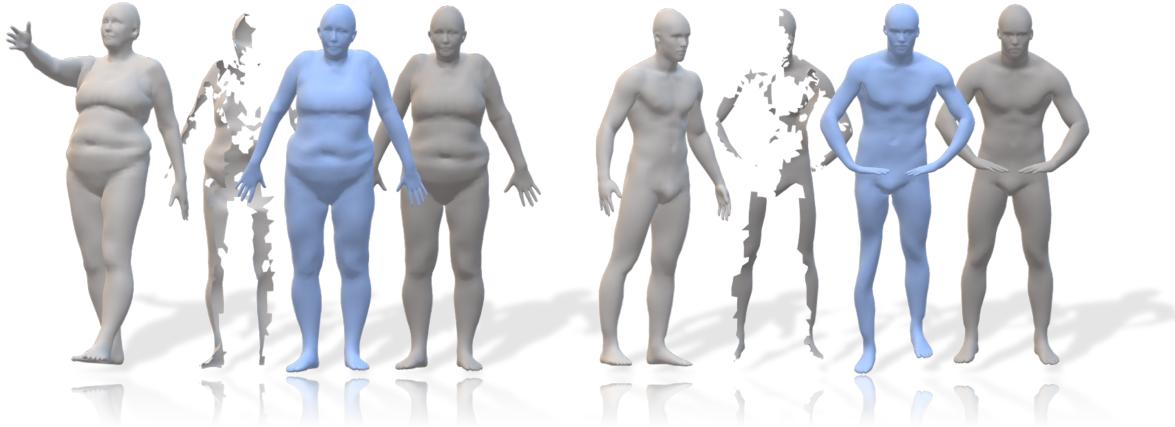


000
001
002
003
004
005
006
007
008
009
010
011
012
013
014
015
016
017
018
019
020
021
022
023
024
025
026
027
028
029
030
031
032
033
034
035
036
037
038
039
040
041
042
043
044
045
046
047
048
049
050
051
052
053

054
055
056
057
058
059
060
061
062
063
064
065
066
067
068
069
070
071
072
073
074
075
076
077
078
079
080
081
082
083
084
085
086
087
088
089
090
091
092
093
094
095
096
097
098
099
100
101
102
103
104
105
106
107

The Whole Is Greater Than the Sum of Its Nonrigid Parts



Left to right: input reference shape, input part, output completion, and the ground truth full model.

Anonymous CVPR submission

Paper ID 5649

Abstract

According to Aristotle, a philosopher in Ancient Greek, “the whole is greater than the sum of its parts”. This observation was adopted to explain human perception by the Gestalt psychology school of thought in the twentieth century. Here, we claim that observing part of an object which was previously acquired as a whole, one could deal with both partial matching and pose reconstruction in a holistic manner. More specifically, given the geometry of a full, articulated object in a given pose, as well as a partial scan of the same object in a different pose, we address the problem of matching the part to the whole while simultaneously reconstructing the new pose from its partial observation. Our approach is data-driven, and takes the form of a Siamese autoencoder without the requirement of a consistent vertex labeling at inference time; as such, it can be used on unorganized point clouds as well as on triangle meshes. We demonstrate the practical effectiveness of our model in the applications of single-view deformable shape completion and dense shape correspondence, both on synthetic and real-world geometric data, where we outperform prior work on these tasks by a large margin.

1. Introduction

Aristotle, a philosopher in Ancient Greek announced that “*the whole is greater than the sum of its parts*”. This axiomatic observation was narrowed down to human perception of planar shapes by Gestalt psychology school of thought in the twentieth century.

From a practical perspective, advances in robotics and augmented and virtual reality technologies often rely upon the ability to render a scene from novel views, manipulate its content, and add physical constraints. This often requires the completion of geometric structures from partial data. To that end, we would like to have a systematic and universal way to complete a partial shape into its full counterpart, a problem frequently referred to as *shape completion* in the geometry processing community.

Shape completion is an ill-posed problem by definition, as one must address the question of what can be considered a legitimate completion. An attempt to answer this question can be addressed in a statistical manner. One can consider many instances of partial shapes and their corresponding completions, defining a statistical relation between a part and the geometric structure from which it was cropped, and by which the whole structure can be recovered in a statistically optimal manner.

In this paper, we introduce a *deterministic* shape comple-

tion framework. Given a partial observation, the method returns a reliable reconstruction of a full, realistic object. At a first glance, at least in a rigid setting, this task may seem impossible: How can we guarantee that such a reconstruction reliably describes the part that was hidden from the viewing direction? However, one soon realizes that the problem can be formulated differently in the *non-rigid* case. Since two non-rigidly related shapes could share the same intrinsic geometry [22, 14, 10], each shape holds information about the other. We therefore pose the alternative question: Given a full object Q and a partial view P in a different pose, can we reconstruct a full version of P by borrowing geometric information from Q ? In this paper we address precisely this question. Our goal can thus be seen as a combination of two complementary tasks: (1) partial shape matching, and (2) non-rigid shape completion. While the two tasks are often addressed separately, we claim that considering their coupling provides powerful means to deal with both.

Contribution. In this paper we propose a novel formulation unifying deformable shape completion and partial shape matching. Specifically, given a full and a partial shape related by a non-rigid transformation, our objective is to deform the full shape to best fit the partial data. To compute this deformation, we train an encoder-decoder network. Once a completion for the part is predicted, the part-to-full correspondence can be trivially recovered using nearest-neighbor search. Our main contributions can be summarized as follows:

1. We introduce a deep Siamese architecture to tackle non-rigid alignment between a shape and its partial scan;
2. To the best of our knowledge, the proposed method is the first that addresses shape completion and partial correspondence in one framework under *extreme* partiality;
3. The proposed method is *efficient*, taking less than a second to provide both outputs without requiring any time-consuming post-processing steps.

2. Related work

Our problem setting is closely related to multiple research directions in the shape analysis and geometric deep learning communities. In an early attempt to use one pose in order to geometrically reconstruct another, Devir *et al.* [20] considered mapping a clean shape in a given pose onto the same noisy shape in a different pose. Elad and Kimmel were the first to treat shapes as metric spaces [21, 22]. In fact, it was the first approach of matching shapes by their spectra, specifically, second order moments of embedding

the intrinsic metric into a Euclidean one via classical scaling. It involved comparing the spectra of the click's laplacian. Bronstein *et al.* [9, 10, 13, 8, 11, 14] dealt with partial matching of articulated objects in various scenarios, including pruning of the intrinsic structure while accounting for cuts.

Shape completion. Recovering a complete shape from partial measurements is a longstanding research problem that comes in many flavors. In the context of deformable shapes, early efforts focused on completion based on geometric priors [36] or reoccurring patterns [13, 38, 60, 40]. These methods are not suited for severe partiality. For such cases model-based techniques are quite popular, *i.e.*, category-specific parametric morphable models that can be fitted to the partial data [4, 23, 44, 1]. Model-based shape completion was demonstrated for keypoint input [2], and was recently proven to be quite useful for recovering 3D body shapes from 2D images [65, 64, 28, 72]. Parametric morphable models [4], coupled with axiomatic image formation models were used to train a network to reconstruct face geometry from images [56, 55, 61]. Still, much less attention has been given to the task of fitting a model to a partial 3D point cloud. Recently, Jiang *et al.* [34] tackled this problem using a skeleton-aware architecture. However, their approach works well when full coverage of the underlying shape is given.

Deep learning of surfaces. Following the huge success of convolutional neural networks in images, in recent years, the geometry processing community has been rapidly adopting and designing computational modules suited for such data. The main challenge is that unlike images, geometric structures like surfaces come in many types of representations, and each requires a unique handling. Early efforts focused on a simple extension from a single image to multi-view representations [62, 68]. Another natural extension are 3D CNN on volumetric grids [69]. A host of techniques for mesh processing were developed as part of a research branch termed *geometric deep learning* [15]. These include graph-based methods [66, 67, 31], intrinsic patch extraction [47, 7, 48], and spectral techniques [41, 29]. Point cloud networks [52, 53] have recently gained much attention. Offering a light-weight computation restricted to sparse points with a sound geometric explanation [35], these networks have shown to provide a good compromise between complexity and accuracy, and are dominating the field of 3D object detection [51, 70], semantic segmentation [24, 71] and even temporal point cloud processing [17, 43]. For generative methods, recent implicit and parametric methods have demonstrated promising results [27, 49].

Inspired by the recent success of [26] in encoding non-rigid shape deformations using a point cloud network, here,

216 we also choose to use a point cloud representation. Import-
 217 antly, while the approach presented in [26] predicts align-
 218 ment of two shapes, it is not designed to handle severe par-
 219 tiality, and assumes a fixed template for the source shape.
 220 Instead, we show how to align arbitrary input shapes and
 221 focus on such a partiality.
 222

223 **Partial shape matching.** Dense non-rigid shape correspon-
 224 dence [37, 16, 41, 30, 58, 12, 18] is a key challenge in 3D
 225 computer vision and graphics, and has been widely explored
 226 in the last few years. A particularly challenging setting
 227 arises whenever one of the two shapes has missing geo-
 228 metry. This setting has been tackled with moderate success
 229 in a few recent works [57, 42, 54], however it largely re-
 230 mains an open problem whenever the partial shape exhibits
 231 severe artifacts or large, irregular missing parts. In this pa-
 232 per we tackle precisely this setting, demonstrating unprece-
 233 dented performance on a variety of real-world and synthetic
 234 datasets.

235 3. Method

236 3.1. Overview

237 We represent shapes as point clouds $S = \{s_i\}_{i=1}^{n_s}$ em-
 238 bedded in \mathbb{R}^3 . Depending on the setting, each point can
 239 carry additional semantic or geometric information encoded
 240 as feature vectors.

241 Given a full shape $Q = \{q_i\}_{i=1}^{n_q}$ and its partial view in a
 242 different pose $P = \{p_i\}_{i=1}^{n_p}$, our goal is to find a nonlinear
 243 function $F : \mathbb{R}^3 \rightarrow \mathbb{R}^3$ aligning Q to P^1 . If $R = \{r_i\}_{i=1}^{n_r}$
 244 is the (unknown) full shape such that $P \subset R$, ideally we
 245 would like to ensure that $F(Q) = R$. Thus, the deformed
 246 shape $F(Q)$ acts as a proxy to solve for the correspondence
 247 between the part P and the whole Q . By calculating for
 248 every vertex in P its nearest neighbor in $R \approx F(Q)$, we
 249 trivially obtain the mapping from P to Q as well.

250 Clearly, the deformation function F depends on the in-
 251 put pair of shapes (P, Q) . We model such dependency by
 252 considering a parametric function $F_\theta : \mathbb{R}^3 \rightarrow \mathbb{R}^3$, where θ
 253 is a latent encoding of the input pair (P, Q) . We implement
 254 this idea via an encoder-decoder neural network, and learn
 255 the space of parametric deformations from example pairs of
 256 partial and complete shapes.

257 Our network is composed of an encoder E and a gen-
 258 erator F_θ . The encoder takes as input the pair (P, Q) and em-
 259 beds it into a latent code θ . To map points from Q to their
 260 new location, we feed them to the generator along with the
 261 latent code. In what follows we first describe each module,
 262 and then give details on the training procedure and the loss
 263 function. We refer to Figure 1 for a schematic illustration
 264 of our learning model.

265 ¹In our setting, we assume that the pose can be inferred from the partial
 266 shape (e.g., an entirely missing limb would make the prediction ambigu-
 267 ous), hence the deformation function F is well defined.

270 3.2. Encoder

271 We propose the adoption of a Siamese pair of encoders,
 272 each producing a global shape descriptor (respectively θ_{part}
 273 and θ_{whole}). The two codes are then concatenated so as to
 274 encode the information of the specific pair of shapes, $\theta =$
 275 $[\theta_{part}, \theta_{whole}]$.

276 Our choice for the internal architecture in the Siamese
 277 pair is based on a preliminary ideal requirement: the en-
 278 coder should *injectively* map each pair of shapes into a lat-
 279 ent code. In other words, we require that each shape should
 280 be accurately reconstructed from its latent code. This re-
 281 quirement guarantees that for a fixed full shape Q and two
 282 different parts P_1 and P_2 , there would be a way to deform
 283 the full shape differently by using $\theta(P_1, Q)$ or $\theta(P_2, Q)$,
 284 based on the input part.

285 Motivated by existing methods [27, 26], we use a
 286 PointNet-3DCODED architecture for our encoder. Specifi-
 287 cally, each encoder within the Siamese pair applies a mul-
 288 tilayer perceptron to each 3D point of the input shape, with
 289 hidden dimensions (64, 128, 1024), followed by a max-
 290 pool operation over the input points, leading to a 1024-
 291 dimensional vector. Finally, we apply a linear layer of size
 292 1024 and a ReLu activation function. Hence, each shape in
 293 the input pair is represented by a latent code $\theta_{whole}, \theta_{part}$
 294 of size 1024 respectively, supplying a representation of the
 295 pair by a latent code θ of size 2048.

296 In practice, we find it helpful to also use the normal vec-
 297 tor coordinates as additional input features to each of the
 298 Siamese network encoders, making each input point a 6D
 299 point. The rationale behind this choice is that, by doing
 300 so, we are able to disambiguate contact points of the
 301 surface and thus to prevent contradicting requirements on the
 302 deformation function.

303 3.3. Generator

304 Given the code θ representing the partial and full shapes,
 305 the generator has to predict the deformation function F_θ
 306 to be applied to the full shape Q . We realize F as a
 307 multi-layer perceptron (MLP) to approximate the func-
 308 tional relation between an input point q_i on the full shape
 309 Q , and the corresponding output point r_i on the ground
 310 truth completed shape. The multi-layer perceptron op-
 311 erates pointwise on each tuple (q_i, θ) , where the shape
 312 context θ is kept fixed for each input pair. The result
 313 is the destination location $F_\theta(q_i) \in \mathbb{R}^3$, for each input
 314 point of the full shape Q . This generator architec-
 315 ture allows, in principle, to calculate the output reconstruc-
 316 tion in a flexible resolution, by providing the generator a
 317 full shape with some desired output resolution. In detail,
 318 the generator consists of 9 layers of hidden dimensions
 319 (2054, 1024, 512, 256, 128, 128, 128, 128, 3), followed by
 320 hyperbolic tangent activation function.

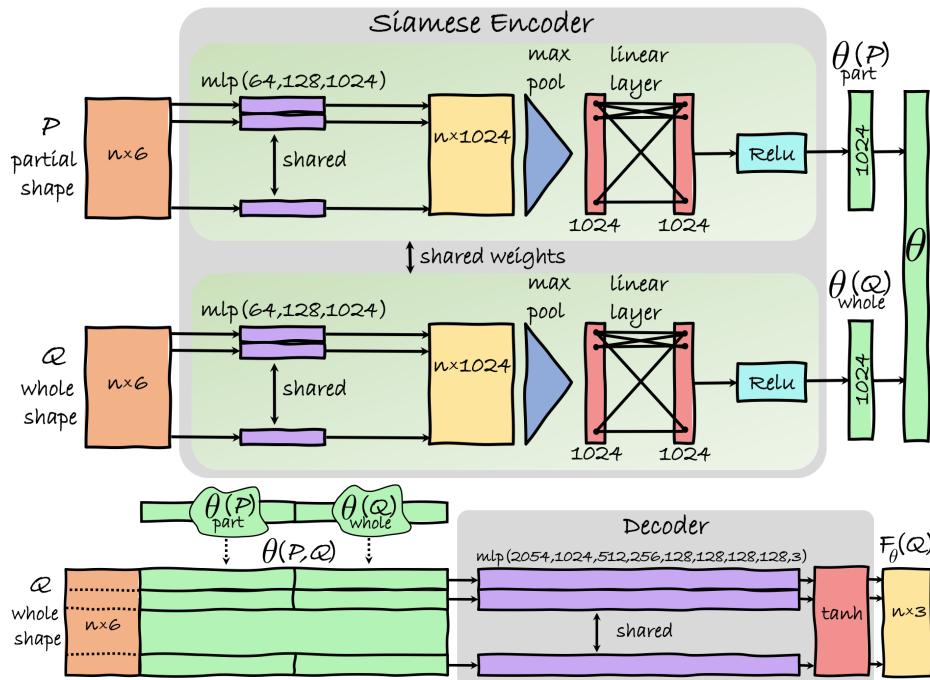


Figure 1. Siamese encoder architecture at the top, and the decoder (generator) architecture at the bottom. A shape is provided to the encoder as a list of 6D points, representing the spatial and unit normal coordinates. The latent codes of the input shapes $\theta_{part}(P)$ and $\theta_{whole}(Q)$ are concatenated to form a latent code θ representing the input pair. Based on this input pair, the decoder deforms the full shape by operating on each of its points with the same function. The result is the deformed full shape $F_\theta(Q)$.

3.4. Training Procedure

In our experiments we used training examples that were prepared from datasets of full human shapes, see details in the data preparation section 4.1. Each of these datasets contain 3D models of different subjects in various poses. Each of our training examples is composed of a triplet (P, Q, R) : a partial shape P , a full shape in a different pose Q and a ground truth completion R . The shape Q and R were sampled from the same subject in two different poses. In each training example, the partial shape P is extracted from R by rendering its depth map in a random viewpoint angle, with the azimuth changing between 0° and 360° and the elevation angle kept at 0° . These synthetic projections aim to approximate partiality realizations commonly occurring in depth sensors. By their nature, these partial shapes still hold most of the information needed to determine the degrees of freedom w.r.t the pose, making the pose reconstruction a well-posed problem. The training examples $(P_n, Q_n, R_n)_{n=1}^N$ were provided in batches to the Siamese Network, where N is the size of the train set. Each input pair is fed to the encoder to receive the latent code $\theta(P_n, Q_n)$ and the reconstruction $F_{\theta(P_n, Q_n)}(Q_n)$. This reconstruction is subsequently compared against the ground-truth reconstruction R_n using the loss defined in the next subsection 3.5.

3.5. Loss function

Essentially, the loss definition should reflect the visual plausibility of the reconstructed shape. Measuring such a quality analytically is a difficult problem. In this paper we adopt a naive measurement based on the Euclidean proximity between the ground-truth and the reconstruction. Formally, we define the loss as,

$$\mathcal{L}(P, Q, R) = \sum_{i=0}^{n_q} \|F_{\theta(P, Q)}(q_i) - r_i\|^2, \quad (1)$$

where $r_i = \pi^*(q_i) \in R$ is the ground-truth matched point of $q_i \in Q$, and $\pi^* : Q \rightarrow R$ is the ground-truth mapping between the full shape Q and the ground-truth reconstruction R . To promote the preservation of shape details we measure the Euclidean distance between the coordinates in \mathbb{R}^3 as well as the Euclidean distance between the normal vectors evaluated at each point, interpreting q_i and r_i in the expression above as points in \mathbb{R}^6 . Specifically, $q_i = (\vec{x}_{qi}, \alpha \vec{n}_{qi})$, $r_i = (\vec{x}_{ri}, \alpha \vec{n}_{ri}) \in \mathbb{R}^6$ are given by the concatenation of the coordinates vectors and the unit normal vector of each point in Q and R , respectively. The constant $\alpha > 0$ scales the normal vectors w.r.t. the Euclidean coordinates. The notation \vec{n}_{qi} denotes the normal of the surface Q at point q_i , while \vec{x}_{qi} denotes the coordinates of that point in \mathbb{R}^3 . Similarly, \vec{n}_{ri} denotes the normal of the surface R at

432 point r_i , while \vec{x}_{ri} denotes the coordinates of that point in
433 \mathbb{R}^3 .

434 The above loss may be refined for better preserving visual
435 appearance. Among the possible reconstructions that
436 admit a low value for \mathcal{L} one could find two shapes where
437 one achieves a lower \mathcal{L} value while the other looks more
438 realistic. The reason is that even though the parts were chosen
439 such that most degrees of freedom for the possible
440 completion w.r.t. pose are determined, still, the exact details
441 such as the accurate position and orientation of missing
442 regions are under-determined. Therefore, there is a narrow
443 margin of acceptable completions in the adjacency of the
444 provided ground-truth shape, implying that the Euclidean
445 proximity is measured w.r.t. to an arbitrary anchor. While
446 for small values of \mathcal{L} the correlation between the loss and the
447 reconstruction quality is indefinite, for high values there
448 is a strong correlation between the suggested loss and realism.
449 Empirically, we observed that this loss’s convergence
450 produces high quality reconstructions.

451 3.6. Implementation considerations

452 The network was trained with each batch containing 10
453 triplet examples (P, Q, R) , using the PyTorch [50] ADAM
454 optimizer with a learning rate of 0.001 and a momentum
455 of 0.9. We used a scale factor of $\alpha = 0.1$ for the normal
456 vector. The network was trained for 50 epochs, each con-
457 taining 10,000 random triplet examples. The input shapes
458 were translated such that their center of mass lies at the ori-
459 gin and the Iterative Closest Point algorithm [3] was further
460 applied on the network output to perfectly align the axes
461 w.r.t. the partial shape. Finally, to calculate the partial cor-
462 respondence for each point in the partial shape, we retrieve
463 its nearest-neighbor in the aligned reconstruction.

464 4. Experiments

465 The proposed method simultaneously tackles two impor-
466 tant tasks in nonrigid shape processing and analysis, shape
467 completion and partial shape matching. We emphasize that
468 the suggested framework gracefully handles severe partial-
469 ity. Prior efforts either addressed one of these tasks or at-
470 tempted to address both only at mild partiality conditions.
471 To thoroughly evaluate our performance, in this section we
472 test our method on each of these tasks and compare with
473 prior art. After a description of the datasets utilized, we
474 present results of shape completion from a single view , and
475 non-rigid partial correspondence. Finally, we show per-
476 formance on real scanned data.

477 4.1. Datasets

478 We utilize two datasets of human shapes for train-
479 ing and evaluation, FAUST [5] and AMASS [45]. Both
480 datasets were generated by fitting SMPL parametric body
481 model [44] to raw scans. The second took the approach

482 one step further by fitting those parameters to motion cap-
483 ture data. These datasets are quite different in size and
484 variability. FAUST is a relatively small set of 10 subjects
485 posing at 10 poses each. We follow previous methods,
486 and test our method in partial shape matching and shape
487 completion tasks using 10 projected views of these models.
488 AMASS is currently the largest and most diverse dataset of
489 human motion designed specifically for deep learning ap-
490 plications. It was generated by curating 15 archived datasets
491 of marker-based optical motion capture data and unifying
492 them into a shared statistical model SMPL+H [59]. As
493 such, it can provide a much richer resource for evaluating
494 the generalization ability of shape alignment and matching
495 techniques. To this end, we used AMASS to create a large
496 set of single-view projections. Specifically, we sampled ev-
497 ery 100th frame from all provided sequences and rendered
498 single-view projections from 10 equally spaced azimuth an-
499 gles (elevation was kept fixed) using pyRender [33]. Keeping
500 the original data splits our dataset consists a total of
501 110K, 10K, and 1K full shapes for train, validation and test,
502 respectively; and 10 times that size for the partial shapes.
503 Note that at train time we randomly mix and match full
504 shapes and their parts which drastically increases the effec-
505 tive set size.

511 4.2. Methods in comparison

512 A recent exploration of the problem of deformable shape
513 completion for arbitrary partiality can be found in Litany *et*
514 *al.* [39]. They suggest to find a completion via optimiza-
515 tion in a learned shape space. Note that in that paper, the
516 task was defined as a completion from a partial view with-
517 out explicit access to a full model. Moreover, their solu-
518 tion requires a preliminary step of running a partial shape
519 matching algorithm, which is slow due to the optimization
520 at inference time. 3D-CODED [26] performs a template
521 based alignment to an input shape in two stages: fast in-
522 ference and a slow refinement through optimization. It is
523 designed for pairs which are either full or witness to mild
524 partiality, hence, to make the comparison more mean-
525 ingful we adjust their loss at the refinement phase to the more
526 suited directional Chamfer distance. FARM [46] is also an
527 alignment-based solution which has shown very impressive
528 results on shape completion and dense correspondences. It
529 builds on the SMPL [44] human body model due to its com-
530 pact parameterization, yet, we found it to be very slow to
531 converge (up to 30 min for a single shape) and prone to
532 getting trapped in local minima. 3D-EPN [19] is a rigid
533 shape completion method based on a voxelization approach,
534 utilizing a 3D CNN. Results are converted to a mesh via
535 computation of an isosurface. A classic Poisson reconstruc-
536 tion [36] is also provided as a naïve baseline with no access
537 to additional data other than the partial input.

540

4.3. Evaluation metrics

As previously discussed, it is challenging to define an analytic measurement for completion quality. Therefore we provide 5 different measurements, each reflecting a different perspective of the final completion as summarized in tables 1,2. First we report the mean square error (MSE) of the Euclidean distance between each point in the reconstructed shape and its ground truth mapping. This measure is reported only for template alignment methods for which the correspondence between the template and the ground truth reconstruction is defined. Next, we report the MSE of directional Chamfer distances: from the ground-truth to the completion and vice versa; The former measures the coverage of the ground-truth shape by the completion while the later penalizes outliers of the completion; We report the sum of both as full the Chamfer distance. Finally we measure the absolute error in the volume of the completion divided by the volume of the ground truth as a measure for volume deformation.

560
561

4.4. Single view completion

We evaluate the proposed method on the task of deformable shape completion on two datasets: FAUST and AMASS.

565

FAUST projections We follow the evaluation protocol proposed in [39] and summarize the completion results of our method and prior art in Table 1. As can be seen our network performs a much more accurate completion. Especially note the small standard deviation in the volumetric error that reflects our network has learned to keep the appearance of the input full shape. Contrary to optimization based methods [39, 25, 46] which are very slow at inference time, our feed-forward network performs inference in less than a second. To better appreciate the quality of our reconstructions, in Figure 3 we visualize several completions attained with various methods. Note the accurate preservation of intricate details which were completely lost in previous methods.

581

AMASS projections Using our generated set of partial shapes from AMASS described in 4.1, we compare our method with two recent methods based on shape alignment: 3D-CODED [26], and FARM [46]. As described in 4.2, 3D-CODED is based on a fixed template, and is not trained to handle severe partiality. It thus serves as a lower bound for our proposed method. FARM, on the other hand, was build for the same setting as ours. We summarize the results in Table 2. As can be seen, our method outperforms the two baselines by a large margin in all reported metrics. Note that on some of the examples (about 30%) FARM crashed during the optimization. We therefore only report the errors

on its successful runs. Particularly interesting is the large error and tiny standard deviation for 3D-CODED. This reflects the flattened results with consistently poor volumetric measure. Visualizations of several completions are shown in Figure 2.

4.5. Non-rigid partial correspondences

Finding dense correspondences between a full shape and its deformed parts is still a very much open research topic. Here we propose a solution in the form of alignment between the full shape and the partial shape, allowing for the recovery of the the correspondence by a simple nearest neighbour search. As before, we evaluate this task on both FAUST and AMASS data.

FAUST projections On the FAUST projections dataset, we compare with both alignment-based methods, FARM and 3D-CODED, as well as 3 methods designed to directly recover correspondences, i.e. without recovering the alignment: MoNet [48], and two 3-layered Euclidean CNN baselines, trained on either SHOT [63] descriptors or depth maps. Results are reported in Figure 4. The test set consists of a total of 200 shapes: 2 subjects at 10 different poses and 10 projected views. The direct matching baselines solve a classification problem for each shape vertex to a template shape. Differently, 3D-CODED has its own template used at train time. Note that since 3D-CODED was not build for severe partiality, we adjust its loss to a one-sided Chamfer distance in the refinement stage. Our method and FARM both require a complete shape, which we chose as the null pose of each of the test examples. Due to slow convergence and unstable behavior of FARM we only kept 20 useful matching results on which we report the performance. As can be seen from Figure 4, our method outperforms prior art by a significant margin. This result is particularly interesting since it demonstrates that even though we solve an alignment problem, which is a strictly harder problem than correspondence, we receive better results than methods that specialize in the latter. At the same time, looking at the poor performance demonstrated by the other alignment based methods, 3D-CODED and FARM, we conclude that simply solving an alignment problem is not enough and the details of our method and training scheme allow for a substantial difference.

4.6. Real scans

To evaluate our method in real world conditions, we test it on raw measurements taken during the preparation of the Dynamic FAUST [6] dataset. This use case nicely matches our setting. These are partial scans of a subject for which we have a complete reference shape at a different pose. We pre-process the input point cloud by extending it with estimation of its point normals using the method presented in

594

595

596

597

598

599

600

601

602

603

604

605

606

607

608

609

610

611

612

613

614

615

616

617

618

619

620

621

622

623

624

625

626

627

628

629

630

631

632

633

634

635

636

637

638

639

640

641

642

643

644

645

646

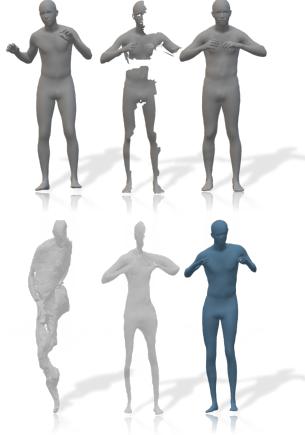
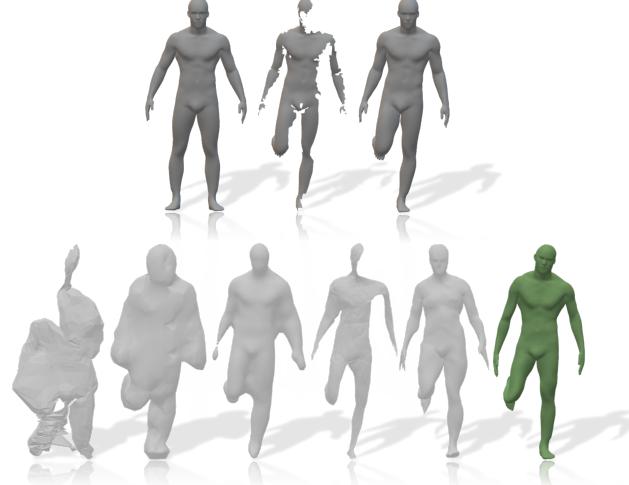
647

648	Error	Euclidean distance	Volumetric err.	Directional Chamfer distance	Directional Chamfer distance	Full Chamfer	702
649		GT and reconstruction [cm]	mean \pm std [%]	GT to reconstruction [cm]	reconstruction to GT [cm]	distance [cm]	703
650	Poisson [36]	23.73	24.8 \pm 23.2	7.3	3.64	10.94	704
651	3D-EPN [19]	23.5	89.7 \pm 33.8	4.52	4.87	9.39	705
652	3D-CODED [25]	35.50	21.8 \pm 0.3	11.15	38.49	49.64	706
653	FARM [46]	35.77	43.08 \pm 20.4	9.5	3.9	13.4	707
654	Litany <i>et al.</i> [39]	7.07	9.24 \pm 8.62	2.84	2.9	5.74	708
655	Ours	2.94	7.05 \pm 3.45	2.42	1.95	4.37	709

Table 1. FAUST Shape Completion. Comparison of different methods with respect to errors in vertex position and shape volume.

659	Error	Euclidean distance	Volumetric err.	Directional Chamfer distance	Directional Chamfer distance	Full Chamfer	710
660		GT and reconstruction [cm]	mean \pm std [%]	GT to reconstruction [cm]	reconstruction to GT [cm]	distance [cm]	711
661	3D-CODED [25]	36.14	14.84 \pm 8.02	13.65	35.35	49	712
662	FARM [46]	27.75	49.42 \pm 29.12	11.17	5.14	16.31	713
663	Ours	6.58	27.62 \pm 15.27	4.86	3.06	7.92	714

Table 2. AMASS Shape Completion. Comparison of different methods with respect to errors in vertex position and shape volume.

Figure 2. AMASS Shape Completion . At the top from left to right: full shape Q , partial shape P , ground truth completion R . At the bottom from left to right: reconstructions of FARM [46], 3D-CODED [25] and our own.Figure 3. FAUST Shape Completion . At the top from left to right: full shape Q , partial shape P , ground truth completion R . At the bottom - reconstructions from FARM [46], 3D-EPN [19], Poisson [36], 3D-CODED [26], Litany *et al.* [39] and our own.

[32]. The point cloud and the reference shape are subsequently inserted into a network pretrained on FAUST. The template, raw scan, and our reconstruction are shown (from left to right) in Figure 6. We show our result both as the recovered point cloud as well as the recovered mesh using the template triangulation. As apparent from the Figure, this is a challenging test case as it introduces several properties not seen at test time: a point cloud without connectivity leads to noisier normals, scanner noise, different point density and extreme partiality (note the missing bottom half of the shapes). Despite all these, our network was able to recover the input quite elegantly, preserving shape details and mimicking the desired pose. In the rightmost column we report a comparison with Litany *et al.* [39]. Note that while

[39] was trained on Dynamic FAUST, our network trained on FAUST which is severely constrained in its pose variability. The result highlights that our method favors realism and details in appearance over pose accuracy.

5. Concluding remarks

We have demonstrated that the problem of partial matching can be treated in a holistic manner when trying to fit a given part to a whole restricted to the pose inflicted by the part. By encoding the space distortion linking between parts at various poses to whole shapes given in other poses, we have been able to train a network to consider couples of parts and shapes at different poses, and match each part

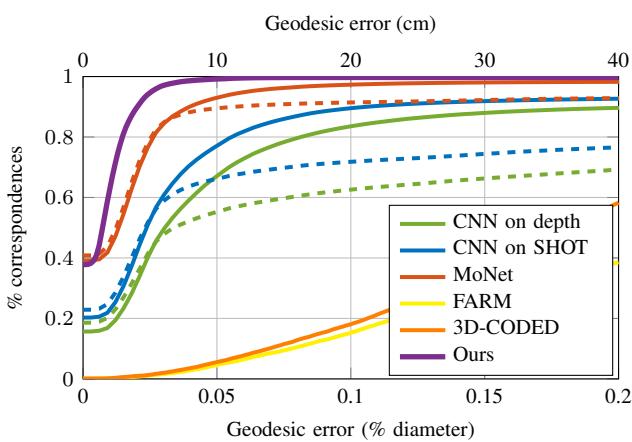


Figure 4. Generalization error on the FAUST dataset. Dashed line and solid line in the same color indicate performance before and after refinement, respectively. Note that our method doesn't require refinement, contributing to its computational speed.

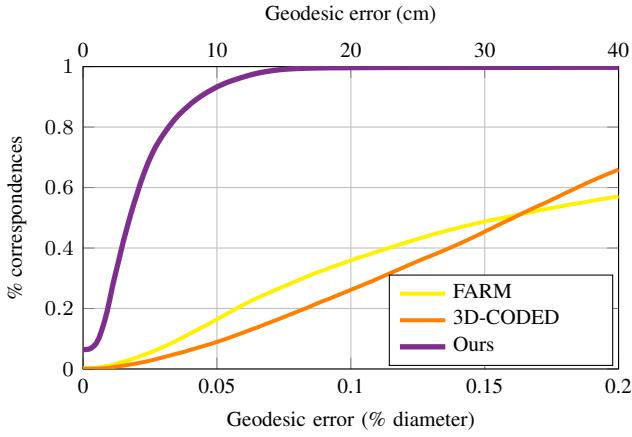


Figure 5. Generalization error on the Amass dataset

to its whole. This was realized by deforming the embedding space such that a resulting new whole shape aligns with, while obtaining the pose of, the given part. From Ancient Greek holistic philosophy, through modern psychology explanations of the human brain perception of shapes, we have demonstrated that computational matching procedures could benefit from the same axiomatic assumption stating that indeed *the whole is larger than the sum of its parts*.

References

- [1] Brett Allen, Brian Curless, Zoran Popović, and Aaron Hertzmann. Learning a correlated model of identity and pose-dependent body shape variation for real-time synthesis. In *Proceedings of the 2006 ACM SIGGRAPH/Eurographics symposium on Computer animation*, pages 147–156. Eurographics Association, 2006. 2
- [2] Dragomir Anguelov, Praveen Srinivasan, Daphne Koller, Sebastian Thrun, Jim Rodgers, and James Davis. Scape: shape completion and animation of people. In *ACM transactions on graphics (TOG)*, volume 24, pages 408–416. ACM, 2005. 2
- [3] Paul J Besl and Neil D McKay. Method for registration of 3-d shapes. In *Sensor fusion IV: control paradigms and data structures*, volume 1611, pages 586–606. International Society for Optics and Photonics, 1992. 5
- [4] Volker Blanz and Thomas Vetter. A morphable model for the synthesis of 3D faces. In *Proc. Computer Graphics and Interactive Techniques*, pages 187–194, 1999. 2
- [5] Federica Bogo, Javier Romero, Matthew Loper, and Michael J Black. FAUST: Dataset and Evaluation for 3d Mesh Registration. In *Proc. CVPR*, 2014. 5
- [6] Federica Bogo, Javier Romero, Gerard Pons-Moll, and Michael J. Black. Dynamic FAUST: Registering human bodies in motion. In *IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, July 2017. 6, 8
- [7] Davide Boscaini, Jonathan Masci, Emanuele Rodolà, and Michael Bronstein. Learning shape correspondence with anisotropic convolutional neural networks. In *Advances in Neural Information Processing Systems*, pages 3189–3197, 2016. 2
- [8] A. M. Bronstein, M. M. Bronstein, A.M. Bruckstein, and R. Kimmel. Matching two-dimensional articulated shapes using generalized multidimensional scaling. In *Proc. of Articulated Motion and Deformable Objects (AMDO)*, 2006. 2
- [9] A. M. Bronstein, M. M. Bronstein, and R. Kimmel. Expression-invariant 3d face recognition. In *Proc. Audio & Video-based Biometric Person Authentication (AVBPA), Lecture Notes in Comp. Science 2688, Springer*, 2003. 2

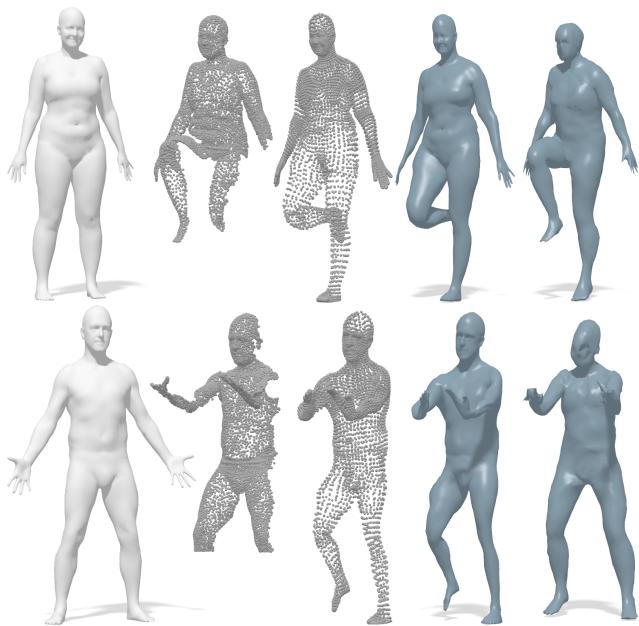


Figure 6. Completion from real scans from the Dynamic Faust dataset[6]. From left to right: Input reference shape; input raw scan; our completed shape as a point cloud; and as mesh; completion from Litany *et al.* [39].

- 864 [10] A. M. Bronstein, M. M. Bronstein, and R. Kimmel. Three- 918
865 dimensional face recognition. *International Journal of Computer 919
866 Vision*, 64(1):5–30, 2005. 2 920
867 [11] A. M. Bronstein, M. M. Bronstein, and R. Kimmel. Face2face: an isometric model for facial animation. In *Conf. 921
868 on Articulated Motion and Deformable Objects (AMDO)*, 922
869 2006. 2 923
870 [12] Alexander M Bronstein, Michael M Bronstein, and Ron 924
871 Kimmel. Generalized multidimensional scaling: a framework 925
872 for isometry-invariant partial surface matching. *PNAS*, 103(5):1168–1172, 2006. 3 926
873 [13] A. M. Bronstein, M. M. Bronstein, and R. Kimmel. Robust 927
874 expression-invariant face recognition from partially missing 928
875 data. In *Proc. ECCV, Graz, Austria*, May 2006. 2 929
876 [14] A. M. Bronstein, M. M. Bronstein, and R. Kimmel. Expression-invariant representations of faces. *IEEE Trans. 930
877 Image Processing*, 16(1):188–197, 2007. 2 931
878 [15] Michael M Bronstein, Joan Bruna, Yann LeCun, Arthur 932
879 Szlam, and Pierre Vandergheynst. Geometric deep learning: 933
880 going beyond euclidean data. *IEEE Signal Processing Magazine*, 34(4):18–42, 2017. 2 934
881 [16] Qifeng Chen and Vladlen Koltun. Robust nonrigid registration 935
882 by convex optimization. In *Proc. ICCV*, 2015. 3 936
883 [17] Christopher Choy, JunYoung Gwak, and Silvio Savarese. 4d 937
884 spatio-temporal convnets: Minkowski convolutional neural 938
885 networks. *arXiv preprint arXiv:1904.08755*, 2019. 2 939
886 [18] Luca Cosmo, Mikhail Panine, Arianna Rampini, Maks 940
887 Ovsjanikov, Michael M Bronstein, and Emanuele Rodolà. 941
888 Isospectralization, or how to hear shape, style, and correspondence. In *Proceedings of the IEEE Conference on 942
889 Computer Vision and Pattern Recognition*, pages 7529–7538, 943
890 2019. 3 944
891 [19] Angela Dai, Charles Ruizhongtai Qi, and Matthias Nießner. 945
892 Shape completion using 3D-encoder-predictor cnns and 946
893 shape synthesis. *arXiv:1612.00101*, 2016. 5, 7 947
894 [20] Y. Devir, G. Rosman, M. M. Bronstein, A. M. Bronstein, 948
895 and R. Kimmel. On reconstruction of non-rigid shapes with 949
896 intrinsic regularization. In *Proc. of Workshop on Nonrigid 950
897 Shape Analysis and Deformable Image Alignment (NOR- 951
898 DIA)*, 2009. 2 952
899 [21] Asi Elad and Ron Kimmel. Bending invariant representations 953
900 for surfaces. In *Proc. of CVPR’01, Hawaii*, December 954
901 2001. 2 955
902 [22] Asi Elad and Ron Kimmel. On bending invariant signatures 956
903 for surfaces. *IEEE Trans. on Pattern Analysis and Machine 957
904 Intelligence (PAMI)*, 25(10):1285–1295, 2003. 2 958
905 [23] Thomas Gerig, Andreas Morel-Forster, Clemens Blumer, 959
906 Bernhard Egger, Marcel Lüthi, Sandro Schönborn, and 960
907 Thomas Vetter. Morphable face models—an open framework. 961
908 *arXiv preprint arXiv:1709.08398*, 2017. 2 962
909 [24] Benjamin Graham, Martin Engelcke, and Laurens van der 963
910 Maaten. 3d semantic segmentation with submanifold sparse 964
911 convolutional networks. In *Proceedings of the IEEE Conference 965
912 on Computer Vision and Pattern Recognition*, pages 9224–9232, 2018. 2 966
913 [25] Thibault Groueix, Matthew Fisher, Vladimir G. Kim, Bryan 967
914 Russell, and Mathieu Aubry. 3D-CODED : 3D correspondences 968
915 by deep deformation. In *ECCV*, 2018. 6, 7 969
916 [26] Thibault Groueix, Matthew Fisher, Vladimir G. Kim, 970
917 Bryan C Russell, and Mathieu Aubry. 3D-coded: 3D cor- 971
918
919
920
921
922
923
924
925
926
927
928
929
930
931
932
933
934
935
936
937
938
939
940
941
942
943
944
945
946
947
948
949
950
951
952
953
954
955
956
957
958
959
960
961
962
963
964
965
966
967
968
969
970
971

- 972 [42] Or Litany, Emanuele Rodolà, Alex M Bronstein, and Michael M Bronstein. Fully spectral partial shape matching. *Computer Graphics Forum*, 36(2):247–258, 2017. 3
- 973 [43] Xingyu Liu, Mengyuan Yan, and Jeannette Bohg. Meteor-net: Deep learning on dynamic 3d point cloud sequences. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 9246–9255, 2019. 2
- 974 [44] Matthew Loper, Naureen Mahmood, Javier Romero, Gerard Pons-Moll, and Michael J Black. Smpl: A skinned multi-person linear model. *ACM transactions on graphics (TOG)*, 34(6):248, 2015. 2, 5
- 975 [45] Naureen Mahmood, Nima Ghorbani, Nikolaus F. Troje, Gerard Pons-Moll, and Michael J. Black. Amass: Archive of motion capture as surface shapes. In *The IEEE International Conference on Computer Vision (ICCV)*, Oct 2019. 5
- 976 [46] Riccardo Marin, Simone Melzi, Emanuele Rodolà, and Umberto Castellani. Farm: Functional automatic registration method for 3d human bodies. In *Computer Graphics Forum*. Wiley Online Library, 2018. 5, 6, 7
- 977 [47] Jonathan Masci, Davide Boscaini, Michael Bronstein, and Pierre Vandergheynst. Geodesic convolutional neural networks on riemannian manifolds. In *Proceedings of the IEEE international conference on computer vision workshops*, pages 37–45, 2015. 2
- 978 [48] Federico Monti, Davide Boscaini, Jonathan Masci, Emanuele Rodolà, Jan Svoboda, and Michael M Bronstein. Geometric deep learning on graphs and manifolds using mixture model cnns. In *Computer Vision and Pattern Recognition (CVPR), 2017 IEEE Conference on*, pages 5425–5434. IEEE, 2017. 2, 6
- 979 [49] Jeong Joon Park, Peter Florence, Julian Straub, Richard Newcombe, and Steven Lovegrove. Deepsdf: Learning continuous signed distance functions for shape representation. *arXiv preprint arXiv:1901.05103*, 2019. 2
- 980 [50] Adam Paszke, Sam Gross, Soumith Chintala, Gregory Chanan, Edward Yang, Zachary DeVito, Zeming Lin, Alban Desmaison, Luca Antiga, and Adam Lerer. Automatic differentiation in pytorch. 2017. 5
- 981 [51] Charles R Qi, Or Litany, Kaiming He, and Leonidas J Guibas. Deep hough voting for 3d object detection in point clouds. *arXiv preprint arXiv:1904.09664*, 2019. 2
- 982 [52] Charles R Qi, Hao Su, Kaichun Mo, and Leonidas J Guibas. Pointnet: Deep learning on point sets for 3d classification and segmentation. In *Proc. CVPR*, 2017. 2
- 983 [53] Charles R Qi, Li Yi, Hao Su, and Leonidas J Guibas. Pointnet++: Deep hierarchical feature learning on point sets in a metric space. *arXiv:1706.02413*, 2017. 2
- 984 [54] Arianna Rampini, Irene Tallini, Maks Ovsjanikov, Alex M Bronstein, and Emanuele Rodolà. Correspondence-free region localization for partial shape similarity via hamiltonian spectrum alignment. *arXiv preprint arXiv:1906.06226*, 2019. 3
- 985 [55] Elad Richardson, Matan Sela, Roy Or-El, and Kimmel Ron. Learning detailed face reconstruction from a single image. In *IEEE Conf. on Computer Vision and Pattern Recognition (CVPR), Hawaii, Honolulu*, 2017. 2
- 986 [56] Elad Richardson, Matan Sela, and Kimmel Ron. 3D face reconstruction by learning from synthetic data. In *4th Int. Conf. on 3D Vision (3DV) Stanford University, CA, USA*, 2016. 2
- 987 [57] Emanuele Rodolà, Luca Cosmo, Michael M Bronstein, Andrea Torsello, and Daniel Cremers. Partial functional correspondence. In *Computer Graphics Forum*, volume 36, pages 222–236. Wiley Online Library, 2017. 3
- 988 [58] Emanuele Rodolà, Samuel Rota Bulo, Thomas Windheuser, Matthias Vestner, and Daniel Cremers. Dense non-rigid shape correspondence using random forests. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 4177–4184, 2014. 3
- 989 [59] Javier Romero, Dimitrios Tzionas, and Michael J. Black. Embodied hands: Modeling and capturing hands and bodies together. *ACM Transactions on Graphics, (Proc. SIGGRAPH Asia)*, 36(6), Nov. 2017. 5
- 990 [60] Kripasindhu Sarkar, Kiran Varanasi, and Didier Stricker. Learning quadrangulated patches for 3D shape parameterization and completion. *arXiv:1709.06868*, 2017. 2
- 991 [61] Matan Sela, Elad Richardson, and Ron Kimmel. Unrestricted facial geometry reconstruction using image-to-image translation. In *Int. Conf. Comp. Vision (ICCV), Venice, Italy*, 2017. 2
- 992 [62] Hang Su, Subhransu Maji, Evangelos Kalogerakis, and Erik Learned-Miller. Multi-view convolutional neural networks for 3d shape recognition. In *Proc. CVPR*, 2015. 2
- 993 [63] Federico Tombari, Samuele Salti, and Luigi Di Stefano. Unique signatures of histograms for local surface description. In *International Conference on Computer Vision (ICCV)*, pages 356–369, 2010. 6
- 994 [64] Gul Varol, Duygu Ceylan, Bryan Russell, Jimei Yang, Ersin Yumer, Ivan Laptev, and Cordelia Schmid. Bodynet: Volumetric inference of 3d human body shapes. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 20–36, 2018. 2
- 995 [65] Gul Varol, Javier Romero, Xavier Martin, Naureen Mahmood, Michael J Black, Ivan Laptev, and Cordelia Schmid. Learning from synthetic humans. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 109–117, 2017. 2
- 996 [66] Nitika Verma, Edmond Boyer, and Jakob Verbeek. Dynamic filters in graph convolutional networks. *arXiv:1706.05206*, 2017. 2
- 997 [67] Yue Wang, Yongbin Sun, Ziwei Liu, Sanjay E Sarma, Michael M Bronstein, and Justin M Solomon. Dynamic graph cnn for learning on point clouds. *ACM Transactions on Graphics (TOG)*, 38(5):146, 2019. 2
- 998 [68] Lingyu Wei, Qixing Huang, Duygu Ceylan, Etienne Vouga, and Hao Li. Dense human body correspondences using convolutional networks. In *Proc. CVPR*, 2016. 2
- 999 [69] Zhirong Wu, Shuran Song, Aditya Khosla, Fisher Yu, Linguang Zhang, Xiaoou Tang, and Jianxiong Xiao. 3D shapenets: A deep representation for volumetric shapes. In *Proc. CVPR*, 2015. 2
- 1000 [70] Danfei Xu, Dragomir Anguelov, and Ashesh Jain. Pointfusion: Deep sensor fusion for 3d bounding box estimation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 244–253, 2018. 2

- 1080 [71] Y Y Ben-Shabat, M Lindenbaum, and Fischer A. 3D 1134
1081 point cloud classification and segmentation using 3D modified 1135
1082 fisher vector representation for convolutional neural networks. *arXiv preprint arXiv:1711.08241*, 2017. 2 1136
1083
1084 [72] Andrei Zanfir, Elisabeta Marinoiu, and Cristian Sminchisescu. Monocular 3d pose and shape estimation of 1137
1085 multiple people in natural scenes-the importance of multiple 1138
1086 scene constraints. In *Proceedings of the IEEE Conference 1139*
1087 on Computer Vision and Pattern Recognition, pages 2148– 1140
1088 2157, 2018. 2 1141
1089
1090
1091
1092
1093
1094
1095
1096
1097
1098
1099
1100
1101
1102
1103
1104
1105
1106
1107
1108
1109
1110
1111
1112
1113
1114
1115
1116
1117
1118
1119
1120
1121
1122
1123
1124
1125
1126
1127
1128
1129
1130
1131
1132
1133