

# Who Left the Dogs Out? 3D Animal Reconstruction with Expectation Maximization in the Loop

Benjamin Biggs<sup>1</sup>, Oliver Boyne<sup>1</sup>, James Charles<sup>1</sup>,  
Andrew Fitzgibbon<sup>2</sup>, and Roberto Cipolla<sup>1</sup>

<sup>1</sup> Department of Engineering, University of Cambridge, Cambridge, UK  
{bjb56,ob312,jjc75,rc10001}@cam.ac.uk

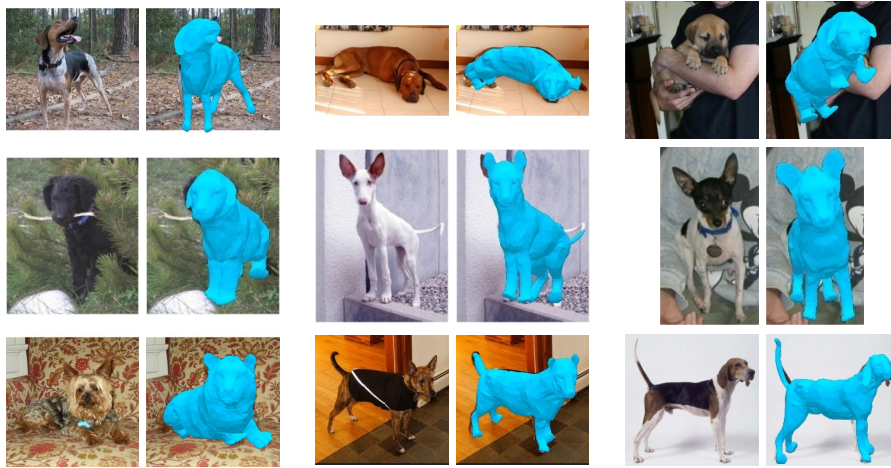
<sup>2</sup> Microsoft, Cambridge, UK awf@microsoft.com

**Abstract.** We introduce an automatic, end-to-end method for recovering the 3D pose and shape of dogs from monocular internet images. The large variation in shape between dog breeds, significant occlusion and low quality of internet images makes this a challenging problem. We learn a richer prior over shapes than previous work, which helps regularize parameter estimation. We demonstrate results on the Stanford Dog Dataset, an “in-the-wild” dataset of 20,580 dog images for which we have collected 2D joint and silhouette annotations to split for training and evaluation. In order to capture the large shape variety of dogs, we show that the natural variation in the 2D dataset is enough to learn a detailed 3D prior through expectation maximisation (EM). As a by-product of training, we generate a new parameterized model (including limb scaling) SMBLD which we release alongside our new annotation dataset *StanfordExtra* to the research community. Code and data are available at <https://sites.google.com/view/wldo>.

## 1 Introduction

Animals contribute greatly to our society, in numerous ways economic and otherwise (there are more than 63 million pet dogs in the US alone [3]). In consequence, there has been considerable attention in the computer vision research community to the interpretation of imagery of animals. Although these techniques share similarities to techniques for understanding images of humans, a key difference is that obtaining labelled training data for animals is more difficult than for humans, because of the wide range of shapes and species of animals, and the difficulty of educating manual labellers in animal physiology.

A particular species of interest is the dog, however it is noticeable that existing work has not yet demonstrated effective 3D reconstruction of dogs over large test sets. We postulate that this is partially because dog breeds are remarkably dissimilar in shape and texture, presenting a challenge to the current state of the art. The methods we propose extend the state of the art in several ways. While each of these qualities exist in some existing works, we believe ours is the



**Fig. 1. End-to-end 3D Dog Reconstruction from monocular images.** We propose a novel method that, given a monocular image of a dog can predict a set of parameters for our SMULD 3D dog model which is consistent with the input. We regularize learning using a multi-modal shape prior, which is tuned during training with an expectation maximization scheme.

first to exhibit this combination, leading to a new state of the art in terms of scale and object diversity.

1. We reconstruct pose and shape on a test set of 1703 low-quality internet images of a complex 3D object class (dogs).
2. We directly regress to object pose and shape from a single image without a model fitting stage.
3. We use easily obtained 2D annotations in training, and none at test time.
4. We incorporate fitting of a new multi-modal prior into the training phase (via EM update steps), rather than fitting it to 3D data as in previous work.
5. We introduce new degrees of freedom to the SMAL model, allowing explicit scaling of subparts.

### 1.1 Related work

The closest work in terms of scale is the category-specific mesh reconstruction of Kanazawa et al. [15], where 2850 images of birds were reconstructed. However doing so for the complex pose and shape variations of dogs required the advances described in this paper.

Table 1 summarizes previous work on animal reconstruction. It is interesting to note that while several papers demonstrate reconstruction across species, which *prima facie* is a richer class than just dogs, the test-time requirements (e.g. manually-clicked keypoints/silhouette segmentations, input image quality etc.) are considerably higher for those systems. Thus we claim that the achievement of reconstructing a full range of dog breeds, with variable fur length, varying shape and pose of ears, and with considerable occlusion, is a significant contribution.

Paper	Animal Class	Training requirements	Template Model	Video required	Test Time Annotation	Model Fitting	Test Size
This paper	Dogs	J2, S2, T3, P3	SMAL	No	None	No	1703
3D-Safari [33]	Zebras, horses	M3 (albeit synthetic), J2, S2, P3	SMAL	3-7 frames / animal	None	Yes	200
Lions and Tigers and Bears (LTB) [34]	MLQ	Not trained	SMAL	3-7 frames / animal	J2, S2	Yes	14
3D Menagerie (3D-M) [35]	MLQ	Not trained	SMAL	No	J2, S2	Yes	48
Creatures Great and SMAL (CGAS) [5]	MLQ	Not trained	SMAL	Yes	S2 (for best results shown)	Yes	9
Category Specific Mesh Reconstructions [15]	Birds	J2, S2	Bird convex hull	No	None	No	2850
What Shape are Dolphins [7]	Dolphins, Pigeons	Not trained	Dolphin Template	25 frames / category	J2, S2	Yes	25
Animated 3D Creatures [29]	MLQ	Not trained	Generalized Cylinders	Yes	J2, S2	Yes	15

**Table 1.** Literature summary: Our paper extends large-scale “in-the-wild” reconstruction to the difficult class of diverse breeds of dogs. MLQ: Medium-to-large quadrupeds. J2: 2D Joints. S2: 2D Silhouettes. T3: 3D Template. P3: 3D Priors. M3: 3D Model.

**Monocular 3D reconstruction of human bodies** The majority of recent work in 3D pose and shape recovery from monocular images tackles the special case of 3D *human* reconstruction. As a result, the research community has collected a multitude of open source human datasets which provide strong supervisory signals for training deep neural networks. These include accurate 3D deformable template models [23] generated from real human scans, 3D motion capture datasets [11,24] and large 2D datasets [22,12,4] which provide keypoint and silhouette annotations.

The abundance of available human data has supported the development of successful monocular 3D reconstruction pipelines [21,13]. Such approaches rely on accurate 3D data to build detailed priors over the distribution of human shapes and poses, and use large 2D keypoints datasets to promote generalization to “in-the-wild” scenarios. Silhouette data has also been shown to assist in accurate reconstruction of clothes, hair and other appearance detail [30,2]. While the dominant paradigm in human reconstruction is now end-to-end deep learning methods, SPIN [20] show impressive improvement by incorporating an energy minimization process within their training loop to further minimize a 2D reprojection loss subject to fixed pose & shape priors. Inspired by this innovation, we learn an iteratively-improving shape prior by applying expectation maximization during the training process.

**Monocular 3D reconstruction of animal categories.** While animals are often featured in computer vision literature, there are still relatively few works that focus on accurate 3D animal reconstruction.

A primary reason for this is absence of large scale 3D datasets<sup>3</sup> stemming from the practical challenges associated with 3D motion capture, as well as a lack of 2D data which captures a wide variety of animals. The recent Animal Pose dataset [6] is one such 2D alternative, but contains significantly fewer labelled images than our new StanfordDogs dataset (4,000 compared to 20,580 in ). On the other hand, animal silhouette data is plentiful [22,9,18].

Zuffi et al. [35] made a significant contribution to 3D animal reconstruction research by releasing SMAL, a deformable 3D quadruped model (analagous to SMPL [23] for human reconstruction) from 41 scans of artist-designed toy figurines. The authors also released shape and pose priors generated from artist data. In this work we develop *SMBLD*, an extension of SMAL that better represents the diverse dog category by adding scale parameters and refining the shape prior using our large image dataset.

While there have been various “model-free” approaches which do not rely on an initial template model to generate the 3D animal reconstruction, these techniques often do not produce a mesh [1,26] or rely heavily on input 2D keypoints or video at test-time [31,28]. An exception is the end-to-end network of Kanazawa et al. [15], although we argue that the bird category exhibits more limited articulation than our dog category.

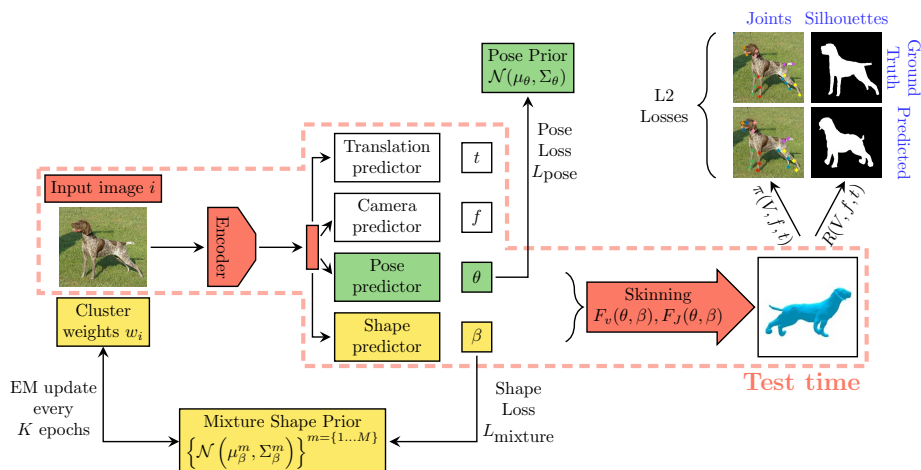
We instead focus on model-based approaches. The SMAL authors [35] demonstrate fitting their deformable 3D model to quadruped species using user-provided keypoint and silhouette dataset. Lions and Tigers and Bears (LTB) [34] then demonstrated fitting to broader animal categories by incorporating multi-view constraints from video sequences. Biggs et al. [5] overcame the need for hand-clicked keypoints by training a joint predictor on synthetic data. 3D-Safari [33] further improve by training a deep network on synthetic data (generated using the LTB method [34]) to recover detailed zebra shapes in the wild.

A drawback of these approaches is their reliance on a test-time energy-based optimization procedure, which is susceptible to failure with poor quality keypoint/silhouette predictions and increases the computational burden. By contrast our method requires no additional energy-based refinement, and is trained purely from single in-the-wild images. The experimental section of this paper contains a robust comparison between our end-to-end method and relevant optimization-based approaches.

A major impediment to research in 3D animal reconstruction has been the lack of a strong evaluation benchmark, with most of the above methods showing only qualitative evaluations or providing quantitative results on fewer than 50 examples. To remedy this, we introduce *StanfordExtra*, a new large-scale dataset which we hope will drive further progress in the field.

---

<sup>3</sup> Released after the submission of this paper, RGBD-Dog dataset [17] is the first open-source 3D motion capture dataset for dogs.



**Fig. 2.** Our method consists of (1) a deep CNN encoder which condenses the input image into a feature vector (2) a set of prediction heads which generate SMBLD parameters for shape  $\beta$ , pose  $\theta$ , camera focal length  $f$  and translation  $t$  (3) skinning functions  $F_v$  and  $F_J$  which construct the mesh from a set of parameters, and (4) loss functions which minimise the error between projected and ground truth joints and silhouettes. Finally, we incorporate a mixture shape prior (5) which regularises the predicted 3D shape and is iteratively updated during training using expectation maximisation. At test time, our system (1) condenses the input image, (2) generates the SMBLD parameters and (3) constructs the mesh.

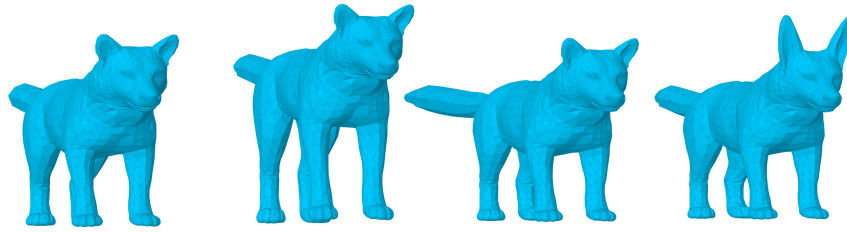
## 2 Parametric animal model

At the heart of our method is a parametric representation of a 3D animal mesh, which is based on the Skinned Multi-Animal Linear (SMAL) model proposed by [35]. SMAL is a deformable 3D animal mesh parameterized by shape and pose. The *shape*  $\beta \in \mathbb{R}^B$  parameters are PCA coefficients of an undeformed template mesh with limbs in default position. The *pose*  $\theta \in \mathbb{R}^P$  parameters meanwhile govern the joint angle rotations ( $35 \times 3$  Rodrigues parameters) which effect the articulated limb movement. The model consists of a linear blend skinning function  $F_v : (\theta, \beta) \mapsto V$ , which generates vertex positions  $V \in \mathbb{R}^{3889 \times 3}$ , and a joint function  $F_J : (\theta, \beta) \mapsto J$ , which generates joint positions  $J \in \mathbb{R}^{35 \times 3}$ .

### 2.1 Introducing scale parameters

While SMAL has been shown to be adequate for representing a variety of quadruped types, we find that the modes of dog variation are poorly captured by the current model. This is unsurprising, since SMAL used only four dogs in its construction.

We therefore introduce a simple but effective way to improve the model’s representational power over this particularly diverse animal category. We augment the set of shape parameters  $\beta$  with an additional set  $\kappa$  which independently



**Fig. 3. Effect of varying SMBLD scale parameters.** *From left to right: Mean SMBLD model, 25% leg elongation, 50% tail elongation, 50% ear elongation.*

scale parts of the mesh. For each model joint, we define parameters  $\kappa_x, \kappa_y, \kappa_z$  which apply a local scaling of the mesh along the local coordinate  $x, y, z$  axes, before pose is applied. Allowing each joint to scale entirely independently can however lead to unrealistic deformations, so we share scale parameters between multiple joints, e.g. leg lengths. The new Skinned Multi-Breed Linear Model for Dogs (SMBLD) is therefore adapted from SMAL by adding 6 scale parameters to the existing set of shape parameters. Figure 3 shows how introducing scale parameters increases the flexibility of the SMAL model. We also extend the provided SMAL shape prior (which later initializes our EM procedure) to cover the new scale parameters by fitting SMBLD to a set of 13 artist-designed 3D dog meshes. Further details left to the supplementary.

### 3 End-to-end dog reconstruction from monocular images

We now consider the task of reconstructing a 3D dog mesh from a monocular image. We achieve this by training an end-to-end convolutional network that predicts a set of SMBLD model and perspective camera parameters. In particular, we train our network to predict pose  $\theta$  and shape  $\beta$  SMBLD parameters together with translation  $t$  and focal length  $f$  for a perspective camera. A complete overview of the proposed system is shown in Figure 2.

#### 3.1 Model architecture

Our network architecture is inspired by the model of 3D-Safari [33]. Given an input image cropped to  $(224, 224)$ , we apply a Resnet-50 [10] backbone network to encode a 1024-dimensional feature map. These features are passed through various linear prediction heads to produce the required parameters. The pose, translation and camera prediction modules follow the design of 3D-Safari, but we describe the differences in our shape module.

**Pose, translation and camera prediction.** These modules are independent multi-layer perceptrons which map the above features to the various parameter types. As with 3D-Safari we use two linear layers to map to a set of  $35 \times 3$  3D pose parameters (three parameters for each joint in the SMBLD kinematic tree) given in Rodrigues form. We use independent heads to predict camera frame

translation  $t_{x,y}$  and depth  $t_z$  independently. We also predict the focal length of the perspective camera similarly to 3D-Safari.

**Shape and scale prediction.** Unlike 3D-Safari, we design our network to predict the set of shape parameters (including scale) rather than vertex offsets. We observe improvement by handling the standard 20 blend-shape parameters and our new scale parameters in separate linear prediction heads. We retrieve the scale parameters by  $\kappa = \exp x$  where  $x$  are the network predictions, as we find predicting log scale helps stabilise early training.

### 3.2 Training losses

A common approach for training such an end-to-end system would be to supervise the prediction of  $(\theta, \beta, t, f)$  with 3D ground truth annotations [20,14,27]. However, building a suitable 3D annotation dataset would require an experienced graphics artist to design an accurate ground truth mesh for each of 20,520 StanfordExtra dog images, a prohibitive expense.

We instead develop a method that instead relies on *weak 2D supervision* to guide network training. In particular, we rely on only 2D keypoints and silhouette segmentations, are significantly cheaper to obtain.

The rest of this section describes the set of losses used to supervise the network at train time.

**Joint reprojection.** The most important loss to promote accurate limb positioning is the joint reprojection loss  $L_{\text{joints}}$  which compares the projected model joints  $\pi(F_J(\theta, \beta), t, f)$  to the ground truth annotations  $\hat{X}$ . Given the parameters predicted by the network, we apply the SMLD model to transform the pose and shape parameters into a set of 3D joint positions  $J \in \mathbb{R}^{35 \times 3}$ , and project them to the image plane using translation and camera parameters. The joint loss  $L_{\text{joints}}$  is given by the  $\ell_2$  error between the ground truth and projected joints:

$$L_{\text{joints}}(\theta, \beta, t, f; \hat{X}) = \|\hat{X} - \pi(F_J(\theta, \beta), t, f)\|_2 \quad (1)$$

Note that many of our training images exhibit significant occlusion, so  $\hat{X}$  contains many invisible joints. We handle this by masking  $L_{\text{joints}}$  to prevent invisible joints contributing to the loss.

**Silhouette loss.** The silhouette loss  $L_{\text{sil}}$  is used to promote shape alignment between the SMLD dog mesh and the input dog. In order to compute the silhouette loss, we define a rendering function  $R : (\nu, t, f) \mapsto S$  which projects the SMLD mesh to produce a binary segmentation mask. In order to allow derivatives to be propagated through  $R$ , we implement  $R$  using the differentiable Neural Mesh Renderer [16]. The loss is computed as the  $\ell_2$  difference between a projected silhouette and the ground truth mask  $\hat{S}$ :

$$L_{\text{sil}}(\theta, \beta, t, f; \hat{S}) = \|\hat{S} - R(F_V(\theta, \beta), t, f)\|_2 \quad (2)$$

**Priors.** In the absence of 3D ground truth training data, we rely on priors obtained from artist graphics models to encourage realism in the network predictions. We model both pose and shape using a multivariate Gaussian prior, consisting of means  $\mu_\theta, \mu_\beta$  and covariance matrices  $\Sigma_\theta, \Sigma_\beta$ . The loss is given as the log likelihood of a given shape or pose vector under these distributions, which corresponds to the Mahalanobis distance between the predicted parameters and their corresponding means:

$$L_{\text{pose}}(\theta; \mu_\theta, \Sigma_\theta) = (\theta - \mu_\theta)^T \Sigma_\theta^{-1} (\theta - \mu_\theta) \quad (3)$$

$$L_{\text{shape}}(\beta; \mu_\beta, \Sigma_\beta) = (\beta - \mu_\beta)^T \Sigma_\beta^{-1} (\beta - \mu_\beta) \quad (4)$$

Unlike previous work, we find there is no need to use a loss to penalize pose parameters if they exceed manually specified joint angle limits. We suspect our network learns this regularization naturally because of our large dataset.

### 3.3 Learning a multi-modal shape prior.

The previous section introduced a unimodal, multivariate Gaussian shape prior, based on mean  $\mu_\beta$  and covariance matrix  $\Sigma_\beta$ . However, we find enforcing this prior throughout training tends to result in predictions which appear similar in 3D shape, even when tested on dog images of different breeds. We propose to improve diversity among predicted 3D dog shapes by extending the above formulation to a Mixture of  $M$  Gaussians prior. The mixture shape loss is then given as:

$$L_{\text{mixture}}(\beta_i; \mu_\beta, \Sigma_\beta, \Pi_\beta) = \sum_{m=1}^M \Pi_\beta^m L_{\text{shape}}(\beta_i; \mu_\beta^m, \Sigma_\beta^m) \quad (5)$$

Where  $\mu_\beta^m, \Sigma_\beta^m$  and  $\Pi_\beta^m$  are the mean, covariance and mixture weight respectively for Gaussian component  $m$ . For each component the mean is sampled from our existing unimodal prior and the covariance is set equal to the unimodal prior i.e.  $\Sigma_\beta^m := \Sigma_\beta$ . All mixture weights are initially set to  $\frac{1}{M}$ .

Each training image  $i$  is assigned a set of latent variables  $\{w_i^1, \dots, w_i^M\}$  encoding the probability of the dog shape in image  $i$  being generated by component  $m$ .

### 3.4 Expectation Maximization in the loop

As previously discussed, our initial shape prior is obtained from artist data which we find is unrepresentative of the diverse shapes present in our real dog dataset. We address this by proposing to recover the latent variables  $w_i^m$  and parameters  $(\mu_\beta^m, \Sigma_\beta^m$  and  $\Pi_\beta^m)$  of our 3D shape prior by learning from monocular images of in-the-wild dogs and their 2D training labels in our training dataset.

We achieve this using Expectation Maximization (EM), which regularly updates the means and variances for each mixture component and per-image mixture weights based on the observed shapes in the training set. While training our 3D reconstruction network, we progressively update our shape mixture model with an alternating ‘E’ step and ‘M’ step described below:



**The ‘E’ Step.** The ‘E’ step computes the expected value of the latent variables  $w_i^m$  assuming fixed  $(\mu_\beta^m, \Sigma_\beta^m, \Pi_\beta^m)$  for all  $i \in \{1, \dots, N\}, m \in \{1, \dots, M\}$ .

The update equation for an image  $i$  with latest shape prediction  $\beta_i$  and cluster  $m$  with parameters  $(\mu_\beta^m, \Sigma_\beta^m, \Pi_\beta^m)$  is given as:

$$w_i^m := \frac{\mathcal{N}(\beta_i | \mu_\beta^m, \Sigma_\beta^m) \Pi_\beta^m}{\sum_{m'}^M \mathcal{N}(\beta_i | \mu_\beta^{m'}, \Sigma_\beta^{m'}) \Pi_\beta^{m'}} \quad (6)$$

**The ‘M’ Step.** The ‘M’ step computes new values for  $(\mu_\beta^m, \Sigma_\beta^m, \Pi_\beta^m)$ , assuming fixed  $w_i^m$  for all  $i \in \{1, \dots, N\}, m \in \{1, \dots, M\}$ .

The update equations are given as follows:

$$\mu_\beta^m := \frac{\sum_i w_i^m \beta_i}{\sum_i w_i^m} \quad \Sigma_\beta^m := \frac{\sum_i w_i^m (\beta_i - \Sigma_\beta^m) (\beta_i - \Sigma_\beta^m)^T}{\sum_i w_i^m} \quad \Pi_\beta^m := \frac{1}{N} \sum_i w_i^m \quad (7)$$

## 4 Experiments

In this section we compare our method to competitive baselines. We begin by describing our new large-scale dataset of annotated dog images, followed by a quantitative and qualitative evaluation.

### 4.1 StanfordExtra: A new large-scale dog dataset with 2D keypoint and silhouette annotations



**Fig. 4. StanfordExtra example images.** *Left:* outlined segmentations and labelled keypoints for 24 representative images. *Right:* heatmap of deviation of worker submitted results from mean for each submission.

In order to evaluate our method, we introduce *StanfordExtra*: a new large-scale dataset with annotated 2D keypoints and binary segmentation masks for dogs. We opted to take source images from the existing Stanford Dog Dataset [19], which consists of 20,580 dog images taken “in the wild” and covers 120 dog

breeds. The dataset contains vast shape and pose variation between dogs, as well as nuisance factors such as self/environmental occlusion, interaction with humans/other animals and partial views. Figure 4 (left) shows samples from the new dataset.

We used Amazon Mechanical Turk to collect a binary silhouette mask and 20 keypoints per image: 3 per leg (knee, ankle, toe), 2 per ear (base, tip), 2 per tail (base, tip), 2 per face (nose and jaw). We can approximate the difficulty of the dataset by analysing the variance between 3 annotators at both the joint labelling and silhouette task. Figure 4 (right) illustrates typical per-joint variance in joint labelling. Further details of the data curation procedure are left to the supplementary materials.

## 4.2 Evaluation protocol

Our evaluation is based on our new StanfordExtra dataset. In line with other “in-the-wild” 3D reconstruction methods tackling articulated subjects [20,21], we filter images from the original dataset of 20,580 for which the majority of dog keypoints are invisible. We consider these images unsuitable for our full-body dog reconstruction task. We also remove images for which the consistency in keypoint/silhouette segmentations between the 3 annotators is below a set threshold. This leaves us with 8,476 images which we divide per-breed into an 80%/20% train and test split.

We consider two primary evaluation metrics. IoU is the intersection-over-union of the projected model silhouette compared to the ground truth annotation and indicates the quality of the reconstructed 3D shape. Percentage of Correct Keypoints (PCK) computes the percentage of joints which are within a normalized distance (based on square root of 2D silhouette area) to the ground truth locations, and evaluates the quality of reconstructed 3D pose. We also produce PCK results on various joint groups (legs, tail, ears, face) to compare the reconstruction accuracy for different parts of the dog model.

## 4.3 Training procedure

We train our model in two stages. The first omits the silhouette loss which we find can lead the network to unsatisfactory local minima if applied too early. With the silhouette loss turned off, we find it satisfactory to use the simple unimodal prior (and without EM) for this preliminary stage since there is no loss to specifically encourage a strong shape alignment. After this, we introduce the silhouette loss, the mixture prior and begin applying the expectation maximization updates over  $M = 10$  clusters. We train the first stage for 250 epochs, the second stage for 150 and apply the EM step every 50 epochs. All losses are weighted, as described in the supplementary. The entire training procedure takes 96 hours on a single P100 GPU.

## 4.4 Comparison to baselines

We first compare our method to various baseline methods. 3D Menagerie (3D-M) [35] is an approach which fits the 3D SMAL model using per-image energy

minimization. Creatures Great and SMAL (CGAS) [5] is a three-stage method, which employs a joint predictor on silhouette renderings from synthetic 3D dogs, applies a genetic algorithm to clean predictions, and finally applies the SMAL optimizer to produce the 3D mesh.

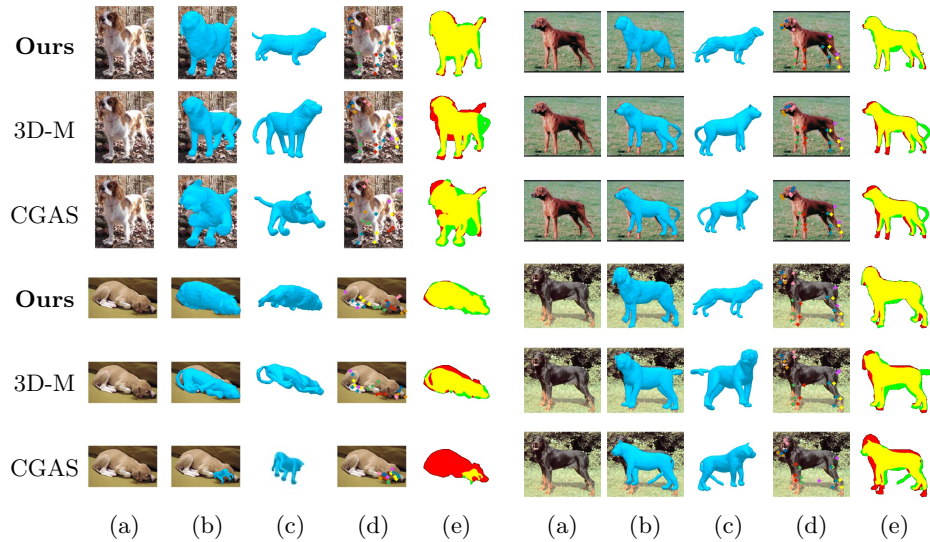
At test-time both 3D-M and CGAS rely on manually-provided segmentation masks, and 3D-M also relies on hand-clicked keypoints. In order to produce a fair comparison, we produce a set of *predicted* keypoints for StanfordExtra by training the Stacked Hourglass Network [25] with 8 stacks and 1 block, and *predicted* segmentation masks using DeepLab v3+ [8]. The Stacked Hourglass Network achieves 71.4% PCK score, DeepLab v3+ achieves 83.4% IoU score and the CGAS joint predictor achieves 41.8% PCK score.

Table 2 and Figure 5 show the comparison against competitive methods. For full examination, we additionally provide results for 3D-M and CGAS in the scenario that ground-truth keypoints and/or segmentations are available at test time.

The results show our end-to-end method outperforms the competitors when they are provided with predicted keypoints/segmentations (white rows). Our method therefore achieves a new state-of-the-art on this 3D reconstruction task. In addition, we show our method achieves improved average IoU/PCK scores than competitive methods, even when they are provided ground truth annotations at test time (grey rows). We also demonstrate wider applicability of two contributions from our work (scale parameters and improved prior) by showing improved performance of the 3D-M method when these are incorporated. Finally, our model’s test-time speed is significantly faster than the competitors as it does not require an optimizer.

Method	Kps	Seg	IoU		PCK @ 0.15			
			Avg	Legs	Tail	Ears	Face	
3D-M [35]	Pred	Pred	69.9	69.7	68.3	68.0	57.8	93.7
3D-M	GT	GT	71.0	75.6	74.2	89.5	60.7	98.6
3D-M	GT	Pred	70.7	75.5	74.1	88.1	60.2	98.7
3D-M	Pred	GT	70.5	70.3	69.0	69.4	58.5	94.0
CGAS [5]	CGAS	Pred	63.5	28.6	30.7	34.5	25.9	24.1
CGAS	CGAS	GT	64.2	28.2	30.1	33.4	26.3	24.5
3D-M + scaling	Pred	Pred	70.4	70.9	69.8	66.9	59.7	94.0
3D-M + scaling + EM prior	Pred	Pred	71.8	73.4	72.5	<b>70.3</b>	62.6	<b>94.1</b>
<b>Ours</b>	—	—	<b>74.2</b>	<b>78.8</b>	<b>76.4</b>	63.9	<b>78.1</b>	92.1

**Table 2. Baseline comparisons.**<sup>4</sup> PCK and silhouette IOU scores are shown for SOTA methods under varying conditions. Directly comparable baseline methods (requiring only an input image) are highlighted. *Pred* keypoints generated with Hourglass-Net [25] and segmentations with DeepLab v3+ [8]. 3D-M/CGAS are also analysed when they have access to ground-truth keypoints and/or segmentation masks. We also analyse adding this paper’s innovations (scale parameters and EM prior) to 3D-M [35].



**Fig. 5. Qualitative comparison to SOTA.** Row 1: **Ours**, Row 2: 3D-M [35], Row 3: CGAS [5]. (a) input image, (b) predicted 3D mesh, (c) canonical view 3D mesh, (d) reprojected model joints and (e) silhouette reprojection error.

#### 4.5 Generalization to unseen dataset

Table 3 shows an experiment to compare how well our model generalizes to a new data domain. We test our model against the 3D-M [35] method (using predicted keypoints and segmentations as above for fairness) on the recent Animal Pose dataset [6]. The data preparation process is the same as for StanfordExtra and no fine-tuning was used for either method. We achieve good results in this unseen domain and still improve over the 3D-M optimizer.

#### 4.6 Ablation study

We also produce a study in which we ablate individual components of our method and examine the effect on the PCK/IoU performance. We evaluate three variants: (1) **Ours w/o EM** that omits EM updates, (2) **Ours w/o MoG** which replaces our mixture shape prior with a unimodal prior, (3) **Ours w/o Scale** which removes the scale parameters.

The results in Table 4 indicate that each individual component has a positive impact on the overall method performance. In particular, it can be seen that the inclusion of the EM and Mixture of Gaussians prior leads to an improvement in IoU, suggesting that the shape prior refinements steps help the model accurately fit the exact dog shape. Interestingly, we notice that adding the Mixture of Gaussians prior but omitting EM steps slightly hinders performance, perhaps due to an sub-optimal initialization for the  $M$  clusters. However, we find adding EM updates to the Mixture of Gaussian model improves all metrics except the ear keypoint accuracy. We observe the error here is caused by the our shape

Method	IoU	PCK @ 0.15				
		Avg	Legs	Tail	Ears	Face
3D-M [35]	64.9	59.2	55.7	56.9	61.3	<b>86.7</b>
<b>Ours</b>	<b>67.5</b>	<b>67.6</b>	<b>60.4</b>	<b>62.7</b>	<b>86.0</b>	<b>86.7</b>

**Table 3. Animal Pose dataset [6] results<sup>4</sup>.** Evaluation on recent Animal Pose dataset with no fine-tuning to our method nor joint/silhouette predictors used for 3D-M.

Method	IoU	PCK @ 0.1				
		Avg	Legs	Tail	Ears	Face
<b>Ours</b>	<b>74.2</b>	<b>63.7</b>	<b>59.5</b>	<b>48.1</b>	60.1	88.0
–EM	68.7	63.2	58.8	44.5	<b>62.6</b>	87.6
–MoG	69.0	63.1	<b>59.5</b>	40.0	60.0	<b>89.5</b>
–Scale	68.3	60.1	58.2	45.2	50.5	88.3

**Table 4. Ablation study.<sup>4</sup>** Evaluation with the following contributions removed: (a) EM updates, (b) Mixture Shape Prior, (c) SMLD scale parameters.

prior learning slightly imprecise shapes for dogs with extremely “floppy” ears. Although there is good silhouette coverage for these regions, the fact our model has only a single articulation point per ear causes a lack of flexibility that results in occasionally misplaced ear tips for these instances. This could be improved in future work by adding additional model joints to the ear. Finally, we find the increased model flexibility afforded by the SMLD scale parameters have a positive effect on IoU/PCK scores.

#### 4.7 Qualitative evaluation

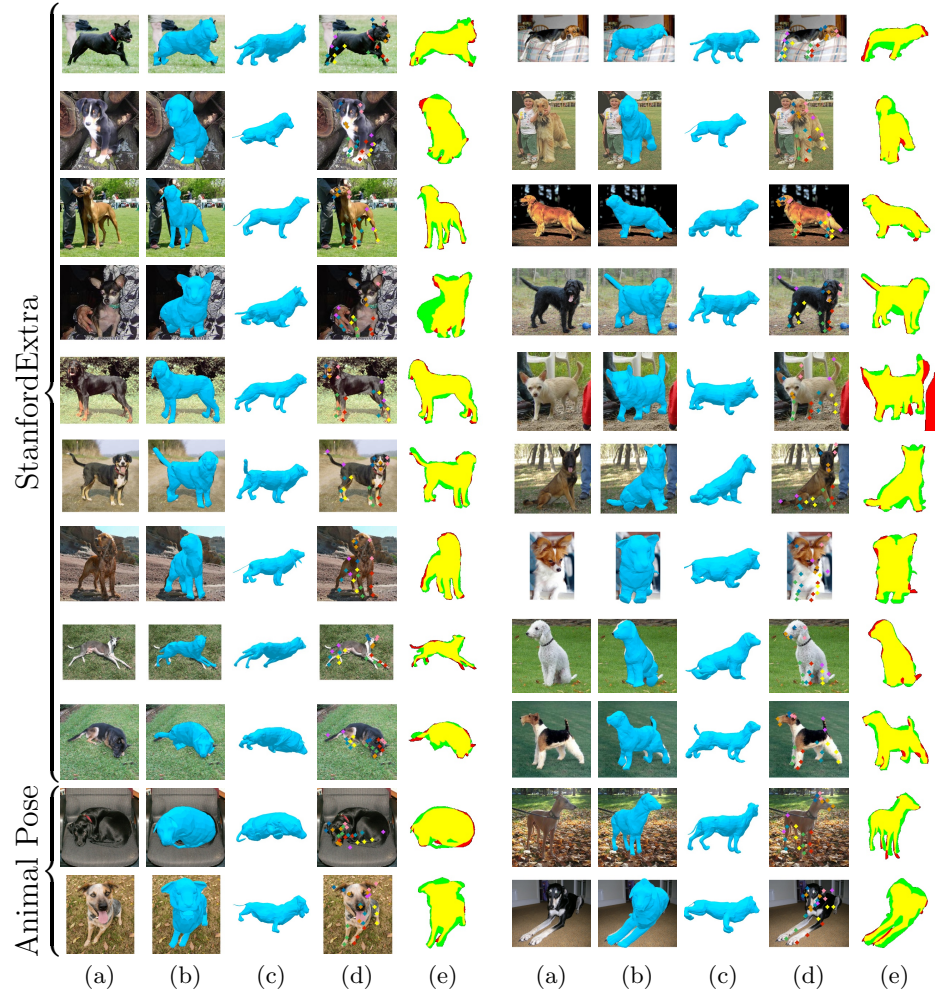
Figure 5 shows a range of example system outputs when tested on range of StanfordExtra and Animal Pose [6] dogs with varying pose and shape and in challenging conditions. Note that only StanfordExtra is used for training.

## 5 Conclusions

This paper presents an end-to-end method for automatic, monocular 3D dog reconstruction. We achieve this using only weak 2D supervision, provided by our novel StanfordExtra dataset. Further, we show we can learn a more detailed shape prior by tuning a gaussian mixture during model training and this leads to improved reconstructions. We also show our method improves over competitive baselines, even when they are given access to ground truth data at test time.

Future work should involve tackling some failure cases of our system, for example handling multiple overlapping dogs or dealing with heavy motion blur. Other areas for research include extending our EM formulation to handle video input to take advantage of multi-view shape constraints, and transferring knowledge accumulated through training on StanfordExtra dogs to other species.

<sup>4</sup> PCK results in tables have been updated to match definitions of Yang and Ramanan [32] normalized by 2D silhouette area. Please see original tables and further details in the appendix.



**Fig. 6. Qualitative results on StanfordExtra and Animal Pose [6].** For each sample we show: (a) input image, (b) predicted 3D mesh, (c) canonical view 3D mesh, (d) reprojected model joints and (e) silhouette reprojection error.

## 6 Acknowledgements

The authors would like to thank the GSK AI team for providing access to their GPU cluster, Michael Sutcliffe, Thomas Roddick, Matthew Allen and Peter Fisher for useful technical discussions, and the GSK TDI group for project sponsorship.

## References

1. Agudo, A., Pijoan, M., Moreno-Noguer, F.: Image collection pop-up: 3d reconstruction and clustering of rigid and non-rigid categories. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (June 2018)
2. Alldieck, T., Magnor, M., Bhatnagar, B.L., Theobalt, C., Pons-Moll, G.: Learning to reconstruct people in clothing from a single rgb camera. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (2019)
3. American Pet Products Association: 2019-2020 APPA National Pet Owners Survey (2020), <http://www.americanpetproducts.org>
4. Andriluka, M., Pishchulin, L., Gehler, P., Schiele, B.: 2d human pose estimation: New benchmark and state of the art analysis. In: IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (June 2014)
5. Biggs, B., Roddick, T., Fitzgibbon, A., Cipolla, R.: Creatures great and SMAL: Recovering the shape and motion of animals from video. In: ACCV (2018)
6. Cao, J., Tang, H., Fang, H., Shen, X., Tai, Y., Lu, C.: Cross-domain adaptation for animal pose estimation. In: 2019 IEEE/CVF International Conference on Computer Vision (ICCV). pp. 9497–9506 (2019)
7. Cashman, T.J., Fitzgibbon, A.W.: What shape are dolphins? Building 3D morphable models from 2D images. *IEEE transactions on pattern analysis and machine intelligence* **35**(1), 232–244 (2013)
8. Chen, L.C., Papandreou, G., Kokkinos, I., Murphy, K., Yuille, A.L.: Deeplab: Semantic image segmentation with deep convolutional nets, atrous convolution, and fully connected crfs. *CoRR abs/1606.00915* (2016)
9. Everingham, M., Van Gool, L., Williams, C.K., Winn, J., Zisserman, A.: The pascal visual object classes (voc) challenge. *International journal of computer vision* **88**(2), 303–338 (2010)
10. He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 770–778 (2016)
11. Ionescu, C., Papava, D., Olaru, V., Sminchisescu, C.: Human3.6m: Large scale datasets and predictive methods for 3d human sensing in natural environments. *IEEE transactions on pattern analysis and machine intelligence* **36**(7) (2013)
12. Johnson, S., Everingham, M.: Clustered pose and nonlinear appearance models for human pose estimation. In: Proceedings of the British Machine Vision Conference (2010), doi:10.5244/C.24.12
13. Kanazawa, A., Black, M.J., Jacobs, D.W., Malik, J.: End-to-end recovery of human shape and pose. In: Computer Vision and Pattern Recognition (CVPR) (2018)
14. Kanazawa, A., Black, M.J., Jacobs, D.W., Malik, J.: End-to-end recovery of human shape and pose. In: Proc. CVPR (2018)
15. Kanazawa, A., Tulsiani, S., Efros, A.A., Malik, J.: Learning category-specific mesh reconstruction from image collections. In: European Conference on Computer Vision. pp. 371–386 (2018)
16. Kato, H., Ushiku, Y., Harada, T.: Neural 3d mesh renderer. In: The IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (2018)
17. Kearney, S., Li, W., Parsons, M., Kim, K.I., Cosker, D.: Rgb-dog: Predicting canine pose from rgb-d sensors. In: IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) (June 2020)
18. Khoreva, A., Benenson, R., Ilg, E., Brox, T., Schiele, B.: Lucid data dreaming for object tracking. The 2017 DAVIS Challenge on Video Object Segmentation - CVPR Workshops (2017)

19. Khosla, A., Jayadevaprakash, N., Yao, B., Fei-Fei, L.: Novel dataset for fine-grained image categorization. In: First Workshop on Fine-Grained Visual Categorization, IEEE Conference on Computer Vision and Pattern Recognition (June 2011)
20. Kolotouros, N., Pavlakos, G., Black, M.J., Daniilidis, K.: Learning to reconstruct 3d human pose and shape via model-fitting in the loop. In: Proceedings of the IEEE International Conference on Computer Vision. pp. 2252–2261 (2019)
21. Kolotouros, N., Pavlakos, G., Daniilidis, K.: Convolutional mesh regression for single-image human shape reconstruction. In: Proc. CVPR (2019)
22. Lin, T.Y., Maire, M., Belongie, S., Hays, J., Perona, P., Ramanan, D., Dollar, P., Zitnick, L.: Microsoft COCO: Common objects in context. In: ECCV. European Conference on Computer Vision (September 2014)
23. Loper, M., Mahmood, N., Romero, J., Pons-Moll, G., Black, M.J.: SMPL: A skinned multi-person linear model. *ACM transactions on graphics (TOG)* **34**(6), 248 (2015)
24. von Marcard, T., Henschel, R., Black, M., Rosenhahn, B., Pons-Moll, G.: Recovering accurate 3d human pose in the wild using imus and a moving camera. In: European Conference on Computer Vision (ECCV) (sep 2018)
25. Newell, A., Yang, K., Deng, J.: Stacked hourglass networks for human pose estimation. In: European Conference on Computer Vision. pp. 483–499. Springer (2016)
26. Novotny, D., Ravi, N., Graham, B., Neverova, N., Vedaldi, A.: C3DPO: Canonical 3d pose networks for non-rigid structure from motion. In: Proc. ICCV (2019)
27. Pavlakos, G., Zhu, L., Zhou, X., Daniilidis, K.: Learning to estimate 3D human pose and shape from a single color image. In: Proc. CVPR (2018)
28. Probst, T., Pani Paudel, D., Chhatkuli, A., Van Gool, L.: Incremental non-rigid structure-from-motion with unknown focal length. In: The European Conference on Computer Vision (ECCV) (2018)
29. Reinert B, Ritschel T, S.H.P.: Animated 3d creatures from single-view video by skeletal sketching. In: Proc. Graphics Interface (2016)
30. Saito, S., , Huang, Z., Natsume, R., Morishima, S., Kanazawa, A., Li, H.: Pifu: Pixel-aligned implicit function for high-resolution clothed human digitization. arXiv preprint arXiv:1905.05172 (2019)
31. Vicente, S., Agapito, L.: Balloon shapes: Reconstructing and deforming objects with volume from images. In: 2013 International Conference on 3D Vision - 3DV 2013. pp. 223–230 (2013)
32. Yang, Y., Ramanan, D.: Articulated human detection with flexible mixtures of parts. *IEEE transactions on pattern analysis and machine intelligence* **35**(12), 2878–2890 (2013)
33. Zuffi, S., Kanazawa, A., Berger-Wolf, T., Black, M.J.: Three-d safari: Learning to estimate zebra pose, shape, and texture from images ”in the wild”. In: The IEEE International Conferene on Computer Vision (ICCV) (2019)
34. Zuffi, S., Kanazawa, A., Black, M.J.: Lions and tigers and bears: Capturing non-rigid, 3D, articulated shape from images. In: IEEE Conference on Computer Vision and Pattern Recognition (CVPR). IEEE Computer Society (2018)
35. Zuffi, S., Kanazawa, A., Jacobs, D., Black, M.J.: 3D menagerie: Modeling the 3D shape and pose of animals. In: IEEE Conf. on Computer Vision and Pattern Recognition (CVPR) (Jul 2017)



# Who Left the Dogs Out?

## Supplementary Material

### A Dataset curation

In this section, we describe our process for obtaining keypoint and segmentation annotations for the Stanford Dog Dataset [19]. We submit the entire set of 20,580 dog images to the Amazon Mechanical Turk crowdsourcing platform to obtain a set of 20 keypoint and segmentation masks. We overlay 1 bounding box, provided with the original dataset, on the submitted images to identify the specific dog for the annotators to label. Each image was sent to 3 independent annotators for collecting keypoints and segmentation masks.

**Keypoints.** To identify keypoints, workers were given a list of 20 keypoints to click: 2 per tail, 3 per leg, 2 per ear, nose and jaw. They were additionally asked to provide a visibility flag per point.

For each keypoint, we process the three clicks to yield a reliable coordinate. From the 3 clicks, we discard clicks that are further than a set tolerance from the mean. If at least 2 clicks remain, we take the mean coordinate as the accepted keypoint position. Otherwise, the point has not been reliably identified between workers, so we set the keypoint as invisible. As described in the main paper, we remove images from train and test splits which have fewer than 8 visible keypoints.

**Segmentation.** For each image, each worker  $w \in \{w_1, w_2, w_3\}$  submits a binary segmentation mask  $\mathbf{A}^w \in \mathbb{R}^{H \times W}$ . We request a re-labelling for any submissions which fail simple criteria, such as if the highlighted area is below a threshold number of pixels.

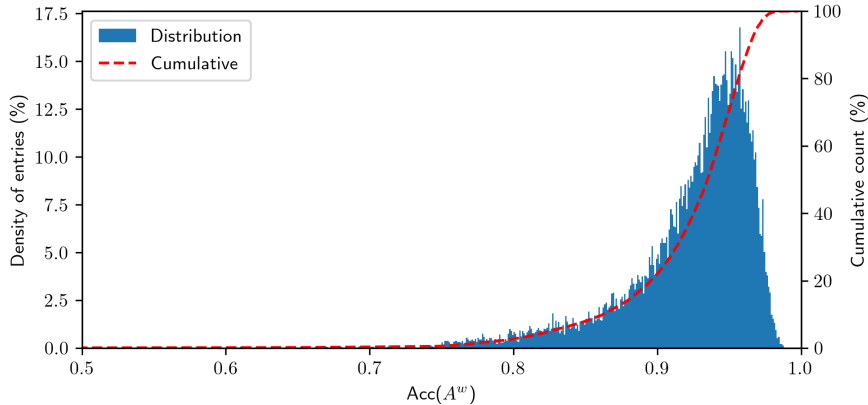
For each image, we generate the most likely segmentation by comparing submissions across workers. For any two workers  $w, w'$  we compute a correlation coefficient:

$$c_{w,w'} = \frac{\sum_i \sum_j [\mathbf{A}^w \odot \mathbf{A}^{w'}]_{i,j}}{\max_{p=\{w,w'\}} \sum_i \sum_j \mathbf{A}^p_{i,j}} \quad (1)$$

Where  $\odot$  denotes the element-wise product of the matrices. We remove a worker’s segmentation  $A^w$  if all correlation coefficients  $c_{w,w'}$  are below a set threshold. The final binary mask is computed from the remaining submissions:

$$\hat{A}_{i,j} = \begin{cases} 1, & \text{if } \sum_w A^w_{i,j} > 1 \\ 0, & \text{otherwise} \end{cases} \quad (2)$$

We can also define the accuracy of a worker’s segmentation, as the largest of their correlation coefficients:  $\text{Acc}(A^w) = \max_{w' \neq w} \{c_{w,w'}\}$ . Figure I shows the set of segmentation annotation accuracies over the entire labelled dataset.



**Fig. I.** Accuracy distribution of all submitted dog segmentations across the entire Stanford Dog Dataset.

## B Fitting SMBLD to 3D animation data

Another method for improving the generalizability of the SMAL model is to improve the 3D shape prior. Such priors are typically used to ensure shape deformation remain within a realistic and anatomically plausible range. Due to the limited diversity of scans used to build the SMAL model, while the shape prior does enforce realism among deformations, it does not allow for a wide enough range to cover the set of dogs in our dataset.

We improve the quality of the prior (and learn a prior over our new scale parameters) by fitting to a set of 13 artist-designed 3D dog meshes, designed for animation use, which are more varied than the original set. We apply an energy minimization scheme which aligns the SMAL vertices to each scan, under smoothing regularizers.

Recall that SMBLD is adapted from the SMAL [35] deformable animal mesh, by including limb scaling parameters. We learn a prior by fitting our SMBLD model, which comprises parameters for pose  $\theta$  and shape  $\beta$  (the latter of which includes our scaling parameters  $\kappa$ ).

Note that fitting SMBLD to 3D scans is significantly easier than to 2D images, since the complete 3D information of the target mesh is available. In addition, our target meshes are not particularly detailed and are already aligned in T-pose, so we avoid need for a complex alignment technique as discussed in, for example SMPL [23] or SMAL [35].

We run an energy minimization process to align the SMBLD mesh to the 3D scans, subject to some smoothing regularizers. We minimize the following energy formulation:

$$E_{\text{opt}} = E_{\text{chamfer}} + E_{\text{laplacian}} + E_{\text{edge}} + E_{\text{normal}} \quad (3)$$

where each of these terms has a scalar weight  $\lambda$ . We set  $\lambda_{\text{chamfer}} = \lambda_{\text{edge}} = 1.0$ ,  $\lambda_{\text{normal}} = 0.01$  and  $\lambda_{\text{laplacian}} = 0.1$ . We run the optimization using SGD, learning rate  $1e - 4$  for 1000 iterations.

**Chamfer energy.** A measure of the average distance between vertices of the SMBLD mesh  $V = F_v(\theta, \beta)$ , and the target mesh vertices  $V'$ , when  $p$  vertices  $v_i, v'_j$  are sampled from each mesh respectively:

$$E_{\text{chamfer}}(V, V') = \frac{1}{p} \sum_{i=1}^p \min_j^p |v_i - v'_j| \quad (4)$$

**Uniform laplacian energy.** A measure of the mesh smoothness.

**Edge energy.** This energy is equal to the average edge length across the mesh, and is used to encourage uniform distribution of vertices.

**Normal energy.** This energy promotes consistency between adjacent faces. It is a measure of the average normal consistency between adjacent faces. For two faces with normals  $\mathbf{n}_0$  and  $\mathbf{n}_1$ , the normal consistency is  $1 - \frac{\mathbf{n}_0 \cdot \mathbf{n}_1}{|\mathbf{n}_0| |\mathbf{n}_1|}$ .

At the end of this process, we have a collection of fits  $(\theta, \beta)_{\{i=1, \dots, 13\}}$  from which we can learn our unimodal pose and shape priors. As discussed, we eventually use this unimodal shape prior to initialize our mixture shape prior, which is tuned with the expectation-maximization step in the training loop.

## C Training procedure

Recall that the training objective for our end-to-end system for predicting SMBLD parameters consistent with a monocular dog input image is given by:

$$L_{\text{opt}} = L_{\text{joints}} + L_{\text{sil}} + L_{\text{pose}} + L_{\text{shape}} + L_{\text{mixture}} \quad (5)$$

As described in the paper, each loss term is weighted with a scalar  $\lambda$  and we train our method in two stages:

**Stage 1.** We set  $\lambda_{\text{joints}} = 10.0$ ,  $\lambda_{\text{pose}} = 1.0$ ,  $\lambda_{\text{shape}} = 1.0$ ,  $\lambda_{\text{sil}} = 0.0$ ,  $\lambda_{\text{mixture}} = 0.0$ . We train this stage for 250 epochs, using the Adam optimizer, with learning rate set to  $10^{-4}$ .

**Stage 2.** In this stage, we introduce the silhouette loss to encourage a shape alignment between the projected model silhouette and the ground truth annotation. We set  $\lambda_{\text{joints}} = 10.0$ ,  $\lambda_{\text{pose}} = 0.5$ ,  $\lambda_{\text{shape}} = 0.0$ ,  $\lambda_{\text{sil}} = 100.0$ ,  $\lambda_{\text{mixture}} = 0.1$ . We train this stage for 150 epochs and run the described EM update step every  $K = 15$  epochs. We selected to use  $M = 10$  clusters based on a grid search over  $M = 1, 5, 10, 25$  and comparing IoU. We again use the Adam optimizer, and set the learning rate to  $10^{-5}$ .

## D Probability of Keypoints Max (PCK-MAX)

In this section, we compare reprojected 2D joint accuracy using the *PCK-MAX* evaluation metric. This protocol is similar to the Percentage of Correct Keypoints (PCK) metric [32] used in the main paper by incorporating ‘invisible’ ground-truth points. The standard PCK metric ignores these points, meaning even correct 3D reconstructions will receive no credit. PCK-MAX instead assumes reconstructed 3D points for missing ground-truth data are correct, providing an interesting upper bound. Results are shown in Table 1 and Table 2.

Method	Kps	Seg	PCK-MAX @ 0.1				
			Avg	Legs	Tail	Ears	Face
3D-M [35]	Pred	Pred	67.1	65.7	79.5	54.9	87.4
3D-M	GT	GT	72.6	69.9	92.0	58.6	96.9
3D-M	GT	Pred	72.6	70.2	91.5	58.1	96.9
3D-M	Pred	GT	67.4	66.0	79.9	55.0	88.2
CGAS [5]	CGAS	Pred	43.7	46.5	64.1	36.5	21.4
CGAS	CGAS	GT	43.6	46.3	64.2	36.3	21.6
3D-M + scaling	Pred	Pred	69.6	69.4	79.3	56.5	87.6
3D-M + scaling + EM prior	Pred	Pred	71.6	71.5	<b>80.7</b>	59.3	88.0
<b>Ours</b>	—	—	<b>75.7</b>	<b>75.0</b>	77.6	<b>69.9</b>	<b>90.0</b>

**Table 1. PCK-MAX baselines.** PCK-MAX scores are shown for SOTA methods under varying conditions. Directly comparable baseline methods (requiring only an input image) are highlighted. *Pred* keypoints generated with Hourglass-Net [25] and segmentations with DeepLab v3+ [8]. 3D-M/CGAS are also analysed when they have access to ground-truth keypoints and/or segmentation masks. We also analyse adding this paper’s innovations (scale parameters and EM prior) to the 3D-M method [35].

Method	PCK-MAX @ 0.1				
	Avg	Legs	Tail	Ears	Face
3D-M [35]	69.1	60.9	83.5	75.0	93.0
<b>Ours</b>	<b>73.8</b>	<b>65.1</b>	<b>85.6</b>	<b>84.0</b>	<b>93.6</b>

**Table 2. PCK-MAX Animal Pose dataset [6].** Evaluation on recent Animal Pose dataset with no fine-tuning to our method nor joint/silhouette predictors used for 3D-M.

Method	PCK-MAX @ 0.1				
	Avg	Legs	Tail	Ears	Face
<b>Ours</b>	<b>75.7</b>	<b>75.0</b>	<b>77.6</b>	69.9	90.0
–EM	74.6	72.9	75.2	<b>72.5</b>	88.3
–MoG	74.9	74.3	73.3	70.0	<b>90.2</b>
–Scale	72.6	72.9	75.3	62.3	89.1

**Table 3. PCK-MAX ablation study.** Evaluation with the following contributions removed: (a) EM updates, (b) Mixture Shape Prior, (c) SMBLD scale parameters.