

---

# DEEP LEARNING FOR 2D AND 3D ROTATABLE DATA: AN OVERVIEW OF METHODS

**Luca DELLA LIBERA, Vladimir GOLKOV,**

**Yue ZHU, Arman MIELKE, and Daniel CREMERS**

{luca.della, vladimir.golkov, yue.zhu, arman.mielke, cremers}@tum.de

Computer Vision Group

Technical University of Munich

## ABSTRACT

One of the reasons for the success of convolutional networks is their equivariance/invariance under *translations*. However, rotatable data such as molecules, living cells, everyday objects, or galaxies require processing with equivariance/invariance under *rotations* in cases where the rotation of the coordinate system does not affect the meaning of the data (e.g. object classification). On the other hand, estimation/processing of rotations is necessary in cases where rotations are important (e.g. motion estimation). There has been recent progress in methods and theory in all these regards. Here we provide an overview of existing methods, both for 2D and 3D rotations (and translations), and identify commonalities and links between them, in the hope that our insights will be useful for choosing and perfecting the methods.

## 1 INTRODUCTION

Rotational and translational equivariance play an important role in image recognition tasks. Convolutional neural networks (CNNs) are already translationally equivariant: the convolution of a translated image with a filter is equivalent to the convolution of the untranslated image with the same filter, followed by a translation. Unfortunately, standard CNNs do not have an analogous property for rotations.

A naive attempt to achieve rotational equivariance/invariance is data augmentation. While this approach can be feasible for 2D input domains, it is less so for 3D, due to the larger number of possible orientations. Methods that achieve rotational equivariance/invariance in more advanced ways have appeared recently. Our goal is to provide an overview of existing methods, both for 2D and 3D, trying to identify meaningful categories and subcategories into which they can be sorted.

Apart from methods that are equivariant under rotations of the input (i.e. where rotation “must not matter”), we also include examples of methods that can return rotations as output, as well as methods that use rotations as input and/or as deep features (i.e. where rotation matters). We do not provide all the details of the listed methods; instead, we explain the basic ideas and intuition behind them.

This work is structured as follows. In Section 2 we introduce important mathematical concepts such as equivariance and steerability. In Sections 3–4 we present the main approaches used to achieve rotational equivariance. In Section 5.1 we categorize concrete methods that use those approaches to achieve equivariance/invariance. We also categorize methods that can return a rotation as output in Section 5.2 and methods that use rotations as input and/or deep features in Section 5.3. Finally, we draw conclusions in Section 6.

The mathematical concepts (Section 2) serve as a foundation for the best (i.e. exact and most general) equivariant approach (Section 3.4). Depending on the reader’s interests, only a subset of the sections can be read:

- Exact equivariance/invariance under rotations of the input: Section 2 (optional), Section 3.4 (the rest of Sections 3 and 4 can be skipped because Section 3.4 describes the best approach for most situations), Section 5.1 (“Exact” rows of Tables 1 and 2), Section 6.
- Rotations as output: Section 5.2.
- Rotations as input or as deep features: Section 5.3.

## 2 FORMAL DEFINITIONS

In this section we define two mathematical concepts that we will use in the next sections, namely *equivariance* and *steerability*.

### 2.1 EQUIVARIANCE

**Definition 1.** A function  $f : \mathcal{X} \rightarrow \mathcal{Y}$  is **equivariant** (see Worrall et al., 2017) under a group  $\Pi$  of transformations if for each transformation  $\pi \in \Pi$  of the input of  $f$ , a transformation  $\psi \in \Psi$  of the output of  $f$  exists such that

$$f(\pi[\mathbf{x}]) = \psi[f(\mathbf{x})] \quad \forall \mathbf{x} \in \mathcal{X}. \quad (1)$$

**Definition 2.** A special case of equivariance is **same-equivariance** (see Dieleman et al., 2016), when  $\psi = \pi$ .

Note that in some sources same-equivariance is called equivariance, and what we call equivariance is called covariance.

**Definition 3.** A special case of equivariance is **invariance**, when  $\psi = \mathbb{I}$ , the identity.

**Definition 4.** Equivariance is **exact** if Eq. (1) holds strictly, **approximate** (for example approximated through learning) if Eq. (1) holds approximatively.

### 2.2 STEERABILITY

**Definition 5.** A function  $f : \mathcal{X} \rightarrow \mathcal{Y}$  is **steerable** (Freeman and Adelson, 1991) if rotated versions of  $f$  can be expressed using linear combinations of a fixed set of basis functions  $g_j$  for  $j = 1, \dots, M$ , that is:

$$f(\pi[\mathbf{x}]) = \sum_{j=1}^M k_j(\pi) g_j(\mathbf{x}), \quad (2)$$

where  $\pi \in \Pi$  is a rotation and  $k_j : \Pi \rightarrow \mathbb{C}$  for  $j = 1, \dots, M$ , are *steering factors*.

For example, if we consider a standard non-normalized 2D Gaussian  $G(x, y) = e^{-(x^2+y^2)}$ , its first derivative  $G_x(x, y) = \frac{\partial G}{\partial x}(x, y)$  in the  $x$  direction can be steered at an arbitrary orientation  $\theta$  through a linear combination of  $G_x^{0^\circ}(x, y) = G_x(x, y) = -2xe^{-(x^2+y^2)}$  and  $G_x^{90^\circ}(x, y) = G_x(\pi^{90^\circ}[x, y]) = -2ye^{-(x^2+y^2)}$ :

$$G_x^\theta(x, y) = G_x(\pi^\theta[x, y]) = \cos(\theta)G_x^{0^\circ}(x, y) + \sin(\theta)G_x^{90^\circ}(x, y). \quad (3)$$

A visualization of the case  $\theta = 30^\circ$  looks as follows:

$$\underbrace{G_x^{30^\circ}(x, y)}_{\text{Image}} = G_x(\pi^{30^\circ}[x, y]) = \underbrace{\cos(30^\circ)}_{\sim 0.87} \underbrace{G_x^{0^\circ}(x, y)}_{\text{Image}} + \underbrace{\sin(30^\circ)}_{0.5} \underbrace{G_x^{90^\circ}(x, y)}_{\text{Image}}. \quad (4)$$

A useful consequence of steerability is that convolution of an image with basis filters  $g_j$  is rotationally equivariant. The mapping  $\psi$  in Eq. (1) in this case corresponds to a linear combination of the feature maps.

As a side note, the mapping  $\psi$  can also be a certain kind of linear combinations even if the filters are not a basis of a steerable filter, see Section 3.3.

---

**Definition 6.** A function  $f : \mathcal{X} \rightarrow \mathcal{Y}$  is **self-steerable** (Liu et al., 2012) if it is steerable and the set of basis functions includes only  $f$  itself, that is:

$$f(\pi[\mathbf{x}]) = k(\pi)f(\mathbf{x}). \quad (5)$$

For example, if we consider the complex exponential  $Z(\phi) = e^{i\phi}$ , it can be self-steered at an arbitrary orientation  $\theta$  as:

$$Z^\theta(\phi) = Z(\pi^\theta[\phi]) = Z(\phi - \theta) = e^{i(\phi - \theta)} = e^{-i\theta}Z(\phi). \quad (6)$$

### 2.2.1 HARMONICS

A harmonic function  $f$  is a twice continuously differentiable function that satisfies Laplace’s equation, i.e.  $\nabla^2 f = 0$ . Circular harmonics and spherical harmonics are defined on the circle and the sphere, respectively, and are similar to the Fourier series (i.e. sinusoids with different frequencies).

2D or 3D rotational equivariance can be hardwired in the network architecture by restricting the filters’ angular component to belong to the circular harmonic or spherical harmonic family, respectively. The proof utilizes the fact that such filters are self-steerable. The radial profile of these filters on the other hand can be learned. There are various techniques to parameterize the radial profile (with learnable parameters). Also discretization requires special techniques. For example, a simple parameterization of learnable radial profiles of 2D filters uses polar coordinates, requiring resampling and bandlimiting when applied to pixel images.

### 2.3 GROUP CONVOLUTION

The group convolution (Cohen and Welling, 2016; Esteves et al., 2018a) is a natural extension of the standard convolution that allows to deal with transformation groups in a structured way.

**Definition 7.** The **group convolution** between a feature map  $\mathbf{F} : \mathcal{X} \rightarrow \mathcal{Y}$  and a filter  $\mathbf{W}$  is defined as

$$(\mathbf{F} \star_{\Pi} \mathbf{W})(\mathbf{x}) = \int_{\pi \in \Pi} \mathbf{F}(\pi[\boldsymbol{\eta}])\mathbf{W}(\pi^{-1}[\mathbf{x}]) d\pi, \quad (7)$$

where  $\Pi$  is a group,  $\pi \in \Pi$  is a transformation, and  $\boldsymbol{\eta}$  is typically a canonical element in the domain of  $\mathbf{F}$  (e.g. the origin if  $\mathcal{X} = \mathbb{R}^n$ ).

The group convolution can be shown (see Esteves et al., 2018b; Kondor and Trivedi, 2018) to be equivariant.

The ordinary convolution is a special case of the group convolution where  $\Pi$  is the additive group of  $\mathbb{R}^n$  (group of translations).

## 3 APPROACHES THAT GUARANTEE EXACT ROTATIONAL EQUIVARIANCE

In this section we list and briefly discuss the approaches used to achieve exact rotational equivariance/invariance. The state-of-the-art approach is described in Section 3.4.

### 3.1 HARDWIRED POSE NORMALIZATION

A basic approach to address the problem of rotational invariance consists in trying to “erase” the effect of rotations by reverting the input to a canonical pose. This can be achieved by *hardwiring* a reversion function at the input layer of the network architecture (e.g. PCA (Vranic, 2003)). See also Section 4.1 for *learned* pose normalization.

The advantage of hardwired pose normalization is the ease of implementation and the effectiveness in most 2D image recognition tasks. Problems are the lack of robustness against outliers of techniques such as PCA, and generally discontinuities and ambiguities of the reversion mapping.

---

### 3.2 HANDCRAFTED FEATURES

Extractors of simple rotationally invariant features can be handcrafted rather than learned. Examples (Liu et al., 2018) include features based on distances between pairs of points (SE( $n$ )-invariant), and/or between each point and the origin (invariant under rotations around the origin).

The advantage of this approach is that it guarantees invariance. The disadvantage is that it can only achieve invariance rather than more general equivariance, thus useful information about the orientation is lost. Also, this approach contradicts the philosophy of deep learning that feature extractors should be learned rather than handcrafted.

### 3.3 ORIENTED RESPONSES

A standard convolutional layer computes a discrete multi-channel convolution between a feature map  $\mathbf{F}$  and a filter  $\mathbf{W}$ . This operation is in general not rotationally equivariant. A simple way to achieve equivariance under rotations by predefined angles is to compute *oriented responses* (Zhou et al., 2017) by either 1) *rotating the filters* or 2) *rotating the feature maps* (Dieleman et al., 2016). The first approach consists in maintaining multiple rotated copies of the filter and convolving each of them with the feature map, the second in rotating the feature map multiple times by a set of fixed angles and convolving it with the filter. The result is similar to data augmentation with random rotations, but in the case of oriented responses equivariance is hardwired. The mapping  $\psi$  in Eq. (1) in this case corresponds to rotating the feature map and cycling its channels, see (Cohen and Welling, 2017, Fig. 3).

The advantage of this simple approach is the ease of implementation and the rotational weight sharing, which reduces the total number of learned parameters. The disadvantage is that since in  $n$ -dimensional Euclidian space the number of rotational degrees of freedom grows as  $\frac{1}{2}n(n-1)$ , for  $n=3$  it is already infeasible to maintain a feature map for each possible orientation if rotation angles are fine. Moreover, in most applications equivariance under rotations by arbitrary angles (rather than a fixed set of angles) is desired.

### 3.4 CONVOLUTION WITH EQUIVARIANT FILTER BANKS

Convolutions with a limited number  $n$  of filters can also achieve equivariance under rotations by arbitrary angles (not only by  $n$  angles as in the general case of oriented responses) if filter banks are constrained to a certain class. In fact, any linear equivariant transformation is equivalent to such convolutions (Kondor and Trivedi, 2018; Cohen et al., 2018a). When using so-called *non-regular representations* of the rotation group, the equivariant mapping corresponds to convolution with steerable filters, Section 2.2. When using *regular representations*, the mapping corresponds to group convolution, Eq. (7).

Note that these equivariant network layers are linear. To learn more complicated equivariant maps, certain equivariant nonlinearities are used between the linear layers. See Section 7 by Cohen et al. (2018a) for an overview of equivariant nonlinearities, equivariant batch normalization, and equivariant residual learning.

For the most common purposes, this is arguably the best of all approaches: It guarantees exact equivariance (unlike Section 4), without an obligation to discretize the transformation group (unlike Section 3.3) or to use untrainable feature extractors (unlike Section 3.2) or unstable pose normalization (unlike Section 3.1).

## 4 APPROACHES TO LEARN APPROXIMATE ROTATIONAL EQUIVARIANCE

Various approaches exist that facilitate the learning of approximated (inexact) rotational equivariance. An example approach is data augmentation by random rotations of the input. In the following, we describe learned pose normalization, soft constraints, and deformable convolution. Note that for datasets where exact equivariance is appropriate, approaches that provide exact equivariance (Section 3.4) usually work better.

---

#### 4.1 LEARNED POSE NORMALIZATION

Instead of *hardwiring* a pose normalization function as described in Section 3.1, it is possible to force or encourage the network to *learn* a reversion function directly from the training data. As an example of the “encourage” case, in spatial transformer networks (Jaderberg et al., 2015), learning a pose normalization is a facilitated but not a guaranteed side effect of learning to classify.

#### 4.2 SOFT CONSTRAINTS

Another approach to let the network learn rotational equivariance/invariance is to introduce additional *soft constraints*, which are typically expressed by auxiliary loss functions that are added to the main loss function. For example, a similarity loss (Coors et al., 2018) can be defined, which penalizes large distances between the predictions or feature embeddings of rotated copies of the input that are simultaneously fed into separate streams of a siamese network.

The advantage of this approach is the ease of implementation. Furthermore, it can be used in combination with other approaches that provide non-exact equivariance (e.g. pose normalization) in order to enhance it. The disadvantage is that equivariance/invariance is only approximative. The quality of the approximation depends on the loss formula, training data, network architecture and optimization algorithm.

#### 4.3 DEFORMABLE CONVOLUTION

In *Deformable Convolutional Networks*, Dai et al. (2017) introduce the *deformable convolution*, a new operation that augments a CNN’s capability of modeling geometric transformations. 2D offsets are added to the regular grid sampling locations in the standard convolution, enabling an input-features-dependent free-form deformation of the sampling grid.

The ordinary discrete convolution on the plane consists of 1) sampling using a regular grid  $\mathcal{R}$  over the input feature map  $\mathbf{x}$  and 2) summation of sampled values weighted by  $\mathbf{W}$ . The grid  $\mathcal{R}$  defines the receptive field size and dilation. For example,

$$\mathcal{R} = \{(-1, -1), (-1, 0), \dots, (0, 1), (1, 1)\} \quad (8)$$

defines a  $3 \times 3$  kernel with dilation 1.

For each location  $\mathbf{x}$  on the output feature map  $\mathbf{H}$ , we have

$$\mathbf{H}(\mathbf{x}) = \sum_{\mathbf{x}_n \in \mathcal{R}} \mathbf{W}(\mathbf{x}_n) \mathbf{F}(\mathbf{x} + \mathbf{x}_n), \quad (9)$$

where  $\mathbf{F}$  is the input feature map and  $\mathbf{x}_n$  enumerates the locations in  $\mathcal{R}$ .

In the deformable convolution, the regular grid  $\mathcal{R}$  is enhanced with offsets  $\Delta \mathbf{x}_n \mid n = 1, \dots, N$ , where  $N = |\mathcal{R}|$ . Equation (9) becomes

$$\mathbf{H}(\mathbf{x}) = \sum_{\mathbf{x}_n \in \mathcal{R}} \mathbf{W}(\mathbf{x}_n) \mathbf{F}(\mathbf{x} + \mathbf{x}_n + \Delta \mathbf{x}_n). \quad (10)$$

The offsets  $\Delta \mathbf{x}_n$  are learned by applying an additional convolutional layer over the same input feature map. Bilinear image interpolation of  $\mathbf{F}$  is used due to non-integer values of  $\Delta \mathbf{x}_n$  (i.e.  $\mathbf{x} + \mathbf{x}_n + \Delta \mathbf{x}_n$  off the discrete grid of  $\mathbf{F}$ ).

The advantage of this approach is that it can learn to handle very general transformations such as rotation, scaling and deformation, if training data encourage this. The disadvantage is that there is no guarantee of invariance.

Method	Input	Approach	Property	Group	Cardinality
Many	*	<b>Learned</b> (Data augmentation)	*	*	*
Transformation Equivariant Boltzmann Machine (Kivinen and Williams, 2011)	Pixel grid	Exact	Equivariance	SE(2)	<b>Discretized</b> (any angle)
Equivariant Filters and Kernel Weighted Mapping (Liu et al., 2012)	Pixel grid	Exact	Equivariance	SE(2)	Continuous
Spatial Transformer Networks (Jaderberg et al., 2015)	Pixel grid	<b>Learned</b> (Learned pose normalization)	Invariance	SE(2)	Continuous
Cyclic Symmetry in CNNs (Dieleman et al., 2016)	Pixel grid	Exact	Equivariance	SE(2)	<b>Discretized</b> (90° angles)
Group Equivariant CNNs (Cohen and Welling, 2016)	Pixel grid	Exact	Equivariance	SE(2)	<b>Discretized</b> (90° angles)
Harmonic Networks (Worrall et al., 2017)	Pixel grid	Exact	Equivariance	SE(2)	Continuous
Vector Field Networks (Marcos et al., 2017)	Pixel grid	Exact	Equivariance	SE(2)	<b>Discretized</b> (any angle)
Oriented Response Networks (Zhou et al., 2017)	Pixel grid	Exact	Equivariance	SE(2)	<b>Discretized</b> (any angle)
Deformable CNNs (Dai et al., 2017)	Pixel grid	<b>Learned</b> (Deformable convolution)	Equivariance	SE(2)	Continuous
Polar Transformer Networks (Esteves et al., 2018b)	Pixel grid <sup>a</sup>	<b>Learned</b> (Learned pose normalization)	Equivariance	SE(2)	Continuous
Steerable Filter CNNs (Weiler et al., 2018b)	Pixel grid	Exact	Equivariance	SE(2)	<b>Discretized</b> (any angle)
Learning Invariance with Weak Supervision (Coors et al., 2018)	Pixel grid	<b>Learned</b> (Soft constraints)	Invariance	SE(2)	Continuous
Rototranslation Covariant CNNs (Bekkers et al., 2018)	Pixel grid	Exact	Invariance	SE(2)	<b>Discretized</b> (any angle)
RotDCF: Decomposition of Convolutional Filters (Cheng et al., 2019)	Pixel grid	Exact	Equivariance	SE(2)	<b>Discretized</b> (any angle)
Siamese Equivariant Embedding (Véges et al., 2018)	Pixel grid	<b>Learned</b> (Soft constraints)	Equivariance	SO(2)	Continuous
CNN model of primary visual cortex (Ecker et al., 2019)	Pixel grid	Exact	Equivariance	SE(2)	<b>Discretized</b> (any angle)

<sup>a</sup> In Polar Transformer Networks, the pixel grid is transformed to a circular signal in an intermediate layer.

Table 1: Methods with equivariance under 2D rotations, classified according to our taxonomy. An overview of the terminology used in the table is given in Section 5. Methods with identical cell entries differ in terms of details. Potential weaknesses are **highlighted**. Harmonic Networks are the “best” neural networks in that they provide continuous exact equivariance. See also Cohen et al. (2018a).

Method	Input	Approach	Property	Group	Cardinality
Many	*	<b>Learned</b> (Data augmentation)	*	*	*
Spherical Harmonic Representation (Kazhdan et al., 2003)	Voxel grid	Exact	Invariance	SO(3)	Continuous
Equivariant Filters and Kernel Weighted Mapping (Liu et al., 2012)	Voxel grid	Exact	Equivariance	SE(3)	Continuous
Spatial Transformer Networks (Jaderberg et al., 2015)	Voxel grid	<b>Learned</b> (Learned pose normalization)	Invariance	SE(3)	Continuous
SO(3) Equivariant Representations (Esteves et al., 2018a)	Spherical signal	Exact	Equivariance	SO(3)	Continuous
Spherical CNNs (Cohen et al., 2018b)	Spherical signal	Exact	Equivariance	SO(3)	Continuous
Tensor Field Networks (Thomas et al., 2018)	Point cloud	Exact	Equivariance	SE(3)	Continuous
N-body Networks (Kondor, 2018)	Point cloud	Exact	Equivariance	SO(3)	Continuous
CubeNet (Worrall and Brostow, 2018)	Voxel grid	Exact	Equivariance	SE(3)	<b>Discretized</b> (90° angles)
3D Steerable CNNs (Weiler et al., 2018a)	Voxel grid	Exact	Equivariance	SE(3)	Continuous
PPF-FoldNet (Deng et al., 2018)	Point cloud	Exact ( <b>hand-crafted features</b> )	Invariance	SE(3)	Continuous

Table 2: Methods with equivariance under 3D rotations, classified according to our taxonomy. An overview of the terminology used in the table is given in Section 5. Potential weaknesses are **highlighted**. Tensor Field Networks are the “best” for point clouds in that they provide continuous exact SE(3)-equivariance. Similarly, 3D Steerable CNNs are the “best” neural networks for voxel grids. All these 3D deep learning methods appeared in 2018.

## 5 OVERVIEW OF METHODS

In this section we list and categorize deep learning methods for handling rotatable data and rotations. This includes methods that are equivariant under rotations of the input, methods that output a rotation, and methods that use rotations as inputs and/or deep features.

### 5.1 EQUIVARIANCE UNDER ROTATIONS OF THE INPUT

Methods that are equivariant under rotations of the input are categorized in Table 1 (for 2D rotations) and Table 2 (for 3D rotations) according to the following criteria:

- Input:
  - Pixel grid: a grid representation of 2D image data
  - Voxel grid: a grid representation of 3D volumetric data
  - Point cloud: a set of 3D point coordinates
  - Multi-view images: multiple images of the same object, taken from different viewpoints
  - Circular signal: a function defined on the circle
  - Spherical signal: a function defined on the sphere
- Approach: see Definition 4 and Sections 3–4
- Property: equivariance (Definition 1) or invariance (Definition 3)
- Group:

- 
- $SO(2)$ : the group of 2D rotations
  - $SE(2)$ : the group of 2D rigid-body motions
  - $SO(3)$ : the group of 3D rotations
  - $SE(3)$ : the group of 3D rigid-body motions
  - Cardinality: continuous (entire group) or discretized to specific angles

## 5.2 ROTATIONS AS OUTPUT

Examples of deep learning methods that output a (3D) rotation are categorized in Table 3 and Table 4 according to the following characteristics:

- Input to the network that outputs the rotation, and according rotation-prediction task:
  - Image. The task is to estimate the orientation of a depicted object relative to the camera.
  - Cropped image: Image cropped to the relevant object, usually automatically.
  - Cropped stereo image: A pair of images is taken at the same time from two cameras that are close together and point in the same direction. The images are cropped in a predetermined fashion. Each pair of cropped images constitutes one input. The position of the cameras relative to each other is fixed. The task is to estimate the orientation of a depicted object relative to the cameras.
  - Volumetric data. The task is to estimate the orientation of an object relative to volume coordinates.
  - Slice of volumetric data. The task is to estimate the orientation of a 2D slice relative to an entire (predefined) 3D object.
  - Video: Two or three or more images (video frames). The task is to estimate the relative rotation and translation of the camera between (not necessarily consecutive) frames.
- Number of objects/rotations: Describes how many rotations the network outputs.
  - One rigid object: One rotation is estimated that is associated with a rigid object. The visible “object” is the entire scene in cases where camera motion relative to a static scene is estimated. Other small moving objects can be additionally accounted for (for example to refine the optical-flow estimation), but their *rotation* is *not* estimated.
  - Hierarchy of object parts: Relative rotations between the objects and their parts are estimated, with several hierarchy levels, i.e. an “object” consisting of parts can itself be one of several parts of a “higher-level” object.
- Specialized on specific object(s)?
  - Specialized on one object: The network can only process one type of object on which it was trained.
  - Specialized on multiple objects: The network can process an arbitrarily large but fixed set of object types on which it was trained.
  - Generalizing to new objects: The network can generalize to unseen types of objects.
- Group:  $SO(3)$  or  $SE(3)$
- Representation (embedding) of the rotation(s):
  - Rotation matrix
  - Quaternion
  - Euler angles
  - Axis-angle representation
  - Discrete bins: Rotations are grouped into a finite set of bins.
  - Transformation matrix



- 
- 3D coordinates of four keypoints on the object (predefined object-specifically, e.g. four of its corners)
  - Eight corners and centroid of 3D bounding box projected into 2D image space
  - Learned representation: The latent space (e.g. of an autoencoder) is used to represent the rotation. There are several interesting aspects at play:
    - \* Learned representations allow for ambiguity: If an object looks very similar from two angles and the loss allows for it, the network can learn to use the same encoding to represent both rotations. On the other hand, unambiguous representations (like the ones listed above) would require generative/probabilistic models to deal with ambiguity.
    - \* Certain representations are encouraged due to the overall network architecture. For example, in capsule networks, learned representations of rotation are processed in a very specific way (multiplied by learned transformation matrices).
    - \* Features other than rotation might be entangled into the learned representation. This is not even always discouraged. For example, in capsule networks, the learned representation may also contain other object features such as color.
  - Loss function. We distinguish the following categories:
    - Rotations are estimated at the output layer. The loss measures the similarity to ground truth rotations of training samples. These methods are listed in Table 3.
      - \* Geodesic distance between prediction and ground truth. This loss is rotationally invariant, i.e. the network miscalculating a rotation by  $10^\circ$  always results in the same loss value, regardless of the ground truth rotation and of the direction into which the prediction is biased.
      - \*  $L^p$  distance in embedding space: This loss value is fast to compute but not rotationally invariant, i.e. an error of  $10^\circ$  yields different loss values depending on the ground truth and on the prediction. Due to this “unfairness”/“arbitrariness”, such losses are **highlighted** in the table.
    - Rotations are estimated in an intermediate layer and used in subsequent layers for a “higher-level” goal of a larger system. Ground truth rotations are not required. These methods are listed in Table 4.
      - \* Object classification: Rotation prediction is trained as part of a larger system for object classification. The estimated rotation is used to rotate the input or feature map (in the case of spatial transformer networks) or predicted poses (in the case of capsule networks) as an intermediate processing step. It is assumed that learning to rotate to a canonical pose (in the case of spatial transformer networks) or to let object parts vote about the overall object pose (in the case of capsule networks) is beneficial for object classification. The estimation of rotation is incidental and encouraged by the overall setup. However, its approximate correctness is not necessary for perfect object classification. Therefore, the “predicted rotation” can be very wrong, and due to this danger this loss is **highlighted** in the table.
      - \* View warping: At least two video frames are used to estimate the scene geometry (depth maps) and camera motion between the views (rotation, translation). These estimates are used to warp one view (image, and possibly depth map) to resemble another view. The loss measures this resemblance. This is a form of self-supervised learning: ground truth geometry and motion are not given, but are estimated such that they cause warping that is consistent with the input images. The rotation estimation can be expected to be good, because it is necessary for good view synthesis.
      - \* Autoencoder reconstruction loss: The network is trained to reconstruct its input (a view of the object) after passing it through a lower-dimensional latent space. The output target has a neutral image background and lacks other objects that were visible in the input image. This allows the network to learn to discard the information about the background and other objects

Method	Input to the network that outputs rotation	Number of objects/rotations	Specialized on specific object(s)?	Group	Representation (embedding) of the rotation(s)	Loss function
PoseNet (Kendall et al., 2015)	Image	One rigid object	Specialized on one object	SE(3)	Quaternion	$L^2$ distance in embedding space
Relative Camera Pose Estimation Using CNNs (Melekhov et al., 2017)	Video (two non-consecutive frames)	One rigid object	Generalizing to new objects	SE(3)	Quaternion	$L^2$ distance in embedding space
3D Pose Regression using CNNs and axis-angle representation (Mahendran et al., 2017)	Image	One rigid object	Specialized on multiple objects	SO(3)	Axis-angle representation	Geodesic distance
3D Pose Regression using CNNs and quaternions (Mahendran et al., 2017)	Image	One rigid object	Specialized on multiple objects	SO(3)	Quaternion	Geodesic distance
Real-Time Seamless Single Shot 6D Object Pose Prediction (Tekin et al., 2017)	Image	One rigid object	Specialized on multiple objects	SE(3)	Eight corners and centroid of 3D bounding box projected into 2D image space	Squared $L^2$ distance in embedding space
Registration of a slice to a predefined volume (Salehi et al., 2018)	Slice of volumetric data	One rigid object	Specialized on one object	SE(3)	Axis-angle representation	Geodesic distance <sup>a</sup>
Registration of a volume to another, predefined volume (Salehi et al., 2018)	Volumetric data	One rigid object	Specialized on one object	SE(3)	Axis-angle representation	Geodesic distance <sup>a</sup>
SSD-AF-3D-AxisBinned (Pandey et al., 2018)	Cropped stereo image	One rigid object	Specialized on multiple objects	SE(3)	Discrete bins	Smoothed $L^1$ distance in embedding space
SSD-AF-MultiplePoint (Pandey et al., 2018)	Cropped stereo image	One rigid object	Specialized on multiple objects	SE(3)	Four keypoint locations in 3D space	Smoothed $L^1$ distance in embedding space
SSD-AF-Stereo6D-Quat (Pandey et al., 2018)	Cropped stereo image	One rigid object	Specialized on multiple objects	SE(3)	Quaternion	Smoothed $L^1$ distance in embedding space
SSD-AF-Stereo6D-Euler (Pandey et al., 2018)	Cropped stereo image	One rigid object	Specialized on multiple objects	SE(3)	Euler angles	Smoothed $L^1$ distance in embedding space
Learning Local RGB-to-CAD Correspondences for Object Pose Estimation (Georgakis et al., 2018)	Image and 3D model	One rigid object	Specialized on multiple objects	SE(3)	Rotation matrix	Squared $L^2$ distance in embedding space

<sup>a</sup> Initially squared  $L^2$  distance in embedding space (fast to compute); then geodesic distance for rotation and squared  $L^2$  distance for translation.

Table 3: Examples of deep learning methods that can output a 3D rotation, where a ground truth rotation is used for training. An overview of the terminology used in the table is given in Section 5.2. Losses that lack rotational invariance are highlighted.

Method	Input to the network that outputs rotation	Number of objects/rotations	Specialized on specific object(s)?	Group	Representation (embedding) of the rotation(s)	Loss function
Capsule Networks (Sabour et al., 2017)	Image	Hierarchy of object parts	Generalizing to new objects	SE(3)	Poses of lowest-level object parts: learned representation; part-to-object pose transformations: transformation matrices (as trainable parameters)	Object classification
Spatial Transformer Networks (Jaderberg et al., 2015)	Volumetric data	One rigid object	Generalizing to new objects	SE(3)	Transformation matrix	Object classification
Unsupervised Learning of Depth and Ego-Motion (Zhou et al., 2017)	Video (three consecutive frames)	One rigid object	Generalizing to new objects	SE(3)	Euler angles	View warping
Learning Implicit Representations of 3D Object Orientations from RGB (Sundermeyer et al., 2018)	Image	One rigid object	Specialized on one object	SO(3)	Learned representation	Autoencoder reconstruction loss
GeoNet (Yin and Shi, 2018)	Video (several consecutive frames)	One rigid object	Generalizing to new objects	SE(3)	Euler angles	View warping

Table 4: Examples of deep learning methods that can output a 3D rotation, where a ground truth rotation is *not* necessary for training. An overview of the terminology used in the table is given in Section 5.2. Losses are highlighted for which a good prediction of rotations is not necessary for a good overall loss value.

before the bottleneck layer. If the network is specialized on one object, then maintaining in the latent space only the information about the object pose is sufficient for such reconstruction. If additionally the latent space is sufficiently low-dimensional, then the learning is encouraged to be economic about the amount of information encoded in the latent space, i.e. to encode nothing but the pose. Thus, a learned representation of rotation emerges.

Estimation of 2D rotations is similar in almost all listed regards, but simpler in terms of representation. Predicting the sine and cosine of the rotation angle (and normalizing the predicted vector to length 1, because otherwise the predicted sine and cosine might slightly contradict each other, or be beyond  $[-1, 1]$ ) is better than predicting the angle, because the latter requires learning a function that has a jump (from  $360^\circ$  to  $0^\circ$ ), which is not easy for (non-generative) neural networks.

### 5.3 ROTATIONS AS INPUT OR AS DEEP FEATURES

Other uses of rotations in deep learning are to take rotations as input, or to restrict deep features to belong to  $SO(n)$  (without requiring them to directly approximate rotations present in the data). For example, Huang et al. (2017) use rotation matrices as inputs and as deep features. They restrict deep features to  $SO(3)$  by using layers that map from  $SO(3)$  to  $SO(3)$ .

---

## 6 CONCLUSIONS

In this work we gathered the current knowledge about deep learning for rotatable data and organized it in a structured way. After summarizing the approaches to achieve rotational equivariance/invariance, we provided a list of methods that use those approaches, and we classified them according to our taxonomy. Furthermore, we categorized some methods that can return a rotation as output.

Among the methods for equivariance/invariance, the most successful ones are based on the exact and most general approach (Section 3.4). They are very effective in 3D input domains as well.

With emerging theory (Kondor and Trivedi, 2018; Cohen et al., 2018a) for exact equivariance and with emerging approaches, it appears to be the perfect time to use the methods in various application domains and to tune them. An example of tuning general ideas to very advanced application-specific properties is given by Anderson et al. (2019). Solutions in many other domains of application can be expected to be less complex than this quantum-mechanics-based example.

Existing pipelines that do not have exact equivariance yet and for example rely on data augmentation are likely to benefit from incorporating exact-equivariance approaches.

## ACKNOWLEDGMENTS

We thank Erik Bekkers and Remco Duits for valuable comments on an early draft of this manuscript, and Antonij Golkov, Christine Allen-Blanchette, and Taco Cohen for mathematical discussions. This manuscript was supported by the ERC Consolidator Grant “3DReloaded”.

## REFERENCES

- Anderson, B. M., Hy, T., and Kondor, R. (2019). Cormorant: Covariant molecular neural networks. *CoRR*, abs/1906.04015. (^12)
- Bekkers, E. J., Lafarge, M. W., Veta, M., Eppenhof, K. A. J., Pluim, J. P. W., and Duits, R. (2018). Roto-translation covariant convolutional networks for medical image analysis. *CoRR*, abs/1804.03393. (^6)
- Cheng, X., Qiu, Q., Calderbank, R., and Sapiro, G. (2019). RotDCF: decomposition of convolutional filters for rotation-equivariant deep networks. In *International Conference on Learning Representations*. (^6)
- Cohen, T., Geiger, M., and Weiler, M. (2018a). A general theory of equivariant CNNs on homogeneous spaces. *CoRR*, abs/1811.02017. (^4, 6, 12)
- Cohen, T. S., Geiger, M., Köhler, J., and Welling, M. (2018b). Spherical CNNs. In *International Conference on Learning Representations*. (^7)
- Cohen, T. S. and Welling, M. (2016). Group equivariant convolutional networks. In *Proceedings of the 33rd International Conference on International Conference on Machine Learning - Volume 48, ICML’16*, pages 2990–2999. JMLR.org. (^3, 6)
- Cohen, T. S. and Welling, M. (2017). Steerable cnns. In *5th International Conference on Learning Representations, ICLR 2017, Toulon, France, April 24-26, 2017, Conference Track Proceedings*. (^4)
- Coors, B., Condurache, A., Mertins, A., and Geiger, A. (2018). Learning transformation invariant representations with weak supervision. In *International Conference on Computer Vision Theory and Applications*. (^5, 6)
- Dai, J., Qi, H., Xiong, Y., Li, Y., Zhang, G., Hu, H., and Wei, Y. (2017). Deformable convolutional networks. In *2017 IEEE International Conference on Computer Vision (ICCV)*, pages 764–773. (^5, 6)

- 
- Deng, H., Birdal, T., and Ilic, S. (2018). PPF-FoldNet: unsupervised learning of rotation invariant 3D local descriptors. In *European Conference on Computer Vision (ECCV)*. Springer. (^7)
- Dieleman, S., De Fauw, J., and Kavukcuoglu, K. (2016). Exploiting cyclic symmetry in convolutional neural networks. In *Proceedings of the 33rd International Conference on International Conference on Machine Learning - Volume 48, ICML'16*, pages 1889–1898. JMLR.org. (^2, 4, 6)
- Ecker, A. S., Sinz, F. H., Froudarakis, E., Fahey, P. G., Cadena, S. A., Walker, E. Y., Cobos, E., Reimer, J., Tolia, A. S., and Bethge, M. (2019). A rotation-equivariant convolutional neural network model of primary visual cortex. In *International Conference on Learning Representations*. (^6)
- Esteves, C., Allen-Blanchette, C., Makadia, A., and Daniilidis, K. (2018a). Learning SO(3) equivariant representations with spherical CNNs. In Ferrari, V., Hebert, M., Sminchisescu, C., and Weiss, Y., editors, *Computer Vision – ECCV 2018*, pages 54–70, Cham. Springer International Publishing. (^3, 7)
- Esteves, C., Allen-Blanchette, C., Zhou, X., and Daniilidis, K. (2018b). Polar transformer networks. In *International Conference on Learning Representations*. (^3, 6)
- Freeman, W. T. and Adelson, E. H. (1991). The design and use of steerable filters. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 13(9):891–906. (^2)
- Georgakis, G., Karanam, S., Wu, Z., and Kosecka, J. (2018). Learning local RGB-to-CAD correspondences for object pose estimation. *CoRR*, abs/1811.07249. (^10)
- Huang, Z., Wan, C., Probst, T., and Gool, L. V. (2017). Deep learning on Lie groups for skeleton-based action recognition. In *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 1243–1252. (^11)
- Jaderberg, M., Simonyan, K., Zisserman, A., and Kavukcuoglu, K. (2015). Spatial transformer networks. In *Proceedings of the 28th International Conference on Neural Information Processing Systems - Volume 2, NIPS'15*, pages 2017–2025, Cambridge, MA, USA. MIT Press. (^5, 6, 7, 11)
- Kazhdan, M., Funkhouser, T., and Rusinkiewicz, S. (2003). Rotation invariant spherical harmonic representation of 3D shape descriptors. In *Proceedings of the 2003 Eurographics/ACM SIGGRAPH Symposium on Geometry Processing, SGP '03*, pages 156–164, Aire-la-Ville, Switzerland, Switzerland. Eurographics Association. (^7)
- Kendall, A., Grimes, M., and Cipolla, R. (2015). PoseNet: Convolutional networks for real-time 6-DOF camera relocalization. *CoRR*, abs/1505.07427. (^10)
- Kivinen, J. J. and Williams, C. K. I. (2011). Transformation equivariant Boltzmann machines. In Honkela, T., Duch, W., Girolami, M., and Kaski, S., editors, *Artificial Neural Networks and Machine Learning – ICANN 2011*, pages 1–9, Berlin, Heidelberg. Springer Berlin Heidelberg. (^6)
- Kondor, R. (2018). N-body networks: a covariant hierarchical neural network architecture for learning atomic potentials. *CoRR*, abs/1803.01588. (^7)
- Kondor, R. and Trivedi, S. (2018). On the generalization of equivariance and convolution in neural networks to the action of compact groups. *ArXiv*, abs/1802.03690. (^3, 4, 12)
- Liu, K., Wang, Q., Driever, W., and Ronneberger, O. (2012). 2D/3D rotation-invariant detection using equivariant filters and kernel weighted mapping. In *2012 IEEE Conference on Computer Vision and Pattern Recognition*, pages 917–924. (^3, 6, 7)
- Liu, Y., Wang, C., Song, Z., and Wang, M. (2018). Efficient global point cloud registration by matching rotation invariant features through translation search. In *The European Conference on Computer Vision (ECCV)*. (^4)

- 
- Mahendran, S., Ali, H., and Vidal, R. (2017). 3D pose regression using convolutional neural networks. *CoRR*, abs/1708.05628. (^10)
- Marcos, D., Volpi, M., Komodakis, N., and Tuia, D. (2017). Rotation equivariant vector field networks. In *2017 IEEE International Conference on Computer Vision (ICCV)*, pages 5058–5067. (^6)
- Melekhov, I., Kannala, J., and Rahtu, E. (2017). Relative camera pose estimation using convolutional neural networks. *CoRR*, abs/1702.01381. (^10)
- Pandey, R., Pidlypenskyi, P., Yang, S., and Kaeser-Chen, C. (2018). Efficient 6-DoF tracking of handheld objects from an egocentric viewpoint. *CoRR*, abs/1804.05870. (^10)
- Sabour, S., Frosst, N., and Hinton, G. E. (2017). Dynamic routing between capsules. In *Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, 4-9 December 2017, Long Beach, CA, USA*, pages 3856–3866. (^11)
- Salehi, S. S. M., Khan, S., Erdogmus, D., and Gholipour, A. (2018). Real-time deep pose estimation with geodesic loss for image-to-template rigid registration. *CoRR*, abs/1803.05982. (^10)
- Sundermeyer, M., Puang, E. Y., Marton, Z.-C., Durner, M., and Triebel, R. (2018). Learning implicit representations of 3D object orientations from RGB. In *ICRA 2018 Workshop on Representing a Complex World*. (^11)
- Tekin, B., Sinha, S. N., and Fua, P. (2017). Real-time seamless single shot 6D object pose prediction. *CoRR*, abs/1711.08848. (^10)
- Thomas, N., Smidt, T., Kearnes, S. M., Yang, L., Li, L., Kohlhoff, K., and Riley, P. (2018). Tensor field networks: rotation- and translation-equivariant neural networks for 3D point clouds. *CoRR*, abs/1802.08219. (^7)
- Véges, M., Varga, V., and Lörincz, A. (2018). 3D human pose estimation with siamese equivariant embedding. *Neurocomputing*, 339:194–201. (^6)
- Vranic, D. V. (2003). An improvement of rotation invariant 3D-shape based on functions on concentric spheres. In *Proceedings 2003 International Conference on Image Processing (Cat. No.03CH37429)*, volume 3, pages III–757. (^3)
- Weiler, M., Geiger, M., Welling, M., Boomsma, W., and Cohen, T. (2018a). 3D steerable CNNs: learning rotationally equivariant features in volumetric data. In Bengio, S., Wallach, H., Larochelle, H., Grauman, K., Cesa-Bianchi, N., and Garnett, R., editors, *Advances in Neural Information Processing Systems 31*, pages 10381–10392. Curran Associates, Inc. (^7)
- Weiler, M., Hamprecht, F. A., and Storath, M. (2018b). Learning steerable filters for rotation equivariant CNNs. *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 849–858. (^6)
- Worrall, D. and Brostow, G. (2018). CubeNet: equivariance to 3D rotation and translation. In *The European Conference on Computer Vision (ECCV)*. (^7)
- Worrall, D. E., Garbin, S. J., Turmukhambetov, D., and Brostow, G. J. (2017). Harmonic networks: deep translation and rotation equivariance. In *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 7168–7177. (^2, 6)
- Yin, Z. and Shi, J. (2018). GeoNet: Unsupervised learning of dense depth, optical flow and camera pose. *CoRR*, abs/1803.02276. (^11)
- Zhou, T., Brown, M., Snavely, N., and Lowe, D. G. (2017). Unsupervised learning of depth and ego-motion from video. *CoRR*, abs/1704.07813. (^11)
- Zhou, Y., Ye, Q., Qiu, Q., and Jiao, J. (2017). Oriented response networks. In *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 4961–4970. (^4, 6)