same random vector leads to linear mode connectivity between different NFs . We believe that the possibility of performing meaningful linear interpolation between the weights occupying the same positions across different NFs can be interpreted by considering corresponding weights as carrying out the same role in terms of feature detection units, explaining why the `nf2vec` encoder succeeds in processing the weights of our NFs .

IX. CONCLUDING REMARKS

We have shown that it is possible to apply deep learning on individual NFs representing 3D shapes and radiance fields. Our approach leverages a task-agnostic encoder which embeds NFs into compact and meaningful latent codes without accessing the underlying function. We have shown that these embeddings can be fed to standard deep-learning machinery to solve various tasks effectively. Moreover, we have introduced the first benchmark for the task of NF classification, showing that our proposal obtains the best score (on par with Functua [15]) while preserving the ability to reconstruct the input dataset with high quality.

We point out two main limitations of our approach: i) Although NFs capture continuous geometric cues, in some cases, deep learning on `nf2vec` embeddings achieve results inferior to state-of-the-art solutions that work on specific discrete representations ii) There is no obvious way to perform online data augmentation on shapes represented as NFs by directly altering their weights.

In the future, we plan to investigate these shortcomings and apply `nf2vec` to NFs encoding other input modalities, like images or audio. We will also investigate weight-space symmetries [61] as a different path to favor the alignment of weights across NFs .

We reckon that our work may foster the adoption of NFs as a unified 3D representation, overcoming the current fragmentation of 3D structures and processing architectures.

REFERENCES

- [1] C. R. Qi, L. Yi, H. Su, and L. J. Guibas, "Pointnet++: Deep hierarchical feature learning on point sets in a metric space," *Advances in neural information processing systems*, vol. 30, 2017.
- [2] Y. Wang, Y. Sun, Z. Liu, S. E. Sarma, M. M. Bronstein, and J. M. Solomon, "Dynamic graph cnn for learning on point clouds," *Acm Transactions On Graphics (tog)*, vol. 38, no. 5, pp. 1–12, 2019.
- [3] S.-M. Hu, Z.-N. Liu, M.-H. Guo, J.-X. Cai, J. Huang, T.-J. Mu, and R. R. Martin, "Subdivision-based mesh convolution networks," *ACM Transactions on Graphics (TOG)*, vol. 41, no. 3, pp. 1–16, 2022.
- [4] Y. Xie, T. Takikawa, S. Saito, O. Litany, S. Yan, N. Khan, F. Tombari, J. Tompkin, V. Sitzmann, and S. Sridhar, "Neural fields in visual computing and beyond," *arXiv preprint arXiv:2111.11426*, 2021.
- [5] J. J. Park, P. Florence, J. Straub, R. Newcombe, and S. Lovegrove, "Deepsdf: Learning continuous signed distance functions for shape representation," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019, pp. 165–174.
- [6] J. Chibane, G. Pons-Moll *et al.*, "Neural unsigned distance fields for implicit function learning," *Advances in Neural Information Processing Systems*, vol. 33, pp. 21 638–21 652, 2020.
- [7] A. Gropp, L. Yariv, N. Haim, M. Atzmon, and Y. Lipman, "Implicit geometric regularization for learning shapes," in *International Conference on Machine Learning*. PMLR, 2020, pp. 3789–3799.
- [8] T. Takikawa, J. Litalien, K. Yin, K. Kreis, C. Loop, D. Nowrouzezahrai, A. Jacobson, M. McGuire, and S. Fidler, "Neural geometric level of detail: Real-time rendering with implicit 3d shapes," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021, pp. 11 358–11 367.
- [9] L. Mescheder, M. Oechsle, M. Niemeyer, S. Nowozin, and A. Geiger, "Occupancy networks: Learning 3d reconstruction in function space," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019, pp. 4460–4470.
- [10] S. Peng, M. Niemeyer, L. Mescheder, M. Pollefeys, and A. Geiger, "Convolutional occupancy networks," in *European Conference on Computer Vision*. Springer, 2020, pp. 523–540.
- [11] B. Mildenhall, P. P. Srinivasan, M. Tancik, J. T. Barron, R. Ramamoorthi, and R. Ng, "Nerf: Representing scenes as neural radiance fields for view synthesis," in *European conference on computer vision*. Springer, 2020, pp. 405–421.
- [12] T. Müller, A. Evans, C. Schied, and A. Keller, "Instant neural graphics primitives with a multiresolution hash encoding," *ACM Trans. Graph.*, vol. 41, no. 4, pp. 102:1–102:15, Jul. 2022. [Online]. Available: <https://doi.org/10.1145/3528223.3530127>
- [13] J. N. P. Martel, D. B. Lindell, C. Z. Lin, E. R. Chan, M. Monteiro, and G. Wetzstein, "Acorn: Adaptive coordinate networks for neural scene representation," *ACM Trans. Graph. (SIGGRAPH)*, vol. 40, no. 4, 2021.
- [14] H.-T. D. Liu, F. Williams, A. Jacobson, S. Fidler, and O. Litany, "Learning smooth neural functions via lipschitz regularization," *arXiv preprint arXiv:2202.08345*, 2022.
- [15] E. Dupont, H. Kim, S. A. Eslami, D. J. Rezende, and D. Rosenbaum, "From data to functa: Your data point is a function and you can treat it like one," in *International Conference on Machine Learning*. PMLR, 2022, pp. 5694–5725.
- [16] V. Sitzmann, J. Martel, A. Bergman, D. Lindell, and G. Wetzstein, "Implicit neural representations with periodic activation functions," *Advances in Neural Information Processing Systems*, vol. 33, pp. 7462–7473, 2020.
- [17] V. Sitzmann, E. Chan, R. Tucker, N. Snavely, and G. Wetzstein, "Metasdf: Meta-learning signed distance functions," *Advances in Neural Information Processing Systems*, vol. 33, pp. 10 136–10 147, 2020.
- [18] E. Dupont, A. Goliński, M. Alizadeh, Y. W. Teh, and A. Doucet, "Coin: Compression with implicit neural representations," *arXiv preprint arXiv:2103.03123*, 2021.
- [19] Y. Strümpfer, J. Postels, R. Yang, L. Van Gool, and F. Tombari, "Implicit neural representations for image compression," *arXiv preprint arXiv:2112.04267*, 2021.
- [20] Y. Zhang, T. van Rozendaal, J. Brehmer, M. Nagel, and T. Cohen, "Implicit neural video compression," *arXiv preprint arXiv:2112.11312*, 2021.
- [21] M. Tancik, P. Srinivasan, B. Mildenhall, S. Fridovich-Keil, N. Raghavan, U. Singhal, R. Ramamoorthi, J. Barron, and R. Ng, "Fourier features let networks learn high frequency functions in low dimensional domains," *Advances in Neural Information Processing Systems*, vol. 33, pp. 7537–7547, 2020.
- [22] L. De Luigi, A. Cardace, R. Spezialetti, P. Zama Ramirez, S. Salti, and L. Di Stefano, "Deep learning on implicit neural representations of shapes," in *International Conference on Learning Representations (ICLR)*, 2023.
- [23] A. Zhou, K. Yang, Y. Jiang, K. Burns, W. Xu, S. Sokota, J. Z. Kolter, and C. Finn, "Neural functional transformers," *Advances in neural information processing systems*, vol. 37, 2023.
- [24] A. Navon, A. Shamsian, I. Achituve, E. Fetaya, G. Chechik, and H. Maron, "Equivariant architectures for learning in deep weight spaces," in *International Conference on Machine Learning*, 2023.
- [25] J. Frankle and M. Carbin, "The lottery ticket hypothesis: Finding sparse, trainable neural networks," *arXiv preprint arXiv:1803.03635*, 2018.
- [26] T. Choudhary, V. Mishra, A. Goswami, and J. Sarangapani, "A comprehensive survey on model compression and acceleration," *Artificial Intelligence Review*, vol. 53, no. 7, pp. 5113–5155, 2020.
- [27] V. Sitzmann, M. Zollhöfer, and G. Wetzstein, "Scene representation networks: Continuous 3d-structure-aware neural scene representations," *Advances in Neural Information Processing Systems*, vol. 32, 2019.
- [28] C. Jiang, A. Sud, A. Makadia, J. Huang, M. Nießner, T. Funkhouser *et al.*, "Local implicit grid representations for 3d scenes," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 6001–6010.
- [29] M. Atzmon and Y. Lipman, "Sal: Sign agnostic learning of shapes from raw data," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 2565–2574.
- [30] Z. Chen and H. Zhang, "Learning implicit fields for generative shape modeling," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019, pp. 5939–5948.
- [31] S. Saito, Z. Huang, R. Natsume, S. Morishima, A. Kanazawa, and H. Li, "Pifu: Pixel-aligned implicit function for high-resolution clothed human digitization," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2019, pp. 2304–2314.
- [32] M. Oechsle, L. Mescheder, M. Niemeyer, T. Strauss, and A. Geiger, "Texture fields: Learning texture representations in function space," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2019, pp. 4531–4540.
- [33] M. Niemeyer, L. Mescheder, M. Oechsle, and A. Geiger, "Differentiable volumetric rendering: Learning implicit 3d representations without 3d supervision," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 3504–3515.
- [34] —, "Occupancy flow: 4d reconstruction by learning particle dynamics," in *Proceedings of the IEEE/CVF international conference on computer vision*, 2019, pp. 5379–5389.
- [35] T. Unterthiner, D. Keysers, S. Gelly, O. Bousquet, and I. O. Tolstikhin, "Predicting neural network accuracy from weights," *arXiv*, vol. abs/2002.11448, 2020.
- [36] K. Schürholt, D. Kostadinov, and D. Borth, "Self-supervised representation learning on neural network weights for model characteristic prediction," in *Advances in Neural Information Processing Systems*, A. Beygelzimer, Y. Dauphin, P. Liang, and J. W. Vaughan, Eds., 2021. [Online]. Available: <https://openreview.net/forum?id=F1D8buayXQT>
- [37] B. Knyazev, M. Drozdal, G. W. Taylor, and A. Romero, "Parameter prediction for unseen deep architectures," in *Advances in Neural Information Processing Systems*, A. Beygelzimer, Y. Dauphin, P. Liang, and J. W. Vaughan, Eds., 2021. [Online]. Available: <https://openreview.net/forum?id=vqHak8NLk25>
- [38] F. Jaecle and M. P. Kumar, "Generating adversarial examples with graph neural networks," in *Uncertainty in Artificial Intelligence*. PMLR, 2021, pp. 1556–1564.
- [39] J. Lu and M. P. Kumar, "Neural network branching for neural network verification," in *International Conference on Learning Representations*, 2020. [Online]. Available: <https://openreview.net/forum?id=B1evfa4tPB>
- [40] R. Hecht-Nielsen, "On the algebraic structure of feedforward network weight spaces," in *Advanced Neural Computers*. Elsevier, 1990, pp. 129–135.
- [41] A. Zhou, K. Yang, K. Burns, A. Cardace, Y. Jiang, S. Sokota, J. Z. Kolter, and C. Finn, "Permutation equivariant neural functionals," *Advances in neural information processing systems*, vol. 37, 2023.
- [42] A. X. Chang, T. Funkhouser, L. Guibas, P. Hanrahan, Q. Huang, Z. Li, S. Savarese, M. Sava, S. Song, H. Su *et al.*, "Shapenet: An information-rich 3d model repository," *arXiv preprint arXiv:1512.03012*, 2015.
- [43] W. E. Lorensen and H. E. Cline, "Marching cubes: A high resolution 3d surface construction algorithm," *ACM siggraph computer graphics*, vol. 21, no. 4, pp. 163–169, 1987.
- [44] Z. Wu, S. Song, A. Khosla, F. Yu, L. Zhang, X. Tang, and J. Xiao, "3d shapenets: A deep representation for volumetric shapes," in *2015 IEEE*

Conference on Computer Vision and Pattern Recognition (CVPR), 2015, pp. 1912–1920.

- [45] C. R. Qi, H. Su, K. Mo, and L. J. Guibas, “Pointnet: Deep learning on point sets for 3d classification and segmentation,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 652–660.
- [46] H. Fan, H. Su, and L. J. Guibas, “A point set generation network for 3d object reconstruction from a single image,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 605–613.
- [47] A. Lahav and A. Tal, “Meshwalker: Deep mesh understanding by random walks,” *ACM Transactions on Graphics (TOG)*, vol. 39, no. 6, pp. 1–13, 2020.
- [48] D. Maturana and S. Scherer, “Voxnet: A 3d convolutional neural network for real-time object recognition,” in *2015 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE, 2015, pp. 922–928.
- [49] A. Dai, A. X. Chang, M. Savva, M. Halber, T. Funkhouser, and M. Nießner, “ScanNet: Richly-annotated 3d reconstructions of indoor scenes,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 5828–5839.
- [50] C. R. Wolfe and K. T. Lundgaard, “E-stitchup: Data augmentation for pre-trained embeddings,” *arXiv preprint arXiv:1912.00772*, 2019.
- [51] L. Yi, V. G. Kim, D. Ceylan, I.-C. Shen, M. Yan, H. Su, C. Lu, Q. Huang, A. Sheffer, and L. Guibas, “A scalable active framework for region annotation in 3d shape collections,” *SIGGRAPH Asia*, 2016.
- [52] P. Achlioptas, O. Diamanti, I. Mitliagkas, and L. Guibas, “Learning representations and generative models for 3d point clouds,” in *International conference on machine learning*. PMLR, 2018, pp. 40–49.
- [53] R. Li, X. Li, K.-H. Hui, and C.-W. Fu, “Sp-gan: Sphere-guided 3d shape generation and manipulation,” *ACM Transactions on Graphics (TOG)*, vol. 40, no. 4, pp. 1–12, 2021.
- [54] L. Pan, X. Chen, Z. Cai, J. Zhang, H. Zhao, S. Yi, and Z. Liu, “Variational relational point completion network,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021, pp. 8524–8533.
- [55] K. He, X. Zhang, S. Ren, and J. Sun, “Deep residual learning for image recognition,” in *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016, pp. 770–778.
- [56] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, Alexander C. Berg and Li Fei-Fei, “Imagenet large scale visual recognition challenge,” *arXiv*, vol. abs/1409.0575, 2015.
- [57] Q. Xu, W. Wang, D. Ceylan, R. Mech, and U. Neumann, “Disn: Deep implicit surface network for high-quality single-view 3d reconstruction,” in *Advances in Neural Information Processing Systems*, H. Wallach, H. Larochelle, A. Beygelzimer, F. d’Alché-Buc, E. Fox, and R. Garnett, Eds., vol. 32. Curran Associates, Inc., 2019. [Online]. Available: https://proceedings.neurips.cc/paper_files/paper/2019/file/39059724f73a9969845dfe4146c5660e-Paper.pdf
- [58] R. Li, H. Gao, M. Tancik, and A. Kanazawa, “Nerfacc: Efficient sampling accelerates nerfs,” *arXiv preprint arXiv:2305.04966*, 2023.
- [59] M. Tatarchenko, S. R. Richter, R. Ranftl, Z. Li, V. Koltun, and T. Brox, “What do single-view 3d reconstruction networks learn?” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019, pp. 3405–3414.
- [60] S. K. Ainsworth, J. Hayase, and S. Srinivasa, “Git re-basin: Merging models modulo permutation symmetries,” *arXiv preprint arXiv:2209.04836*, 2022.
- [61] R. Entezari, H. Sedghi, O. Saukh, and B. Neyshabur, “The role of permutation invariance in linear mode connectivity of neural networks,” in *International Conference on Learning Representations*, 2021.



Luca De Luigi received his PhD in Computer Science and Engineering in 2023. He is currently a computer vision engineer at eyecan.ai. His research focuses on deep learning for computer vision problems, especially in 3D geometry and neural fields.



Daniele Sirocchi received his master’s degree in Artificial Intelligence in 2023. He currently works as a software engineer, mainly focused on creating applications that blend various technologies, encompassing a wide spectrum of technological domains. His research interests primarily revolve around machine and deep learning for computer vision tasks, where he extensively engages with neural fields.



Adriano Cardace is a third-year PhD student at the Computer Vision Laboratory (CVLab), University of Bologna. He has been a Research Intern at Mitsubishi Electric Research Laboratories (MERL) and authored several research papers covering a range of subjects, including semantic segmentation, Domain Adaptation, and Neural Fields.



Riccardo Spezialetti received his PhD degree in Computer Science and Engineering from University of Bologna in 2020. After two years as a Post-doc researcher at the Department of Computer Science and Engineering, University of Bologna, he recently joined eyecan.ai as a computer vision engineer. His research interest concerns machine/deep learning for 3D computer vision problems.



Francesco Ballerini received his Master’s degree in Artificial Intelligence in 2023 and is currently a PhD student in Computer Science and Engineering, both at the University of Bologna. His main research interests include the study of neural architectures aimed at processing neural fields, specifically those representing 3D data, with a focus on permutation symmetries and the effect of initialization on the input neural fields themselves.



Samuele Salti is currently an associate professor at the Department of Computer Science and Engineering (DISI) of the University of Bologna, Italy. His main research interest is computer vision, mainly 3D computer vision and machine/deep learning applied to computer vision problems. Dr. Salti has co-authored more than 60 publications and 8 international patents. In 2020, he co-founded the start-up eyecan.ai.



Luigi Di Stefano received a PhD degree in electronic engineering and computer science from the University of Bologna in 1994. He is a full professor at the Department of Computer Science and Engineering, University of Bologna, where he founded and led the Computer Vision Laboratory (CVLab). His research interests include image processing, computer vision, and machine/deep learning. He is the author of more than 150 papers and several patents. He has been a scientific consultant for major computer vision and machine learning companies. He is a member of the

Pierluigi Zama Ramirez received his PhD in Computer Science and Engineering in 2021. He has been a Research Intern at Google for 6 months and is currently a Post-Doc at the University of Bologna. He co-authored more than 20 publications on computer vision research topics such as semantic segmentation, depth estimation, optical flow, domain adaptation, virtual reality, and 3D computer vision.

