

# Canonical Capsules: Unsupervised Capsules in Canonical Pose

Weiwei Sun<sup>1,5,\*</sup>, Andrea Tagliasacchi<sup>3,4,\*</sup>, Boyang Deng<sup>4</sup>, Sara Sabour<sup>3,4</sup>,  
Soroosh Yazdani<sup>4</sup>, Geoffrey Hinton<sup>3,4</sup>, Kwang Moo Yi<sup>1,5</sup>

<sup>1</sup>University of British Columbia, <sup>3</sup>University of Toronto,  
<sup>4</sup>Google Research, <sup>5</sup>University of Victoria, \*equal contributions

## Abstract

We propose an unsupervised capsule architecture for 3D point clouds. We compute capsule decompositions of objects through permutation-equivariant attention, and self-supervise the process by training with pairs of randomly rotated objects. Our key idea is to aggregate the attention masks into semantic keypoints, and use these to supervise a decomposition that satisfies the capsule invariance/equivariance properties. This not only enables the training of a semantically consistent decomposition, but also allows us to learn a canonicalization operation that enables object-centric reasoning. In doing so, we require neither classification labels nor manually-aligned training datasets to train. Yet, by learning an object-centric representation in an unsupervised manner, our method outperforms the state-of-the-art on 3D point cloud reconstruction, registration, and unsupervised classification. We will release the code and dataset to reproduce our results as soon as the paper is published.

## 1. Introduction

Understanding objects is one of the core problems of computer vision [31, 15, 35]. While this task has traditionally relied on large annotated datasets [38, 21], unsupervised approaches [7] have emerged to remove the need for labels. Recently, researchers have attempted to extend these methods to work on 3D point clouds [53], but the field of unsupervised 3D learning remains relatively uncharted. Conversely, researchers have been extensively investigating 3D deep representations for shape auto-encoding (also at times referred to as reconstruction) [55, 19, 32, 17], making one wonder whether these discoveries can now benefit from unsupervised learning for tasks *other* than auto-encoding.

Importantly, these recent methods for 3D deep representation learning are not entirely unsupervised. Whether using point clouds [55], meshes [19], or implicits [32], they owe

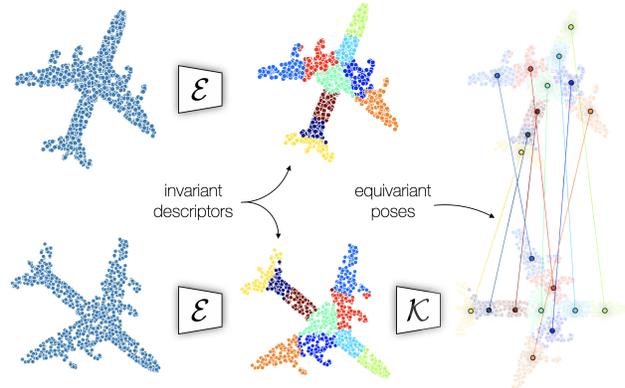


Figure 1. **Teaser** – We design a capsule architecture  $\mathcal{E}$  that is trained in an unsupervised fashion – by only observing *pairs* of randomly rotated 3D point clouds (of the *same* object). The architecture builds a  $K$ -fold decomposition that estimates primary capsules whose descriptors are *invariant* to rigid transformations, and whose poses – i.e. keypoints that summarize the location of a part – are transformation *equivariant*. By relating capsule poses to capsule descriptors, the network  $\mathcal{K}$  introduces the concept of a “mental picture” of an object: an object-centric coordinate frame for the object that empowers downstream tasks such as auto-encoding and classification.

their success to the *inductive bias of the training dataset*. Specifically, all 3D models in the popular ShapeNet [5] dataset are “object-centric” – they are pre-canonicalized to a unit bounding box, and, even more importantly, with an orientation that synchronizes object semantics to Euclidean frame axes (e.g. airplane cockpit is always along  $+y$ , car wheels always touch  $z = 0$ ). Differentiable 3D decoders are heavily affected by the consistent alignment of their output with an Euclidean frame [10, 17]: local-to-global transformations *cannot* be easily learnt by fully connected layers. Thus, these methods fail in the absence of pre-alignment.

With our method, *Canonical Capsules*, we tackle the problem of training 3D deep representation in a truly unsupervised fashion through a capsule formulation [22]. In

capsule networks, a scene is perceived via its decomposition into *part* hierarchies, and each part is represented with a (pose, descriptor) pair: ① The capsule *pose* specifies the frame of reference of a part, and hence should be transformation equivariant; ② The capsule *descriptor* specifies the appearance of a part, and hence should be transformation invariant. In our evaluation we show how these representations are effective in a variety of tasks, but “*can we train them effectively in an unsupervised setting?*”

In *Canonical Capsules*, we propose a novel architecture to compute a K-part decomposition of a point cloud via an attention mechanism; see Figure 1. Our network is trained by feeding pairs of a randomly rotated copies of the same shape. We aggregate the K parts into K keypoints, so that keypoints across shapes are in semantic one-to-one correspondence. Equivariance is then enforced by requiring the two keypoint sets to only differ by the known (relative) transformation; to realize invariance, we simply ask that the descriptors of the two instances match. Note that we train such a decomposition in a fully unsupervised fashion, and that the network only ever sees randomly rotated point clouds.

In *Canonical Capsules*, we exploit our decomposition to recover a canonical frame that allows unsupervised “object-centric” learning of 3D deep representations *without* requiring a semantically aligned dataset. We achieve this task by regressing *canonical* capsule poses from capsule descriptors via a deep network, and computing a canonicalizing transformation by solving, yet again, a shape-matching problem. This not only allows more effective shape auto-encoding, but our experiments confirm this results in a latent representation that is more effective in unsupervised classification tasks. Note that, like our decomposition, our canonicalizing transformations are also learnt in a fully unsupervised fashion, by only training on randomly rotated point clouds.

**Contributions.** In summary, in this paper we:

- propose an architecture for unsupervised learning with 3D point clouds based on capsules;
- demonstrate how capsule decompositions can be learnt via straightforward transformation augmentation;
- enable unsupervised learning to be object-centric by introducing a learned canonical frame of reference;
- achieve state-of-the-art performance in unsupervised 3D point cloud registration, auto-encoding/reconstruction, and unsupervised classification.

## 2. Related works

Our technique proposes a *capsule* architecture for learning of 3D representations that are usable across a range of unsupervised tasks: from classification [55] to *reconstruction* and *registration*.

**Capsule Networks.** Convolutional Neural Networks lack equivariance to rigid transformations, despite their pivotal role in describing the structure of the 3D scene behind a 2D image. One promising approach to overcome this shortcoming is to add equivariance under a group action in each layer. In [44] an  $SE(3)$ -equivariant network is introduced by interpreting each layer as a sum of  $SO(3)$  representations and convolutions as tensor product of different representations. To the best of our knowledge, these models have not been extended to more general rigid transformations. In our work, we remove the need for a global  $SE(3)$ -equivariant network by canonicalizing the input.

Capsule Networks [22], on the other hand, have been proposed to overcome this issue towards a relational and hierarchical understanding of natural images. Techniques such as Dynamic Routing [37, 47] and EM-algorithms [24] have been proposed as potential architectures, and have found applications ranging from medical imaging [1] to language understanding [57]. Of particular relevance to our work, are methods that apply capsule networks to 3D input data [58, 59, 41], but note these methods are not unsupervised, as they either rely on classification supervision [59], or on datasets that present a significant inductive bias in the form of pre-alignment [58]. In this paper, we take inspiration from the recent Stacked Capsule Auto-Encoders [29], where the authors have shown how capsule-style reasoning is effective as far as *primary* capsules can be trained in an unsupervised fashion – by only using unsupervised reconstruction losses. The natural question, which we answer in this paper, is “*how can we engineer networks that generate primary 3D capsules in an unsupervised fashion?*”

**Deep 3D representations.** Reconstructing 3D objects requires effective inductive biases about 3D vision *and* 3D geometry. When the input is images, the core challenge is how to encode *3D projective geometry* concepts into the model. This can be achieved by explicitly modeling multi-view geometry [27], by attempting to learn it [14], or by hybrid solutions [54]. But even when input is 3D, there are still significant challenges. It is still not clear which is the *3D representation* that is most amenable to deep learning. Researchers proposed the use of meshes [49, 30], voxels [50, 51], surface patches [19, 13, 11], and implicit functions [33, 32, 9]. Unfortunately, the importance of geometric structures (i.e. *part-to-whole* relationships) is often overlooked. Recent works have tried to close this gap by using part decomposition consisting of oriented boxes [45], ellipsoids [18, 17], convex polytopes [12], and grids [4]. However, as previously discussed, most of these still heavily rely on a pre-aligned training dataset; our paper attempts to bridge this gap, allowing learning of *structured* 3D representations *without* requiring pre-aligned data.

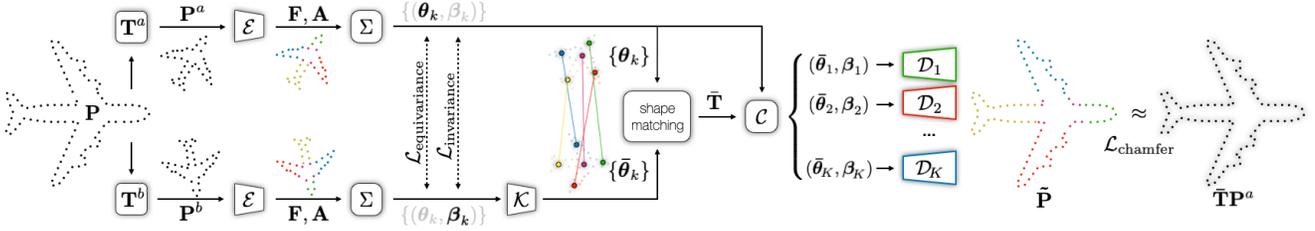


Figure 2. **Framework** – We learn a capsule encoder for 3D point clouds by relating the decomposition result of two random rigid transformations  $\mathbf{T}^a$  and  $\mathbf{T}^b$ , of a given point cloud, *i.e.*, a Siamese training setup. We learn the parameters of an encoder  $\mathcal{E}$ , a per-capsule decoder  $\mathcal{D}_k$ , as well as a network that represents a learnt canonical frame  $\mathcal{K}$ . For illustrative purposes, we shade-out the outputs that do not flow forward, and with  $\Sigma$  summarize the aggregations in (2).

**Registration.** One way to circumvent the requirement of pre-aligned datasets is to rely on methods capable of registering a point cloud into a canonical frame. The recently proposed CaSPR [36] fulfills this premise, but requires ground-truth canonical point clouds in the form of normalized object coordinate spaces [48] for supervision. Similarly, [20] regresses each view’s pose relative to the canonical pose, but still requires weak annotations in the form of multiple partial views. In contrast to these methods, our solution is completely unsupervised. Several registration techniques based on deep learning have been proposed [52, 56], even using semantic keypoints and symmetry to perform the task [16]. These methods typically register a *pair* of instances from the same class, but lack the ability to consistently register *all* instances to a shared canonical frame.

### 3. Method

Our network trains on unaligned point clouds as illustrated in Figure 2: we train a network that *decomposes* point clouds into parts, and enforce invariance/equivariance through a Siamese training setup [43]. We then *canonicalize* the point cloud to a learnt frame of reference, and perform *auto-encoding* in this coordinate space. The losses employed to train  $\mathcal{E}$ ,  $\mathcal{K}$ , and  $\mathcal{D}$ , will be covered in Section 3.1, while the details of their architecture are in Section 3.2.

**Decomposition.** In more detail, given a point cloud  $\mathbf{P} \in \mathbb{R}^{P \times D}$  of  $P$  points in  $D$  dimensions, we perturb it with two random transformations  $\mathbf{T}^a, \mathbf{T}^b \in \mathbf{SE}(D)$  to produce point clouds  $\mathbf{P}^a, \mathbf{P}^b$ . We then use a shared permutation-equivariant capsule encoder  $\mathcal{E}$  to compute a  $K$ -fold attention map  $\mathbf{A} \in \mathbb{R}^{P \times K}$  for  $K$  capsules, as well as per-point feature map  $\mathbf{F} \in \mathbb{R}^{P \times C}$  with  $C$  channels:

$$\mathbf{A}, \mathbf{F} = \mathcal{E}(\mathbf{P}), \quad (1)$$

where we drop the superscript indexing the Siamese branch for simplicity. From these attention masks, we then compute, for the  $k$ -th capsule its pose  $\theta_k \in \mathbb{R}^3$  parameterized by its location in 3D space, and the corresponding capsule

descriptor  $\beta_k \in \mathbb{R}^C$ :

$$\theta_k = \frac{\sum_p \mathbf{A}_{p,k} \mathbf{P}_p}{\sum_p \mathbf{A}_{p,k}}, \quad \beta_k = \frac{\sum_p \mathbf{A}_{p,k} \mathbf{F}_p}{\sum_p \mathbf{A}_{p,k}}. \quad (2)$$

Hence, as long as  $\mathcal{E}$  is invariant w.r.t. rigid transformations of  $\mathbf{P}$ , the pose  $\theta_k$  will be transformation equivariant, and the descriptor  $\beta_k$  will be transformation invariant. Note that this simplifies the design (and training) of the encoder  $\mathcal{E}$ , which only needs to be invariant, rather than equivariant [44, 41].

**Canonicalization.** Simply enforcing invariance and equivariance with the above framework is not enough to learn 3D representations that are object-centric, as we lack an (unsupervised) mechanism to bring information into a shared “object-centric” reference frame. Furthermore, the “right” canonical frame is nothing but a convention, thus we need a mechanism that allows the network to make a choice – a choice, however, that must then be consistent across all objects. For example, a learnt canonical frame where the cockpit of airplanes is consistently positioned along  $+z$  is *just as good* as a canonical frame where it is positioned along the  $+y$  axis. To address this, we propose to link the capsule descriptors to the capsule poses in canonical space, that is, we ask that objects with similar appearance to be located in similar Euclidean neighborhoods in canonical space. We achieve this by regressing canonical capsules poses (*i.e.* canonical keypoints)  $\bar{\theta} \in \mathbb{R}^{K \times 3}$  using the descriptors  $\beta \in \mathbb{R}^{K \times C}$  via a fully connected deep network  $\mathcal{K}$ :

$$\bar{\theta} = \mathcal{K}(\beta) \quad (3)$$

As fully connected layers are biased towards learning low-frequency representations [26], this regressor also acts as a regularizer that enforces semantic locality.

**Auto-encoding.** Finally, in the learnt *canonical* frame of reference, to train the capsule descriptors via auto-encoding, we reconstruct the point clouds with per-capsule decoders  $\mathcal{D}_k$ :

$$\tilde{\mathbf{P}} = \cup_k \{ \mathcal{D}_k(\bar{\mathbf{R}}\theta_k + \bar{\mathbf{t}}, \beta_k) \}, \quad (4)$$

where  $\cup$  denotes the union operator. The canonicalizing transformation  $\bar{\mathbf{T}} = (\bar{\mathbf{R}}, \bar{\mathbf{t}})$  can be readily computed by solving a shape-matching problem [40], thanks to the property that our capsule poses and regressed keypoints are in *one-to-one correspondence*:

$$\bar{\mathbf{R}}, \bar{\mathbf{t}} = \arg \min_{\mathbf{R}, \mathbf{t}} \frac{1}{K} \sum_k \|(\mathbf{R}\boldsymbol{\theta}_k + \mathbf{t}) - \bar{\boldsymbol{\theta}}_k\|_2^2. \quad (5)$$

While the reconstruction in (4) is in canonical frame, note it is trivial to transform the point cloud back to the original coordinate system after reconstruction, as the transformation  $\bar{\mathbf{T}}^{-1}$  is available.

### 3.1. Losses

As common in unsupervised methods, our framework relies on a number of losses that control the different characteristics we seek to obtain in our representation. Note how all these losses are unsupervised, and require no labels. We organize the losses according to the portion of the network they supervise: *decomposition*, *canonicalization*, and *reconstruction*.

**Decomposition.** While a transformation invariant encoder architecture should be sufficient to realize the desired equivariant/invariant properties, this does not prevent the encoder from producing trivial solutions/decompositions once trained. As capsule poses should be *transformation equivariant*, the poses of the two rotation augmentations  $\boldsymbol{\theta}_k^a$  and  $\boldsymbol{\theta}_k^b$  should only differ by the (known) relative transformation:

$$\mathcal{L}_{\text{equivariance}} = \frac{1}{K} \sum_k \|\boldsymbol{\theta}_k^a - (\mathbf{T}^a)(\mathbf{T}^b)^{-1}\boldsymbol{\theta}_k^b\|_2^2. \quad (6)$$

Conversely, capsule descriptors should be *transformation invariant*, and as the two input points clouds are of the *same* object, the corresponding capsule descriptors  $\boldsymbol{\beta}$  should be identical:

$$\mathcal{L}_{\text{invariance}} = \frac{1}{K} \sum_k \|\boldsymbol{\beta}_k^a - \boldsymbol{\beta}_k^b\|_2^2. \quad (7)$$

We further regularize the capsule decomposition to ensure each of the  $K$  heads roughly represent the same ‘‘amount’’ of the input point cloud, hence preventing degenerate (zero attention) capsules. This is achieved by penalizing the attention *variance*:

$$\mathcal{L}_{\text{equilibrium}} = \frac{1}{K} \sum_k \|a_k - \frac{1}{K} \sum_k a_k\|_2^2. \quad (8)$$

where  $a_k = \sum_p (\mathbf{A}_{p,k})$  denotes the total attention exerted by the  $k$ -th head on the point cloud.

Finally, to facilitate the training process, we ask for capsules to learn a localized representation of geometry. We

express the spatial extent of a capsule by computing first-order moments of the represented points with respect to the capsule pose  $\boldsymbol{\theta}_k$ :

$$\mathcal{L}_{\text{localization}} = \frac{1}{K} \sum_k \frac{1}{a_k} \sum_p \mathbf{A}_{p,k} \|\boldsymbol{\theta}_k - \mathbf{P}_p\|_2^2. \quad (9)$$

**Canonicalization.** To train our canonicalizer  $\mathcal{K}$ , we relate the predicted capsule poses to regressed canonical capsule poses via the *optimal* rigid transformation from (5):

$$\mathcal{L}_{\text{canonical}} = \frac{1}{K} \sum_k \|(\bar{\mathbf{R}}\boldsymbol{\theta}_k + \bar{\mathbf{t}}) - \bar{\boldsymbol{\theta}}_k\|_2^2. \quad (10)$$

Recall that  $\bar{\mathbf{R}}$  and  $\bar{\mathbf{T}}$  are obtained through a purely differentiable process. Thus, this loss is forcing the aggregated pose  $\boldsymbol{\theta}_k$  to agree with the one that goes through the regression path,  $\bar{\boldsymbol{\theta}}_k$ . Here,  $\bar{\boldsymbol{\theta}}_k$  is regressed solely from the set of capsule descriptors, hence similar shapes will result in similar canonical keypoints, and the coordinate system of  $\bar{\boldsymbol{\theta}}_k$  is one that employs Euclidean space to encode semantics.

**Reconstruction.** To learn canonical capsule descriptors in an unsupervised fashion, we rely on an auto-encoding task. We train the decoders  $\{\mathcal{D}_k\}$  by minimizing the *Chamfer Distance* (CD) between the (canonicalized) input point cloud and the reconstructed one, as in [55, 19]:

$$\mathcal{L}_{\text{recon}} = \text{CD}(\bar{\mathbf{R}}\mathbf{P} + \bar{\mathbf{t}}, \tilde{\mathbf{P}}). \quad (11)$$

### 3.2. Network Architectures

We briefly summarize our implementation details, including the network architecture; for further details, please refer to the **supplementary material**.

**Encoder –  $\mathcal{E}$ .** Our architecture is based on the one suggested in [42]: a pointnet-like architecture with residual connections and attentive context normalization. We utilize Batch Normalization instead of the Group Normalization, which trained faster in our experiments. We further extend their method to have *multiple* attention maps, where each attention map corresponds to a capsule.

**Decoder –  $\mathcal{D}$ .** The decoder from (4) operates on a per-capsule basis. Our decoder architecture is similar to AtlasNetV2 [13] (with trainable grids). The difference is that we translate the per-capsule decoded point cloud by the corresponding capsule pose.

**Regressor –  $\mathcal{K}$ .** We simply concatenate the descriptors and apply a series of fully connected layers with ReLU activation to regress the  $P$  capsule locations. At the output layer, we use a linear activation, and further subtract the mean of the outputs to make our regressed locations zero-centered in the canonical frame.

**Canonicalizing the descriptors.** As our descriptors are only *approximately* rotation invariant (via augmentation), we found it useful to re-extract the capsule descriptors  $\beta_k$  after canonicalization. Specifically, we compute  $\bar{\mathbf{F}}$  with the same encoder setup, but with  $\bar{\mathbf{P}}=\bar{\mathbf{R}}\mathbf{P}+\bar{\mathbf{T}}$  instead of  $\mathbf{P}$  and use it to compute  $\bar{\beta}_k$ ; we validate this empirically in Section 4.5.

## 4. Results

We first discuss the experimental setup in Section 4.1, and then validate our method on a *variety* of tasks: auto-encoding (Section 4.2), registration (Section 4.3), and unsupervised classification (Section 4.4). Importantly, while the task differs, our learning process remains the same: we learn capsules by reconstructing objects in a learnt canonical frame. We conclude with ablation studies in Section 4.5. For additional ablations and qualitative results, please refer to our supplementary material.

### 4.1. Experimental setup

To evaluate our method, we rely on the ShapeNet (Core) dataset [5]. We follow the category choices from AtlasNetV2 [13], using the airplane and chair classes for *single-category* experiments, while for *multi-category* experiments we use all 13 classes: airplane, bench, cabinet, car, chair, monitor, lamp, speaker, firearm, couch, table, cellphone, and watercraft. To make our results most compatible with those reported in the literature, we also use the same splits as in AtlasNetV2 [13]: 31747 shapes in the train, and 7943 shapes in the test set.<sup>1</sup> Unless noted otherwise, we randomly sample 1024 points from the object surface for each shape to create our 3D point clouds.

**De-canonicalizing the dataset.** As discussed in the introduction, ShapeNet (Core) contains substantial inductive bias in the form of consistent semantic alignment. To remove this bias, we create random  $\text{SE}(3)$  transformations, and apply them to each point cloud. We first generate uniformly sampled random rotations, and add uniformly sampled random translations within the range  $[-0.2, 0.2]$ , where the bounding volume of the shape ranges in  $[-1, +1]$ . Note the relatively limited translation range is chosen to give state-of-the-art methods a chance to compete with our solution. We then use the relative transformation between the point clouds extracted from this ground-truth transformation to evaluate our methods. We refer to this unaligned version of the ShapeNet Core dataset as the *unaligned* setup, and using the vanilla ShapeNet Core dataset as the *aligned* setup. For the *aligned* setup, as there is no need for equivariance adaptation, we simply train our method without the random transformations, and so  $\mathcal{L}_{\text{equivariance}}$  and  $\mathcal{L}_{\text{invariance}}$  are not used. This setup is to

<sup>1</sup>Note the numbers are slightly smaller than in [13], as they ignore the fact that the last batch is smaller than the rest (i.e. incomplete).

	Aligned			Unaligned		
	Airplane	Chair	Multi	Airplane	Chair	Multi
3D-PointCapsNet [58]	1.94	3.30	2.49	5.58	7.57	4.66
AtlasNetV2 [13]	1.28	2.36	2.14	2.80	3.98	3.08
Our method	<b>0.96</b>	<b>1.99</b>	<b>1.76</b>	<b>1.08</b>	<b>2.65</b>	<b>2.25</b>

Table 1. **Auto-encoding / quantitative** – Performance in terms of Chamfer distance – metric is multiplied by  $10^3$  as in [13].

simply demonstrate how Canonical Capsules would perform in the presence of a dataset bias.

We emphasize here that a proper generation of random rotation is important. While some existing works have generated them by uniformly sampling the degrees of freedom of an Euler-angle representation, this is known to be an incorrect way to sample random rotations [2], leading to biases in the generated dataset; see supplementary material.

**Implementation details.** For all our experiments we use the Adam optimizer [28] with an initial learning rate of 0.001 and decay rate of 0.1. We train for 325 epochs for the *aligned* setup to match the AtlasNetV2 [13] original setup. For the *unaligned* setting, as the problem is harder, we train for a longer number of 450 epochs. Unless stated otherwise, we use  $k=10$  capsules and capsule descriptors of dimension  $C=128$ . We train three models with our method: two that are single-category (*i.e.*, for airplane and chairs), and one that is multi-category (*i.e.*, all 13 classes). To set the weights for each loss term, we rely on the reconstruction performance (CD) in the training set. We set weights to be one for all terms except for  $\mathcal{L}_{\text{equivariance}}$  (5) and  $\mathcal{L}_{\text{equilibrium}}$  ( $10^{-3}$ ). In the aligned case, because  $\mathcal{L}_{\text{equivariance}}$  and  $\mathcal{L}_{\text{invariance}}$  are not needed (always zero), we reduce the weights for the other decomposition losses by  $10^3$ ;  $\mathcal{L}_{\text{localization}}$  to  $10^{-3}$  and  $\mathcal{L}_{\text{equilibrium}}$  to  $10^{-6}$ .

### 4.2. Auto-encoding – Figure 3 and Table 1

We evaluate the performance of our method for the task that was used to train the network – reconstruction / auto-encoding – against two baselines (trained in both single-category and multi-category variants):

- AtlasNetV2 [13], the state-of-the-art auto-encoder which utilizes a multi-head patch-based decoder;
- 3D-PointCapsNet [58], an auto-encoder for 3D point clouds that utilize a capsule architecture.

We do not compare against [41], as unfortunately no code is publicly available, and the paper is yet to be peer-reviewed.

**Quantitative analysis – Table 1.** We achieve state-of-the-art performance in both the *aligned* and *unaligned* settings. The wider margin in the *unaligned* setup indicates tackling the realistic, unbiased case degrades both AtlasNetV2 [13] and 3D-PointCapsNet [58] more than our method. We note

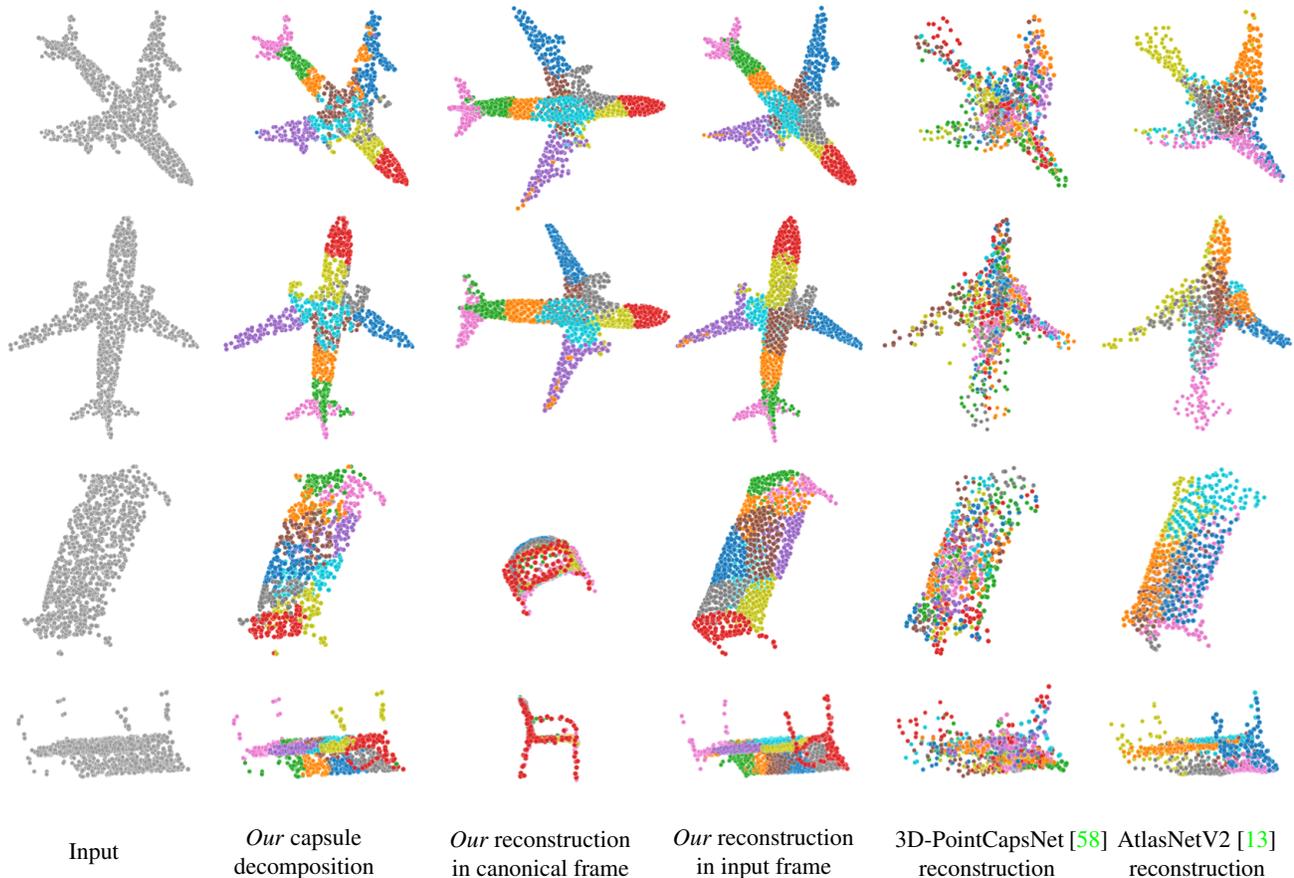


Figure 3. **Auto-encoding / qualitative** – Example decomposition and reconstruction results using Canonical Capsules on several point cloud instances from the test set. We color each Canonical Capsule with a unique colour, and similarly color “patches” from the reconstruction heads of 3D-PointCapsNet [58] and AtlasNetV2 [13]. Canonical Capsules provide semantically consistent decomposition that is aligned in canonical frame, leading to improved reconstruction quality.

that the results in this table differ slightly from what is reported in the original papers as we use 1024 points to speed-up our experiments; However, we show in Section 4.5 that the same trends hold regardless of the number of points, and match with what is reported in the original papers when 2500 points are used.

**Qualitative analysis – Figure 3.** We illustrate our decomposition-based reconstruction of 3D point clouds, as well as the reconstructions of 3D-PointCapsNet [58] and AtlasNetV2 [13]. As shown, even in the *unaligned* setup, our method is able to provide *semantically consistent* capsule decompositions – *e.g.* the wings of the airplane have consistent colours, and when aligned in the canonical frame, the different airplane instances are well-aligned. Compared to AtlasNetV2 [13] and 3D-PointCapsNet [58], the reconstruction quality is also visibly improved: we better preserve details along the engines of the airplane, or the thin structures of the bench; note also that the decompositions are semantically

consistent among the examples we show. Results are better appreciated in our supplementary material, where we visualize the performance as we continuously traverse  $SE(3)$ .

### 4.3. Registration – Table 2

We now evaluate the performance of our method on its capability to *register* 3D point clouds, and compare against three baselines:

- Deep Closest Points (DCP) [52], a deep learning-based point cloud registration method;
- DeepGMR–RRI [56], a state-of-the-art registration method that decomposes clouds into Gaussian mixtures given rotation-invariant features; see Rigorously Rotation-Invariant (RRI) features from [6];
- DeepGMR–XYZ [56], where *raw* XYZ coordinates are used as input instead of rotation-invariant features;
- Our method–RRI, a variant of our technique where we use RRI features [6] as the *sole* input of our architecture.

	Airplane	Chair	Multi
Deep Closest Points [52]	0.318	0.160	0.131
DeepGMR-XYZ [56]	0.079	0.082	0.077
Our method-XYZ	<b>0.024</b>	<b>0.027</b>	<b>0.070</b>
DeepGMR-RRI [56]	<b>0.0001</b>	<b>0.0001</b>	<b>0.0001</b>
Our method-RRI	0.0006	0.0009	0.0016

Table 2. **Registration** – Performance in terms of root mean-square error between registered and ground-truth points. Note that the methods we compare to are *pure* registration algorithms, whereas ours tackles the more generic problem of unsupervised representation learning.

	Aligned		Unaligned	
	SVM	K-Means	SVM	K-Means
AtlasNetV2	94.07	61.66	71.13	14.59
3D-PointCapsNet	93.81	65.87	64.85	17.12
Our method	<b>94.21</b>	<b>69.82</b>	<b>87.17</b>	<b>43.86</b>

Table 3. **Classification** – Top-1 accuracy (%)

For our method using RRI features, we follow the DeepGMR training protocol and train for 100 epochs, while for DCP and DeepGMR we use the authors’ official implementation.

**Quantitative analysis.** We report our results in Table 2. When RRI is used as input, our method is on par with DeepGMR, up to a level where registration is near perfect – alignment differences when errors are in the  $10^{-4}$  ballpark are indiscernible. Moreover, our method achieves the best performance when RRI is not used. We note that the performance of Deep Closest Points [52] is not as good as reported in the original paper, as we uniformly draw rotations from  $\mathbf{SO}(3)$ . When a sub-portion of  $\mathbf{SO}(3)$  is used, *e.g.* a quarter of what we are using, DCP performs relatively well (0.008 in the multi-class experiment). While curriculum learning could be used to enhance the performance of DCP, our technique does not need to rely on these more complex training techniques.

**Canonicalization and RRIs.** We further note that, while RRI delivers good registration performance, using RRI features cause the learnt canonicalization to fail –  $\mathcal{L}_{\text{canonical}}$  does not converge. This hints that RRI features may be throwing away too much information to achieve transformation invariance. Our method using raw XYZ coordinates as input, on the other hand, provides comparable registration performance, and is able to do significantly more than just registration (*i.e.* classification, reconstruction).

#### 4.4. Unsupervised classification – Table 3

Beyond reconstruction and registration, which are tasks that are directly relevant to the losses used for training,

	Full	$\neg\mathcal{L}_{\text{invar}}$	$\neg\mathcal{L}_{\text{canonical}}$	$\neg\mathcal{L}_{\text{equiv}}$	$\neg\mathcal{L}_{\text{localization}}$	$\neg\mathcal{L}_{\text{equilibrium}}$
CD	<b>1.08</b>	1.09	1.09	1.16	1.45	1.61

Table 4. **Effect of losses** – Reconstruction performance in terms of the Chamfer Distance CD (multiplied by  $10^3$ ) when loss terms are removed; *unaligned* setup and training on *airplanes* only.

we evaluate the usefulness of our method via a classification task that is not related in *any* way to the losses used for training. We compute the features from the auto-encoding methods compared in Section 4.2 – AtlasNetV2 [13], 3D-PointCapsNet [58], and our method (where we build features by combining pose with descriptors) – and use them to perform 13-way classification with two different techniques:

- We train a *supervised* linear Support Vector Machine (SVM) on the extracted features [3, Ch. 7];
- We perform *unsupervised* K-Means clustering [3, Ch. 9] and then label each cluster via bipartite matching with the actual labels through the Hungarian algorithm.

Note the former provides an upper bound for unsupervised classification, while better performance on the latter implies that the learnt features are able to separate the classes into clusters that are compact (in an Euclidean sense). We report classification performance in Table 3.

**Analysis of results – SVM.** Note how our method provides best results in all cases, and when the dataset is not unaligned the difference is significant. This shows that, while 3D-PointCapsNet and AtlasNetV2 are able to somewhat auto-encode point clouds in the *unaligned* setup, what they learn does not translate well to classification. However, the features learned with Canonical Capsules are more related to the semantics of the object, which helps classification.

**Analysis of results – K-Means.** The performance gap becomes wider when K-Means is used – even in the *aligned* case. This could mean that the features extracted by Canonical Capsules are better suited for other unsupervised tasks, having a feature space that is close to being Euclidean in terms of semantics. The difference is striking in the *unaligned* setup. We argue that these results emphasize the importance of the capsule framework – jointly learning the invariances and equivariances in the data – is cardinal to unsupervised learning [25, 23].

#### 4.5. Ablation studies

We further analyze different components of Canonical Capsules that affect the performance. To make the computational cost manageable, we perform all experiments in this section with the *airplane* category, and with the *unaligned* setup, unless otherwise noted.

	Airplane	All
One shot alignment	1.12	2.27
Our method	<b>1.08</b>	<b>2.25</b>

Table 5. **One-shot canonicalization** – We compare the reconstruction performance of our method against a naive one-shot alignment, where an arbitrary point cloud is selected as reference; *unaligned* setup and training on *airplanes* only.

	AtlasNetV2 [13]	Ours ( $\beta_k$ )	Ours ( $\bar{\beta}_k$ )
Aligned	1.28	<b>0.96</b>	0.99
Unaligned	2.80	2.12	<b>1.08</b>

Table 6. **Effectiveness of canonical descriptors** – Auto-encoding performance (Chamfer Distance) of our method with descriptors computed directly on the input cloud  $\beta_k$  or in canonical pose  $\bar{\beta}_k$ .

**Effect of losses – Table 4.** We first analyze the importance of each loss term – with the exception of  $\mathcal{L}_{\text{recon}}$  – which is necessary for training. Among our losses, the most important appears to be  $\mathcal{L}_{\text{equilibrium}}$  from (8), as distributing approximately same number of points to each capsule enables the model to fully utilize its capacity.

**One-shot canonicalization – Table 5.** A naive alternative to our learnt canonicalizer would be to use one point cloud as a reference point cloud to align to. Using the canonicalizer provides improved reconstruction performance over this naïve approach, removes the dependency on the choice of the reference point cloud, and allows our method to work effectively when dealing with multi-class canonicalization.

**Effectiveness of canonical descriptors – Table 6.** We evaluate the effectiveness of the descriptor enhancement strategy described in Section 3.2. We report the reconstruction performance with and without the enhancement. Recomputing the descriptor in canonical frame helps when dealing with the *unaligned* setup. Note that even without this enhancement, our method outperforms the state-of-the-art.

**Supervising the attention’s invariance.** Since  $\theta$  is inferred by weighted averaging of  $\mathbf{P}$  with the attention map  $\mathbf{A}$  in (2), we also considered directly adding a loss on  $\mathbf{A}$  that enforces invariance instead of a loss on  $\theta$ . This variant provides slightly degraded performance of  $CD=1.11$  compared to our baseline  $CD=1.08$ . We hypothesize that this is because  $\mathcal{L}_{\text{equivariance}}$  directly supervises the end-goal (capsule pose equivariance) whereas supervising  $\mathbf{A}$  is an indirect one.

**Encoder architecture.** We further note that using a permutation equivariant attention architecture (ACNe) [42] for the encoder is essential. When ACNe is replaced with the permutation invariant PointNet [34], the auto-encoding performance drops to  $CD=1.52$ , which is *significantly* lower

	1024 pts	2500 pts
3D-PointCapsNet [58]	2.49	1.49
AtlasNetV2 [13]	2.14	1.22
Our method	<b>1.76</b>	<b>0.97</b>

Table 7. **Number of points  $P$**  – Auto-encoding performance (Chamfer distance) as we vary the input point cloud cardinality; *aligned* setup for both training and testing.

compared to our method which give  $CD=1.08$ .

**Number of points  $P$  – Table 7.** To speed-up experiments we have mostly used  $P=1024$ , but in the table we show that our findings are consistent regardless of the number of points used. Note that the AtlasNetV2 [13] results are *very* similar to what is reported in the original paper. The slight differences exist due to random subsets that were used in AtlasNetV2 and not fully reproducible.<sup>2</sup>

## 5. Conclusions and future work

In this paper, we provide an unsupervised framework to train capsule decompositions for 3D point clouds. We rely on a Siamese training setup, circumventing the customary need to train on pre-aligned datasets. Despite being trained in an unsupervised fashion, our representation achieves state-of-the-art performance across auto-encoding/reconstruction, registration and classification tasks. These results are made possible by allowing the network to learn a canonical frame of reference. We interpret this result as giving our neural networks a mechanism to construct a “mental picture” of a given 3D object – so that downstream tasks are executed within an object-centric coordinate frame.

**Future work.** There are many ways in which our work can be extended. As many objects have natural symmetries [13], providing our canonicalizer a way to encode such a prior is likely to further improve the representation. It would be interesting to investigate whether, by providing some part annotations (via few-shot learning), the network could learn to favor decompositions with higher-level semantics (*e.g.*, where the wing of an airplane is a *single* capsule). Our decomposition has a single layer, and it would be interesting to investigate how to effectively engineer multi-level decompositions [46]; one way could be to over-decompose the input in a redundant fashion (with large  $K$ ), and use a downstream layers that “selects” the decomposition heads to be used [8]. We would also like to extend our results to more “in-the-wild” 3D computer vision and understand whether learning object-centric representations is possible when *incomplete* (*i.e.*, single view [33]) data is given in in-

<sup>2</sup>And to a minor bug in the evaluation code (*i.e.* non deterministic test set creation) that we have already communicated to the authors of [19].

put, when an entire scene with potentially *multiple* objects is given [39], or where our measurement of the 3D world is a single 2D image [43], or by exploiting the persistence of objects in video.

## Acknowledgements

This work was supported by the Natural Sciences and Engineering Research Council of Canada (NSERC) Discovery Grant, NSERC Collaborative Research and Development Grant, Google, Compute Canada, and Advanced Research Computing at the University of British Columbia.

## References

- [1] Parnian Afshar, Arash Mohammadi, and Konstantinos N Plataniotis. Brain Tumor Type Classification via Capsule Networks. In *International Conference on Image Processing*, 2018. 2
- [2] James Arvo. Fast Random Rotation Matrices. In *Graphics Gems III (IBM Version)*, pages 117–120. Elsevier, 1992. 5, 13
- [3] Christopher M Bishop. *Pattern Recognition and Machine Learning*. springer, 2006. 7
- [4] Rohan Chabra, Jan Eric Lenssen, Eddy Ilg, Tanner Schmidt, Julian Straub, Steven Lovegrove, and Richard Newcombe. Deep Local Shapes: Learning Local SDF Priors for Detailed 3D Reconstruction. In *European Conference on Computer Vision*, 2020. 2
- [5] Angel X Chang, Thomas Funkhouser, Leonidas Guibas, Pat Hanrahan, Qixing Huang, Zimo Li, Silvio Savarese, Manolis Savva, Shuran Song, Hao Su, et al. Shapenet: An Information-Rich 3D Model Repository. *arXiv Preprint*, 2015. 1, 5
- [6] Chao Chen, Guanbin Li, Ruijia Xu, Tianshui Chen, Meng Wang, and Liang Lin. Clusternet: Deep Hierarchical Cluster Network with Rigorously Rotation-Invariant Representation for Point Cloud Analysis. In *Conference on Computer Vision and Pattern Recognition*, 2019. 6
- [7] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A Simple Framework for Contrastive Learning of Visual Representations. *International Conference on Machine Learning*, 2020. 1
- [8] Zhiqin Chen, Andrea Tagliasacchi, and Hao Zhang. Bsp-net: Generating compact meshes via binary space partitioning. In *Conference on Computer Vision and Pattern Recognition*, 2020. 8
- [9] Zhiqin Chen and Hao Zhang. Learning Implicit Fields for Generative Shape Modeling. In *Conference on Computer Vision and Pattern Recognition*, 2019. 2
- [10] Boyang Deng, JP Lewis, Timothy Jeruzalski, Gerard Pons-Moll, Geoffrey Hinton, Mohammad Norouzi, and Andrea Tagliasacchi. NASA: Neural Articulated Shape Approximation. In *European Conference on Computer Vision*, 2020. 1
- [11] Zhantao Deng, Jan Bednárík, Mathieu Salzmann, and Pascal Fua. Better Patch Stitching for Parametric Surface Reconstruction. *arXiv Preprint*, 2020. 2
- [12] Deng, Boyang and Genova, Kyle and Yazdani, Soroosh and Bouaziz, Sofien and Hinton, Geoffrey and Tagliasacchi, Andrea. CvxNet: Learnable Convex Decomposition. In *Conference on Computer Vision and Pattern Recognition*, 2020. 2
- [13] Theo Deprelle, Thibault Groueix, Matthew Fisher, Vladimir Kim, Bryan Russell, and Mathieu Aubry. Learning Elementary Structures for 3D Shape Generation and Matching. In *Advances in Neural Information Processing Systems*, 2019. 2, 4, 5, 6, 7, 8, 12, 13
- [14] Haoqiang Fan, Hao Su, and Leonidas J Guibas. A Point Set Generation Network for 3D Object Reconstruction from a Single Image. In *Conference on Computer Vision and Pattern Recognition*, 2017. 2
- [15] Pedro Felzenszwalb, David McAllester, and Deva Ramanan. A Discriminatively Trained, Multiscale, Deformable Part Model. In *Conference on Computer Vision and Pattern Recognition*, 2008. 1
- [16] Clara Fernandez-Labrador, Ajad Chhatkuli, Danda Pani Paudel, Jose J Guerrero, Cédric Demonceaux, and Luc Van Gool. Unsupervised Learning of Category-Specific Symmetric 3D Keypoints from Point Sets. In *European Conference on Computer Vision*, 2020. 3
- [17] Kyle Genova, Forrester Cole, Avneesh Sud, Aaron Sarna, and Thomas Funkhouser. Deep Structured Implicit Functions. In *Conference on Computer Vision and Pattern Recognition*, 2020. 1, 2
- [18] Kyle Genova, Forrester Cole, Daniel Vlasic, Aaron Sarna, William T Freeman, and Thomas Funkhouser. Learning Shape Templates with Structured Implicit Functions. In *International Conference on Computer Vision*, 2019. 2
- [19] Thibault Groueix, Matthew Fisher, Vladimir G Kim, Bryan C Russell, and Mathieu Aubry. A Papier-Mâché Approach to Learning 3D Surface Generation. In *Conference on Computer Vision and Pattern Recognition*, 2018. 1, 2, 4, 8
- [20] Jiayuan Gu, Wei-Chiu Ma, Sivabalan Manivasagam, Wenyan Zeng, Zihao Wang, Yuwen Xiong, Hao Su, and Raquel Urtasun. Weakly-Supervised 3D Shape Completion in the Wild. In *European Conference on Computer Vision*, 2020. 3
- [21] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep Residual Learning for Image Recognition. In *Conference on Computer Vision and Pattern Recognition*, 2016. 1
- [22] Geoffrey E Hinton, Alex Krizhevsky, and Sida D Wang. Transforming Auto-Encoders. In *International Conference on Artificial Neural Networks*, 2011. 1, 2
- [23] Geoffrey E Hinton and Kevin J Lang. Shape Recognition and Illusory Conjunctions. In *International Joint Conference on Artificial Intelligence*, 1985. 7
- [24] Geoffrey E Hinton, Sara Sabour, and Nicholas Frosst. Matrix Capsules with EM Routing. In *International Conference on Learning Representations*, 2018. 2
- [25] Geoffrey F Hinton. A Parallel Computation that Assigns Canonical Object-based Frames of Reference. In *International Joint Conference on Artificial Intelligence*, 1981. 7
- [26] Arthur Jacot, Franck Gabriel, and Clément Hongler. Neural Tangent Kernel: Convergence and Generalization in Neural Networks. In *Advances in Neural Information Processing Systems*, 2018. 3

- [27] Abhishek Kar, Christian Häne, and Jitendra Malik. Learning a Multi-View Stereo Machine. In *Advances in Neural Information Processing Systems*, 2017. 2
- [28] D.P. Kingma and J. Ba. Adam: A Method for Stochastic Optimisation. In *International Conference on Learning Representations*, 2015. 5
- [29] Adam Kosiorek, Sara Sabour, Yee Whye Teh, and Geoffrey E Hinton. Stacked Capsule Autoencoders. In *Advances in Neural Information Processing Systems*, 2019. 2
- [30] Or Litany, Alex Bronstein, Michael Bronstein, and Ameesh Makadia. Deformable Shape Completion with Graph Convolutional Autoencoders. In *Conference on Computer Vision and Pattern Recognition*, 2018. 2
- [31] David G Lowe. Distinctive Image Features from Scale-Invariant Keypoints. *International Journal of Computer Vision*, 60:91–110, 2004. 1
- [32] Lars Mescheder, Michael Oechsle, Michael Niemeyer, Sebastian Nowozin, and Andreas Geiger. Occupancy Networks: Learning 3D Reconstruction in Function Space. In *Conference on Computer Vision and Pattern Recognition*, 2019. 1, 2
- [33] Jeong Joon Park, Peter Florence, Julian Straub, Richard Newcombe, and Steven Lovegrove. DeepSDF: Learning Continuous Signed Distance Functions for Shape Representation. In *Conference on Computer Vision and Pattern Recognition*, 2019. 2, 8
- [34] Charles R Qi, Hao Su, Kaichun Mo, and Leonidas J Guibas. Pointnet: Deep Learning on Point Sets for 3D Classification and Segmentation. In *Conference on Computer Vision and Pattern Recognition*, 2017. 8
- [35] Joseph Redmon, Santosh Divvala, Ross Girshick, and Ali Farhadi. You Only Look Once: Unified, Real-Time Object Detection. In *Conference on Computer Vision and Pattern Recognition*, 2016. 1
- [36] Davis Rempe, Tolga Birdal, Yongheng Zhao, Zan Gojcic, Srinath Sridhar, and Leonidas J. Guibas. CaSPR: Learning Canonical Spatiotemporal Point Cloud Representations. In *Advances in Neural Information Processing Systems*, 2020. 3
- [37] Sara Sabour, Nicholas Frosst, and Geoffrey E Hinton. Dynamic Routing Between Capsules. In *Advances in Neural Information Processing Systems*, 2017. 2
- [38] Karen Simonyan and Andrew Zisserman. Very Deep Convolutional Networks for Large-Scale Image Recognition. *International Conference on Learning Representations*, 2015. 1
- [39] Shuran Song, Fisher Yu, Andy Zeng, Angel X Chang, Manolis Savva, and Thomas Funkhouser. Semantic scene completion from a single depth image. *Proceedings of 30th IEEE Conference on Computer Vision and Pattern Recognition*, 2017. 9
- [40] Olga Sorkine-Hornung and Michael Rabinovich. Least-Squares Rigid Motion Using SVD. *Computing*, 2017. 4
- [41] Nitish Srivastava, Hanlin Goh, and Ruslan Salakhutdinov. Geometric Capsule Autoencoders for 3D Point Clouds. *arXiv Preprint*, 2019. 2, 3, 5
- [42] Weiwei Sun, Wei Jiang, Eduard Trulls, Andrea Tagliasacchi, and Kwang Moo Yi. ACNe: Attentive Context Normalization for Robust Permutation-Equivariant Learning. In *Conference on Computer Vision and Pattern Recognition*, 2020. 4, 8, 12
- [43] Supasorn Suwajanakorn, Noah Snively, Jonathan J Tompson, and Mohammad Norouzi. Discovery of Latent 3D Keypoints via End-to-End Geometric Reasoning. In *NIPS*, 2018. 3, 9
- [44] Nathaniel Thomas, Tess Smidt, Steven Kearnes, Lusann Yang, Li Li, Kai Kohlhoff, and Patrick Riley. Tensor Field Networks: Rotation-and Translation-Equivariant Neural Networks for 3D Point Clouds. *arXiv Preprint*, 2018. 2, 3
- [45] Shubham Tulsiani, Hao Su, Leonidas J Guibas, Alexei A Efros, and Jitendra Malik. Learning shape abstractions by assembling volumetric primitives. In *Conference on Computer Vision and Pattern Recognition*, 2017. 2
- [46] Oliver van Kaick, Kai Xu, Hao Zhang, Yanzhen Wang, Shuyang Sun, Ariel Shamir, and Daniel Cohen-Or. Co-hierarchical analysis of shape structures. *ACM SIGGRAPH*, 2013. 8
- [47] Dilin Wang and Qiang Liu. An Optimization View on Dynamic Routing Between Capsules. In *International Conference on Learning Representations*, 2018. 2
- [48] He Wang, Srinath Sridhar, Jingwei Huang, Julien Valentin, Shuran Song, and Leonidas J Guibas. Normalized Object Coordinate Space for Category-level 6d Object Pose and Size Estimation. In *Conference on Computer Vision and Pattern Recognition*, 2019. 3
- [49] Nanyang Wang, Yinda Zhang, Zhuwen Li, Yanwei Fu, Wei Liu, and Yu-Gang Jiang. Pixel2mesh: Generating 3d mesh models from single rgb images. In *European Conference on Computer Vision*, 2018. 2
- [50] Peng-Shuai Wang, Yang Liu, Yu-Xiao Guo, Chun-Yu Sun, and Xin Tong. O-CNN: Octree-Based Convolutional Neural Networks for 3d Shape Analysis. *ACM Transactions on Graphics*, 36(4):1–11, 2017. 2
- [51] Peng-Shuai Wang, Chun-Yu Sun, Yang Liu, and Xin Tong. Adaptive o-cnn: A patch-based deep representation of 3d shapes. *ACM Transactions on Graphics (TOG)*, 37(6):1–11, 2018. 2
- [52] Yue Wang and Justin M Solomon. Deep Closest Point: Learning Representations for Point Cloud Registration. In *International Conference on Computer Vision*, 2019. 3, 6, 7
- [53] Saining Xie, Jiatao Gu, Demi Guo, Charles R. Qi, Leonidas J. Guibas, and Or Litany. PointContrast: Unsupervised Pre-training for 3D Point Cloud Understanding. In *European Conference on Computer Vision*, 2020. 1
- [54] Xinchun Yan, Jimei Yang, Ersin Yumer, Yijie Guo, and Honglak Lee. Perspective Transformer Nets: Learning Single-View 3D Object Reconstruction without 3D Supervision. In *Advances in Neural Information Processing Systems*, 2016. 2
- [55] Yaoqing Yang, Chen Feng, Yiru Shen, and Dong Tian. FoldingNet: Point Cloud Auto-Encoder via Deep Grid Deformation. In *Conference on Computer Vision and Pattern Recognition*, 2018. 1, 2, 4
- [56] Wentao Yuan, Ben Eckart, Kihwan Kim, Varun Jampani, Dieter Fox, and Jan Kautz. DeepGMR: Learning Latent Gaussian Mixture Models for Registration. In *European Conference on Computer Vision*, 2020. 3, 6, 7
- [57] Wei Zhao, Haiyun Peng, Steffen Eger, Erik Cambria, and Min Yang. Towards Scalable and Reliable Capsule Networks for Challenging NLP Applications. In *Annual Meeting of the Association for Computational Linguistics*, 2019. 2

- [58] Yongheng Zhao, Tolga Birdal, Haowen Deng, and Federico Tombari. 3D Point Capsule Networks. In *Conference on Computer Vision and Pattern Recognition*, 2019. [2](#), [5](#), [6](#), [7](#), [8](#), [13](#)
- [59] Yongheng Zhao, Tolga Birdal, Jan Eric Lenssen, Emanuele Menegatti, Leonidas Guibas, and Federico Tombari. Quaternion Equivariant Capsule Networks for 3D Point Clouds. In *European Conference on Computer Vision*, 2020. [2](#)

# Canonical Capsules: Unsupervised Capsules in Canonical Pose

## Supplementary Material

### A. Architectural details

To facilitate reproduction of our results, we detail our architecture design for the capsule encoder  $\mathcal{E}$ , the decoder  $\mathcal{D}$  and the canonicalizer  $\mathcal{K}$ .

#### A.1. Capsule Encoder – $\mathcal{E}$

The capsule encoder  $\mathcal{E}$  takes in a point cloud  $\mathbf{P} \in \mathbb{R}^{P \times D}$  and outputs a  $K$ -head attention map  $\mathbf{A} \in \mathbb{R}^{P \times K}$  and the feature map  $\mathbf{F} \in \mathbb{R}^{P \times C}$ . Specifically, the  $\mathcal{E}$  is composed of 3 residual blocks. Similar to ACNe [42], each residual block consists of two hidden MLP layers with 128 neurons each. Each hidden MLP layer is followed by Attentive Context Normalization (ACN) [42], batch normalization [A1], and ReLU activation. To be able to attend to multiple capsules, we extend ACN layer into the multi-headed ACN.

**Multi-headed attentive context normalization.** The idea of multi-headed ACN starts from the observation that different capsules may require attending to different parts. Hence, for each MLP layer where we apply ACN, if we denote this layer as  $i$ , we first train a fully-connected layer that creates a  $K$ -headed attention map  $\mathbf{A}^i \in \mathbb{R}^{P \times K}$  given the  $\mathbf{F}^i \in \mathbb{R}^{P \times C}$  of this layer. This is similar to ACN, but instead of a single attention map, we now have  $K$ . The normalization process is similar to ACN afterwards – utilizing weighted moments of  $\mathbf{F}^i$  with  $\mathbf{A}^i$  – but results in  $K$  normalized outcomes instead of one. We then aggregate these  $K$  normalization results into one by summing.

Specifically, given the  $\mathbf{F}_p^i \in \mathbb{R}^{P \times C}$ , if we denote the weights and biases to be trained for the  $k^{th}$  attention head to be  $\mathbf{W}_k^i \in \mathbb{R}^{C \times 1}$  and  $b_k^i \in \mathbb{R}^1$  we write:

$$\mathbf{A}_{p,k}^i = \frac{\exp(\mathbf{F}_p^i \mathbf{W}_k^i + b_k^i)}{\sum_k \exp(\mathbf{F}_p^i \mathbf{W}_k^i + b_k^i)}. \quad (12)$$

We then compute the moments that are used to normalize in ACN, but now for each attention head:

$$\boldsymbol{\mu}_k = \sum_p \frac{\mathbf{A}_{p,k}^i \mathbf{F}_p^i}{\sum_p \mathbf{A}_{p,k}^i}, \quad \boldsymbol{\sigma}_k = \sum_p \frac{\mathbf{A}_{p,k}^i (\mathbf{F}_p^i - \boldsymbol{\mu}_k)^2}{\sum_p \mathbf{A}_{p,k}^i}, \quad (13)$$

which we then use to normalize and aggregate (sum) to get our final normalized feature map:

$$\mathbf{F}_p^i = \sum_k \mathbf{A}_{p,k}^i \frac{(\mathbf{F}_p^i - \boldsymbol{\mu}_k)}{\sqrt{\boldsymbol{\sigma}_k + \epsilon}}, \quad (14)$$

where  $\epsilon = 0.001$  is a very small value to avoid numerical instability.

#### A.2. Capsule Decoder – $\mathcal{D}$

The decoder  $\mathcal{D}$  is composed of  $K$  per-capsule decoders  $\mathcal{D}_k$ . Each per-capsule decoder  $\mathcal{D}_k$  maps the  $k^{th}$  capsule in canonical frame  $(\mathbf{R}\boldsymbol{\theta}_k + \bar{\mathbf{t}}, \boldsymbol{\beta}_k)$  to a group of points  $\tilde{\mathbf{P}}_k \in \mathbb{R}^{M \times D}$  which should correspond to a part of the entire object. We then obtain the auto-encoded (reconstructed) point clouds  $\tilde{\mathbf{P}}$  by collecting the outputs of  $K$  per-capsule decoders and taking their union as in (4).

Specifically, each  $\mathcal{D}_k$  consists of 3 hidden MLP layers of (1280, 640, 320) neurons, each followed by batch normalization and a ReLU activation. At the output layer, we use a fully connected layer, followed by the Tanh activation to regress the coordinates for each point. Similarly to AtlasNetV2 [13],  $\mathcal{D}_k$  additionally receives a set of trainable grids of size  $(M \times 10)$  and deforms them into  $\tilde{\mathbf{P}}_k$ , based on  $\boldsymbol{\beta}_k$ . Finally, as each output point of  $\mathcal{D}_k$  should be relative with respect to the capsule pose, we translate the generated points by the corresponding capsule’s canonicalized pose.

#### A.3. Canonicalizer – $\mathcal{K}$

The canonicalizer  $\mathcal{K}$  learns to regress the canonical pose for each capsule  $\bar{\boldsymbol{\theta}} \in \mathbb{R}^{K \times D}$  from their descriptors  $\boldsymbol{\beta}$ . In order to do so, we concatenate  $\boldsymbol{\beta}$  into a single, global descriptor, which we then feed into a fully connected layer with  $128 \times K$  neurons and a ReLU activation, followed by an additional fully connected layer with  $D \times K$  neurons, creating  $K$   $D$ -dimensional outputs – the poses. We do not apply any activation on the second layer. To make  $\bar{\boldsymbol{\theta}}$  zero-centered in the canonical frame, we further subtract the mean of the outputs.

### B. Qualitative results – auto-encoding with aligned data

For completeness, we further show qualitative results for auto-encoding on an aligned dataset, the common setup in prior work. As shown in Figure 4, our method provides best reconstruction performance even in this case; for quantitative results, see Table 1. Interestingly, while our decoder architecture is very similar to AtlasNetV2 [13], our reconstructions are of higher quality; our methods provides finer details at the propellers on the airplane, at the handle of the firearm, and at the back of the chair. This further supports the effectiveness of our capsule encoding.

### C. Additional ablation study

**Effect of number of capsules – Table 8.** To verify how the number of capsules affect performance, we test with

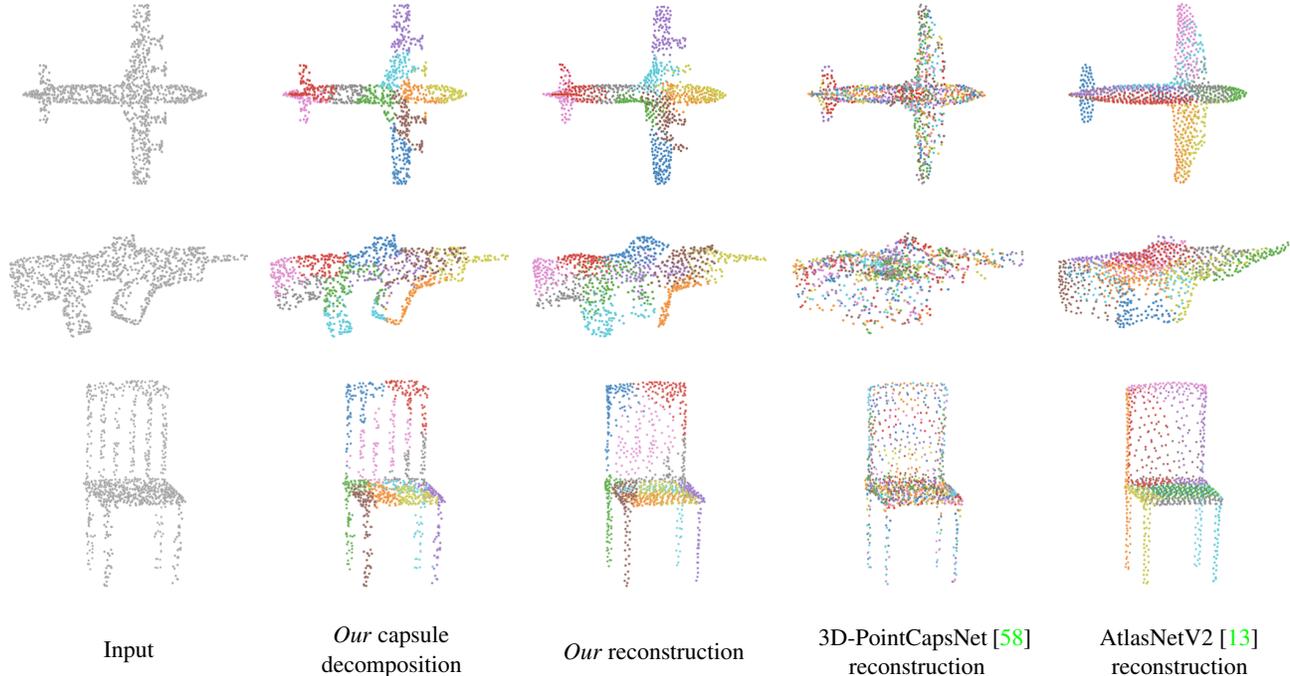


Figure 4. **Auto-encoding / qualitative** – Example decomposition results using Canonical Capsules on the test set, with the *aligned* setup. We color each decomposition (capsule) with a unique color – for 3D-PointCapsNet [58] and AtlasNetV2 [13], these correspond to “patches” in the reconstruction network. Our method provides best results.

AtlasNetV2 [13]	5 capsules	10 capsules	20 capsules
2.80	1.25	<b>1.08</b>	1.13

Table 8. **Ablation study on the number of capsules** – We show the reconstruction performance (Chamfer distance – metric multiplied by  $10^3$ ) with varying number of capsules. While they all perform better than competitors, 10 capsules give best performance.

varying number of capsules; 5, 10, and 20. As the number of capsules is the *only* factor that we wish the vary, we keep everything else identical, including the representation power by reducing the dimension of the descriptor as more capsules are used. For example, with 10 capsules we use a 128-dimensional descriptor, with 20 we use 64, and with 5, we use 256. Our experimental results show that representation with 10 capsules achieves the best performance. Note that our method, even with the sub-optimal number of capsules, still outperforms compared methods by a large margin.

**Random sampling of rotations** – Figure 5. Lastly, we revisit how rotations are randomly sampled to generate the augmentations used by Siamese training. Uniform sampling of Euler angles (i.e., yaw, pitch, roll) leads to a non-uniform coverage of the  $SO(3)$  manifold as shown in Figure 5 (a). Due to this non-uniformity, the reconstruction quality is biased with respect to test-time rotations; see Figure 5 (b).

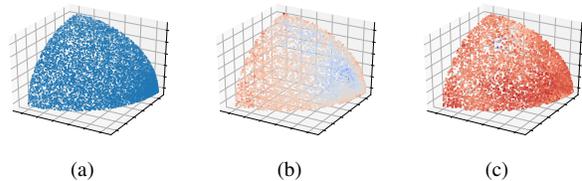


Figure 5. **Random sampling of rotations** – (a) Sampling Euler angles uniformly results in a non-uniform coverage of  $SO(3)$  (we sample one-eighth of a sphere – 0 to 90 degrees for each Euler angle – for easy visualization of the 3D space on paper). (b) This results in auto-encoding error to be biased w.r.t rotations (we use a cold-warm colormap to visualize the Chamfer Distance error). (c) By properly sampling rotations [2], this bias can be alleviated.

Instead, by properly sampling [2] the reconstruction performance is much more *uniform* across the manifold; see see Figure 5 (c) In our experiments, this leads to a significant difference in auto-encoding performance;  $CD=1.08$  with proper uniform sampling vs  $CD=1.17$  with the Euclidean random sampling.

## Additional References

- [A1] Sergey Ioffe and Christian Szegedy. Batch Normalization: Accelerating Deep Network Training by Reducing Internal Covariate Shift. In *International Conference on Machine Learning*, 2015.