PartImageNet: A Large, High-Quality Dataset of Parts

Ju He¹, Shuo Yang², Shaokang Yang³, Adam Kortylewski^{1,4,5}, Xiaoding Yuan¹, Jie-Neng Chen¹, Shuai Liu³, Cheng Yang³, Qihang Yu¹, and Alan Yuille¹

¹Johns Hopkins University ²University of Technology Sydney ³ByteDance Inc. ⁴Max Planck Instutite for Informatics ⁵University of Freiburg

Abstract. It is natural to represent objects in terms of their parts. This has the potential to improve the performance of algorithms for object recognition and segmentation but can also help for downstream tasks like activity recognition. Research on part-based models, however, is hindered by the lack of datasets with per-pixel part annotations. This is partly due to the difficulty and high cost of annotating object parts so it has rarely been done except for humans (where there exists a big literature on part-based models). To help address this problem, we propose PartImageNet, a large, high-quality dataset with part segmentation annotations. It consists of 158 classes from ImageNet with approximately 24,000 images. PartImageNet is unique because it offers part-level annotations on a general set of classes including non-rigid, articulated objects, while having an order of magnitude larger size compared to existing part datasets (excluding datasets of humans). It can be utilized for many vision tasks including Object Segmentation, Semantic Part Segmentation, Few-shot Learning and Part Discovery. We conduct comprehensive experiments which study these tasks and set up a set of baselines.

1 Introduction

When humans observe objects we can effortlessly parse them into their component parts. Studies in cognitive psychology show that humans learn rich hierarchical representations of objects [20] and can decompose objects into parts taking into account their spatial relationships [1]. Partly inspired by these findings, computer vision researchers have studied how to model parts and to represent objects in terms of parts. The big literature on these topics includes deformable templates [45], pictorial structures [11], constellation models [38, 10] and grammar-based models [50, 15]. In particular, there have been work [40, 37] on segmenting object parts and studying whether this helps improve the segmentation of objects, and later studies of the use of these part segmentation for downstream tasks, e.g., describing pedestrians and facilitating person retrieval [32]. Most recently, part representations have been proposed to improve panoptic segmentation [14]. It has also been argued that object-part models will play an important role in few-shot learning [18, 39, 43] and hence help to alleviate the the computer vision communities dependence on annotated large-scale datasets. In



Fig. 1. Example figures of annotated images in PartImageNet. We offer high-quality precise and dense part segmentation annotation on a wide range of general species including both non-rigid and rigid objects. In total there are around 24,000 images of 158 classes from ImageNet annotated.

short, object-part models are promising for improving computer vision on many different tasks.

But in the big data era, research on part-based models and their applications is hindered due to the shortage of datasets with per-pixel part annotations. Current datasets are almost always restricted to humans or to a small number of object categories, e.g., PASCAL-Part [5]. There is a need to extend this datasets to include many more object categories, as suggested by [14]. Existing part datasets almost always focus on humans or a few rigid classes such as cars. This is partly caused by the difficulty of per-pixel part annotations compared to other types of annotation like bounding boxes. It requires much more effort to ensure accuracy and quality consistency, especially for non-rigid objects. The few works [32, 43, 17] that attempt to use parts as a mid-level representation, often learnt without supervision, also suffer from this lack of annotated datasets which makes it hard to evaluate whether they have actually captured meaningful object parts.

This motivates us to introduce PartImageNet — a large, high-quality dataset with part annotation on a general set of object classes. Concretely speaking, we manually select 158 classes from ImageNet [9] and group them into 11 supercategory following the WordNet hierarchy of ImageNet. Part labels are designed according to the super-category while can be elaborated to fine-grained classes in a hierarchical way. A carefully designed pipeline is taken to ensure the high quality of our PartImageNet annotations. As far as we know, this is the only dataset after PASCAL-Part [5] that offers part-level annotations on more general classes instead of just humans and rigid objects. Compared to PASCAL-Part [5], we annotate much more images (24k v.s. 10k) on much more classes (158 v.s. 20). Extensive experiments on PartImageNet are conducted to show that parts could help general object segmentation and few-shot learning and set up a set of baselines of different downstream tasks on this benchmark. We believe that with this dataset, the research on part-based models and their applications will be facilitated a lot. In summary, we make the following contributions in this work:

- 1. We briefly review the history of part-based models and introduce their potential applications in downstream tasks.
- 2. We introduce PartImageNet a large, high-quality dataset with part annotations on a general set of classes. From our perspective, part-level annotation, especially on non-rigid objects, is very rare and valuable.
- 3. We set up a set of baselines on PartImageNet in different vision tasks with state-of-the-art methods which shows the broad usage of the dataset.
- 4. We conduct experiments to show that introducing parts annotations is beneficial to the object segmentation and few-shot learning which points out to a promising future direction.

2 Related Work

2.1 Part-based Models

Modeling objects in terms of parts is a long-standing problem in computer vision and there is rich history of research on this topic. Starting from Pictorial Structures in the early 1970's [13], plenty of different models [45, 11, 38, 10, 50, 15] have been proposed to explicitly model parts and their spatial relations to the whole object. There have been a variety of technical approaches but a common theme is that object-part models provide rich representations of objects and help interpretablity. In recent years, partly due to the availability of big data, research on part-based models also includes part segmentation and unsupervised part exploitation.

Supervised Part Segmentation Annotated human parts datasets have long existed and there has been much work [6, 40] on human part detection and semantic segmentation of human parts. Sun et al. [32] proposed to exploit part-level features for pedestrian image description and thus facilitate person retrieval. There has also been some work [40, 37] on semantic segmentation of parts of a limited class of other objects.

Unsupervised Part Exploitation on general tasks Due to the lack of annotated data on a more general set of classes, research on parts for non-human object classes is often unsupervised. Thewlis et al. [33] proposed to enforce the equivariance of landmark locations under artificial transformations of images to generate semantic meaningful parts of objects. Lorenz et al. [25] improved part discovery by simultaneously exploiting the invariance and equivariance constraints between synthetically transformed images and disentangling the shape and appearance of objects. Recently, parts have shown to be beneficial to unsupervised, or few-shot, object learning since the modeling of parts helps alleviates the scarcity of training data provided the parts can be shared between different classes. He et al. [18] exploited the fact that the feature vectors in the CNN can be viewed approximately as part detectors and their geometry relations could be

estimated by clustering. Additionally, it has been shown that few-shot segmentation benefits from the modeling of parts. Liu et al. [24] decomposes the holistic class representation into a set of part-aware prototypes, capable of capturing diverse and fine-grained object features, which benefits semantic segmentation.

However, the part modeling of these methods mainly relies on unsupervised clustering combined with attention mechanisms. Without strong supervision, the results produced by these methods are not very satisfactory, because they only generate meaningful part segmentation in a few simple scenarios. It remains unclear if these approaches can really lead to the learning of semantically meaningful parts without evaluation. Thus it is important to introduce dataset with part annotation to analyze the actual effectiveness of such part modeling and promote the further research on this promising direction.

2.2 Part Datasets

2D Part Datasets There exists multiple ways to annotate object parts in images. Among them, bounding boxes and keypoints are relatively easy to annotate while per-pixel segmentation is much harder due to the extreme fine-grained difficulty and high cost. The type of objects (i.e. rigid v.s. non-rigid) also plays an important role in deciding the difficulty of annotation. Wang et al. [36] provides dense bounding box annotation for parts on 6 vehicle categories. PASCAL3D+[41], CarFusion [27] and Apollocar3D [31] offer keypoint annotation on a few rigid object categories especially on cars. ADE20K [49] contains part segmentation annotation on many rigid object categories. As for non-rigid objects, CUB-200-2011 [35] provides keypoint location for birds parts. LIP [21], MHP [48], CIHP [16] include instance-aware, part-level annotations for human.

As far as we are concerned, PASCAL-Part [5] is the only existing dataset that offers part-level annotation on a more general set of categories in a per-pixel segmentation format. However, it contains relatively small number of images and only a small set of classes. We introduce a much larger, high-quality dataset of parts to support more research on part-based models.

3D Part Dataset 3D Part datasets also play a crucial role in the advances of 3D shape understanding tasks. Yi et al. [44] first takes an active learning approach to annotate the 3D models selected from 16 categories in ShapeNet [2]. Part-Net [26] further provides hierarchical part annotations on 3D models covering 24 object categories, most of which are indoor furniture. Recently, CGPart [23] proposes to use computer graphics model to efficiently generate a large-scale vehicle dataset which offers part annotation. Although these datasets have shown their effectiveness in helping data-driven models, they still suffer from the problem of a lacking annotation of non-rigid object categories. The different domains also limit their usage in the image part-level parsing.



Fig. 2. Overview of the PartImageNet dataset. Accurate part segmentation masks on both non-rigid and rigid objects are offered. The part labels are defined on the supercategory (e.g. Quadruped) level but can be easily transferred to mid-level or class-level part labels according to the need (e.g. Quadruped Head \rightarrow Dog Head \rightarrow Gordon setter Head) following the WordNet hierarchy as shown by the red dotted arrow.

3 PartImageNet Dataset

In this section, we present the details of how we collect and annotate the data, followed by statistics and analyze on the quality of the PartImageNet dataset. The overview of the PartImageNet is shown in Figure 2.

3.1 Data Collection

Data Source As suggested by the name, we collect images from the ILSVRC-12 dataset [9]. There are in total 158 object classes selected to be annotated in our dataset. All the images conform to licensing for research purposes.

Object Categories Analogous to tieredImageNet [28], which is also a subset of ILSVRC-12 [9], we group classes into super-categories corresponding to higher-level nodes in the ImageNet hierarchy. There are 11 super-categories in total. To make the dataset more challenging and more valuable, we pick up fewer rigid objects such as vehicle but choose more animals to annotate. Thus for super-categories like Quadruped, they contain around 40 classes while for super-categories such as vehicle, boat, they only have less than 10 classes. In total, there exists 118 classes out of 158 which are non-rigid objects.

Filtering Unsuitable Images As our PartImageNet dataset focuses on part segmentation, we eliminate those images which have no proper parts to annotate due to possible occlusion or improper viewpoints. Besides, to simplify the annotation process and avoid ambiguity during annotation, we discard all the images that contain more than one desired object with qualified parts to be annotated. (i.e. all annotated parts in an image are guaranteed to belong to one object).

3.2 Annotation

Instead of annotating bounding boxes or keypoints which is much easier, we aim at annotating high-quality part segmentation mask at pixel-level. The part labels for classes are determined by the corresponding super-categories (i.e. classes under the super-category share the same part labels). It is difficult to ensure that all the classes under the same semantic category have the same detailed parts so we only annotate those that are most important. Take horns for example, some mammals have horns while others do not, to simplify the definition and the annotation process, we do not create a horn label for mammals. Instead, only head label exists and for mammals with horns, the horns are also counted as part of the head during annotation. The detailed part labels for different categories are shown in Table 1. Notice that though the part labels are defined at the super-category level, they can easily be transferred to mid-level or class-level part labels according to the need following the WordNet hierarchy of the ImageNet as shown in Figure 2.

Annotation Pipeline Due to the extreme difficulty of annotating part segmentation mask at pixel level, a good annotation pipeline is of great importance to ensure the high quality and maintain the consistency of the annotation. Motivated by existing works on datasets such as ImageNet [9], COCO [22], Objects365 [29], we divide our annotation pipeline into three steps. First as we sample our images from ImageNet, we already have the class and super-category label information for the images, thus we split the annotation task into 11 (equals to the number of super-categories) sub-tasks to alleviate the workload of annotators. The second step is to choose whether to keep the image for the PartImageNet or not following the requirements in Sec 3.1. In the last step, the annotators are going to annotate the segmentation mask of specific parts for the corresponding super-category. The specific part labels to be annotated are automatically generated according to the super-category label, thus the annotators do not need to pick the right parts among all possible parts of other super-categories which significantly improves the overall accuracy and efficiency.

Annotation Team We divide our annotation team into three groups: annotators, inspectors, examiners. All the images are first annotated by annotators, then a subset of annotated images will be randomly chosen to be checked by the inspectors. In the end, the examiners can check as many images as they want and see if the quality meets the requirement. Any failures in the above steps will result in re-annotation of the job.

Annotator The annotators' job is to annotate all the images following the pipeline introduced in Sec 3.2. Before starting the official annotation, all annotators are going to take courses from the inspectors and go for a test annotation round. During the whole annotation process, the annotators can contact the inspectors at any time if they have questions regarding the current task.

Inspector The duty of the inspectors is to ensure the quality of the annotation made by the annotators. The inspectors will first have a meeting with the annotators before annotation to teach them the annotation rules. Then they will review all the annotated images during the test annotation round and provide feedback to the annotators. For the official annotation, a random subset of annotated images will also be reviewed by them. If there is an obvious error or the annotation fails to meet the quality requirement, the task will be rejected and re-annotated. If the rejection ratio of an annotator is too high, then all his annotated tasks will be discarded and assigned to other annotators.

Examiner Examiners design the annotation rules and discuss with the inspectors to make the quality requirement. They are also responsible for answering all the ambiguity questions. After all annotations are done, they review most of the annotated images to ensure the quality and unqualified ones will be rejected and re-annotated.

3.3 Annotation Quality and Consistency

Part segmentation annotation is very hard due to the variance of objects pose, orientation, occlusion. Besides, another important problem that does not exist in other annotation tasks is the ambiguity of how to define the separation of different parts (i.e. how do we define the boundary between head and body). It is impossible to make clear annotation rules that can keep all annotators consistent on the boundary annotation and handle so many different variations. Thus, to make the annotation of parts as accurate and consistent as we can, we make the following efforts:

Maximize possible information The annotators are asked to annotate all visible and distinguishable parts no matter how small they are in order to keep as much information as we can. We believe that such kind of small parts also need to be studied and should be handled during the algorithm design stage instead of the annotation stage.

Accurate boundary segmentation To keep the annotations accurate so they do not contain too many background pixels. We set strict annotation rules to guide the annotators to create tight segmentation masks in most situations (the requirements are relaxed under fuzzy situations). The occluders are also required to be masked out of the annotation.

Consistent Annotation Task Assignment As our part labels are set according to the super-categories of the images. To keep the consistency among images of the same super-category, we divide the annotators into sub-groups where each sub-group is responsible for annotating one super-category and can communicate freely within the group. In this way, we alleviate the inconsistency of boundary separation especially among the same super-category.

3.4 Statistics

Category and Class Mapping As introduced in Sec 3.1, our PartImageNet focuses more on the challenging animals categories instead of rigid object cat-

Table 1. Number of classes and annotated parts for each category in PartImageNet. The number in the brackets after the category name indicates the total number of classes under the category.

Category	Annotated Parts
Quadruped (46)	Head, Body, Foot, Tail
Biped (17)	Head, Body, Hand, Foot, Tail
Fish (10)	Head, Body, Fin, Tail
Bird (14)	Head, Body, Wing, Foot, Tail
Snake (15)	Head, Body
Reptile (20)	Head, Body, Foot, Tail
Car(23)	Body, Tire, Side Mirror
Bicycle (6)	Head, Body, Seat, Tire
Boat (4)	Body, Sail
Aeroplane (2)	Head, Body, Wing, Engine, Tail
Bottle (5)	Body, Mouth

egories. Thus we pick far more classes in animals than cars, boats, planes etc. The detailed number of classes per category is shown in Tab 1. In summary, we annotated 158 classes and 118 of them are non-rigid objects. By contrast, PASCAL-Part [5] annotated 20 classes and 12 of them are non-rigid objects.

Total Annotated Images Based on our proposed annotation pipeline, around 24,000 images are annotated in the PartImageNet dataset. We randomly sampled 85%, 5% and 10% images per class into training, validation and testing set. The detailed numbers of images and annotated parts in each set is shown in Table 2.

 Table 2. The annotation split of the Table 3. The annotation density of the PartImageNet dataset.

 PartImageNet dataset.

0			0	
	Images	Parts	Number of parts	Proportion (%)
Train	20481	95059	1-2	22.00
Validation	1206	5626	3-6	57.65
Test	2408	11275	7-9	18.61
All	24095	111960	10 +	1.74

Annotation Density Though we directly offer more high-level part annotation compared to PASCAL-Part [5], PartImageNet still has a high density of part annotations with 111960 part instances in total. We compute the proportion of number of annotations per image as shown in Table 3. Most images contain around 3-6 annotations which is quite dense and considered our overall dataset size, it should provide enough training examples for most algorithms.

Class Distribution To dive into the details of the PartImageNet, we provide the overall class distribution in Fig 3. Though the original number of images per class in ImageNet [9] is roughly the same. After our annotation process,



Fig. 3. Number of images per class in PartImageNet and PASCAL-Part [5]. Class index is sorted according to the number of images in the class. Our PartImageNet exhibits a more balanced distribution with a few tail classes that contain few images. By contrast, PASCAL-Part exhibits a more sharpen drop tendency.

some classes will have more images be ruled out due to multiple objects in one image or occlusion (e.g. bottle). Thus PartImageNet naturally has a few tail classes that contain few images. However, compared to the PASCAL-Part [5] dataset, PartImageNet exhibits a much more balanced distribution. The number of images of the most majority class in PASCAL-Part is roughly 9 times more than that of the middle class, while that ratio in PartImageNet is only 1.25.

4 Experiments

In this section, we conduct extensive experiments on our proposed PartImageNet for different tasks including semantic part segmentation, object segmentation and few-shot learning. We broadly evaluate classic methods along with some of the state-of-the-art models to set up a set of baselines on this benchmark. While all these methods do not take part annotations into account, we further show that by exploiting annotated parts, the performance on object segmentation and few-shot learning could get a non-trivial improvement.

4.1 Semantic Part Segmentation

We conduct experiments on semantic part segmentation PartImageNet using Semantic FPN [19], Deeplabv3+ [4] and SegFormer[42]. Semantic FPN [19]

and Deeplabv3+ [4] are classic convolution-based semantic segmentation methods. SegFormer[42] is one of the state-of-the-art transformer-based semantic segmentation frameworks. We use Resnet-50 as backbone for Semantic FPN[19] and Deeplabv3[3]. MiT-b2 (Mix Transformer encoders) is adopted for Segformer [42].

Table 4 summaries the results for semantic part segmentation and Figure 4 shows the qualitatively visualizations. As can be observed, methods with strong supervision can generally produce satisfactory results on the semantic part segmentation especially when the background and the shape are relatively easy. However, they still suffer from three main challenges: 1) inaccurate boundary between semantic parts, 2) wrong label assignments on similar semantic parts, 3) ignoring small semantic parts (e.g. See row3 - row5 of Fig 4). Besides, we also notice that the recent progress in methods for general object segmentation does not seem to bring expected improvement in the context of semantic part segmentation (e.g. SegFormer only bring limited improvement based on Deeplabv3+). This reveals the difficulty of part segmentation and indicates that special architecture or module design is needed for effectively solving the challenge.

Table 4. Experimental results of Part Segmentation on PartImageNet. Performance are evaluated in terms of mIoU and scores on validation and testing set are reported.

Model	Backbone	Crop Size	Val mIoU	Test mIoU
Semantic FPN [19]	$\operatorname{ResNet-50}$	512x512	60.36	58.69
Deeplab v3 $+$ [4]	$\operatorname{ResNet-50}$	512x512	64.20	64.83
SegFormer [42]	MiT-B2	512x512	65.27	65.17

4.2 Object Segmentation

While we offer part-level segmentation masks, they could easily be integrated to offer a full object segmentation mask and thus serves as a benchmark for Salient Object Segmentation. As this role, PartImageNet is a relatively easy one compared to popular MS-COCO [22] and CityScapes [8] since it is a salient holistic segmentation of a single object. However, it is unique as it offers the opportunity to conduct research on the relation between parts and whole objects as shown below. Besides, it also has the advantage that it could be evaluated at a hierarchy level (i.e. super-category & fine-grained class).

Baseline We still adopt Semantic FPN [19], Deeplabv3+ [4] and SegFormer [42] for object segmentation. The experimental setup and training pipeline are in line with those in semantic part segmentation.

Table 5 summarizes the results of Object Segmentation on PartImageNet. Here the mIoU on the fine-grained classes is reported as it is more challenging than segmenting the super-category object labels. As can be observed, existing methods already achieve quite good results on the benchmark as expected. We are more interested in whether parts can help the general object segmentation.



Fig. 4. Example figures of semantic part segmentation results. The quality of the results highly depend on the background, shape of the objects, occlusion situation and the class itself. More visualizations are shown in supplementary materials.

Can Parts help Object Segmentation? Motivated by the intuition that it would be a natural way for models to first learn to group similar pixels together at the early stage followed by gradually forming a whole object mask, we design experiments to validate whether object segmentation can be improved by introducing part annotations as deep supervision. We take Deeplabv3+ [4] here for concrete analysis and ablation study.

Table 6 summaries the experimental results of adding part annotations as deep supervision in Deeplabv3+ [4] at different stages of the backbone. We observe that adding part annotations at early stage such as stage 2 does not have an obvious effect on the results probably because the features are still too local without the ability to capture semantic meanings. While starting from stage 3, the deep supervision gradually increases the performance of object segmentation. When adding it at the stage 4, the model achieves largest improvement by 1.63% and 1.05% mIoU on the validation and testing set respectively.

Table 7 further conducts ablation study on the type of deep supervision when adding at the 4th stage. We first show that by adding object masks as

Table 5. Experimental results of Object Segmentation on PartImageNet. Performance are evaluated in terms of mIoU and scores on validation and testing set are reported.

Model	Backbone	Crop Size	Val mIoU	Test mIoU
Semantic FPN [19]	$\operatorname{ResNet-50}$	512x512	63.09	60.50
Deeplab v3+ $[4]$	$\operatorname{ResNet-50}$	512x512	74.04	71.79
SegFormer [42]	MiT-B2	512 x 512	81.22	78.56

deep supervision, the performance decreases a little which reveals that by simply introducing object segmentation mask as deep supervision, the performance can not be improved. Then we show that when adding part annotations, it should be supervised with binary cross entropy loss (i.e. only the current part class is considered) in the sense that we should encourage pixels that belong to the same part to be similar while avoiding penalizing wrong classification of the part classes. The reasons for that are two-folded: firstly semantic part segmentation is a harder task compared to object segmentation, we should not do the full part of it at the early stage of the network, secondly some pixels of different parts are very similar at the pixel level which makes them hard to be distinguished when semantic meanings have not be well-learned at the shallow stage.

Table 6. Experimental results of exploit- **Table 7.** Exploiting different kinds of ing part segmentation as deep supervision part segmentation as deep supervision at at different stages of the backbone. DS Stage 4. CE stands for training with Cross stands for Deep Supervision. mIoU on val- Entropy loss and BCE stands for training idation and testing sets are reported. with Binary Cross Entropy loss.

teren and testing sets are reperted.			 en Binarj e	TODD LINE	99 10000
DS Stage	Val mIoU	Test mIoU	DS Type	Val mIoU	Test $mIoU$
None	74.04	71.79	None	74.04	71.79
Stage 2	74.21	71.75	Object	72.45	71.02
Stage 3	74.47	72.26	Part CE	73.88	71.61
Stage 4	75.67	72.84	Part BCE	75.67	72.84

4.3 Few-shot Learning

We could also organize our PartImageNet in an another way by splitting non-overlapping classes into training, validation and testing set, thus it naturally becomes a few-shot learning and transfer learning benchmark. The new split especially designed for Few-shot Learning contains 109 classes in training set, 19 classes in validation set and 30 classes in testing set. Unlike tieredImageNet [28], we do not try to avoid semantic overlap between training and testing sets. The details of the split will be presented in the supplementary materials. By converting PartImageNet into a few-shot benchmark, it offers the community a chance to validate and research on the effects of parts in this domain.

Baseline We follow the conventional setting in few-shot classification to resize images to 84 * 84 pixels and adopt Conv4 and ResNet-12 as backbones with

respect to different methods. We conduct experiments on PartImageNet using MAML [12], Prototypical Networks [30], RFS [34], Meta-Baseline [7], COMPAS [18] and DeepEMD [47]. Among them, MAML [12] and Prototypical Networks [30] are classic few-shot classification methods, RFS [34] and Meta-Baseline [7] are representative works for large-training-corpus methods and meta-training methods respectively while COMPAS [18] and DeepEMD [47] are recent works based on exploitation of parts or key regions to facilitate few-shot learning.

Table 8 summaries the results of selected methods on PartImageNet and miniImageNet. As can be observed, recent methods can obtain similar performances on PartImageNet as miniImageNet which shows that PartImageNet itself can serves as a good benchmark for evaluating few-shot algorithms. Besides, we also observe that though PartImageNet theoretically contains more classes with part structures, models [18, 47] that claim to exploit part information do not show obvious advantages compared to others when training without directly using the part annotations.

Table 8. Experimental results of Few-shot Learning on PartImageNet and miniImageNet. Average classification accuracies(%) with 95% confidence intervals are reported.

Model	Backbone	PartImageNet 5-way miniImageNet 5-way				
model		1-shot	5-shot	1-shot	5-shot	
MAML $[12]$	Conv4	46.9 ± 1.4	58.1 ± 0.7	48.7 ± 1.8	63.1 ± 0.9	
Prototypical Networks [30]	Conv4	50.0 ± 0.6	65.4 ± 0.5	49.4 ± 0.8	68.2 ± 0.7	
RFS [34]	ResNet-12	66.8 ± 0.9	81.7 ± 0.6	64.8 ± 0.6	82.1 ± 0.4	
Meta-Baseline [7]	ResNet-12	68.0 ± 0.3	82.7 ± 0.2	63.2 ± 0.2	79.3 ± 0.2	
COMPAS [18]	ResNet-12	67.1 ± 0.5	82.3 ± 0.2	65.7 ± 0.5	82.0 ± 0.3	
DeepEMD $[43]$	$\operatorname{ResNet-12}$	67.3 ± 0.6	82.7 ± 0.4	65.9 ± 0.8	82.4 ± 0.5	

Can Parts help Few-shot Learning? Based on the initial results, we are interested in whether parts can actually help few-shot learning. We take two representative methods-COMPAS [18] and DeepEMD [47] here to validate it.

COMPAS [18] originally constructs a part dictionary D of important regions by using K-Means on the feature representations of the backbone and further builds a map dictionary S of the spatial activation distribution of these regions. To exploit the part annotation on COMPAS, we first convert the part segmentation annotations into bounding boxes followed by using pre-trained backbone to extract feature representations of these parts. Then we apply K-Means on the feature representations to obtain a better initialization of the part dictionary D. Similarly, we directly calculate the spatial distribution of these bounding boxes to offer a better initialization of the map dictionary S.

For DeepEMD [47], it originally tries to get a dense representations of images and then compute the Earth Mover's Distance to generate the optimal matching flows between the representation sets of images. The optimal matching cost is further used as the distance metric to measure the similarity of two images. In the

updated version DeepEMD V2 [46], the authors find that instead of generating dense representations of images, it is better to randomly sample a set of regions in the images and only compute the EMD between these regions. Thus it is natural to replace these randomly generated regions with annotated part regions and their concatenated regions to see if the results get better.

Table 9 summaries the results of COMPAS and DeepEMD w/ and w/o using part annotations on PartImageNet meta-testing set. Both methods achieve non-trivial improvement with explicit exploitation of annotated parts. Concretely speaking, COMPAS gets a 0.9% and 0.6% performance gain and DeepEMD achieves a 1.2% and 0.9% performance gain on 1-shot and 5-shot scenarios respectively. This reveals the great potential of introducing parts into few-shot learning. Potential research direction lies at how to integrate part detector into current few-shot learning pipeline. We leave such interesting works to the future.

Table 9. Experimental results of COMPAS and DeepEMD w & w/o using part annotations. Average classification accuracies(%) with 95% confidence intervals are reported.

Model	Backhono	PartImageNet 5-way		
Model	Dackbolle	1-shot	5-shot	
COMPAS [18]	$\operatorname{ResNet-12}$	67.1 ± 0.5	82.3 ± 0.2	
COMPAS w/ Part Annotations	$\operatorname{ResNet-12}$	68.0 ± 0.5	82.9 ± 0.3	
DeepEMD [47]	$\operatorname{ResNet-12}$	67.3 ± 0.6	82.7 ± 0.4	
DeepEMD w/ Part Annotations	$\operatorname{ResNet-12}$	68.5 ± 0.7	83.6 ± 0.3	

5 Conclusion

Parts provide a good intermediate representation of objects that have many advantages. Once obtained, they can be exploited to increase the accuracy of recognition, localization and benefit many downstream tasks such as pose estimation. In this work, we introduce PartImageNet—a large, high-quality dataset with part annotation on a general set of classes. A set of new baselines are further set of different vision tasks including semantic part segmentation, object segmentation and few-shot learning. We further show that introducing parts is beneficial to object segmentation and few-shot learning. We also reveal that existing works have certain limitations which hinder them to produce satisfactory semantic part segmentation results under complex backgrounds and variations. We hope that with the propose of our PartImageNet, we could attract more attention to the research of part-based models to address these difficulties and make parts great again.

Acknowledgements. The authors gratefully acknowledge supports from NSF BCS-1827427 and ONR N00014-21-1-2812. AK acknowledges support via his Emmy Noether Research Group funded by the German Science Foundation (DFG) under Grant No. 468670075.

References

- 1. Biederman, I.: Recognition-by-components: a theory of human image understanding. Psychological review **94**(2), 115 (1987)
- Chang, A.X., Funkhouser, T., Guibas, L., Hanrahan, P., Huang, Q., Li, Z., Savarese, S., Savva, M., Song, S., Su, H., et al.: Shapenet: An information-rich 3d model repository. arXiv preprint arXiv:1512.03012 (2015)
- 3. Chen, L.C., Papandreou, G., Schroff, F., Adam, H.: Rethinking atrous convolution for semantic image segmentation. arXiv preprint arXiv:1706.05587 (2017)
- Chen, L.C., Zhu, Y., Papandreou, G., Schroff, F., Adam, H.: Encoder-decoder with atrous separable convolution for semantic image segmentation. In: Proceedings of the European conference on computer vision (ECCV). pp. 801–818 (2018)
- Chen, X., Mottaghi, R., Liu, X., Fidler, S., Urtasun, R., Yuille, A.: Detect what you can: Detecting and representing objects using holistic models and body parts. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 1971–1978 (2014)
- Chen, X., Yuille, A.L.: Articulated pose estimation by a graphical model with image dependent pairwise relations. Advances in neural information processing systems 27 (2014)
- Chen, Y., Liu, Z., Xu, H., Darrell, T., Wang, X.: Meta-baseline: Exploring simple meta-learning for few-shot learning. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 9062–9071 (2021)
- Cordts, M., Omran, M., Ramos, S., Rehfeld, T., Enzweiler, M., Benenson, R., Franke, U., Roth, S., Schiele, B.: The cityscapes dataset for semantic urban scene understanding. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 3213–3223 (2016)
- Deng, J., Dong, W., Socher, R., Li, L.J., Li, K., Fei-Fei, L.: Imagenet: A largescale hierarchical image database. In: 2009 IEEE conference on computer vision and pattern recognition. pp. 248–255. Ieee (2009)
- Fei-Fei, L., Fergus, R., Perona, P.: One-shot learning of object categories. IEEE transactions on pattern analysis and machine intelligence 28(4), 594–611 (2006)
- Felzenszwalb, P.F., Huttenlocher, D.P.: Pictorial structures for object recognition. International journal of computer vision 61(1), 55–79 (2005)
- Finn, C., Abbeel, P., Levine, S.: Model-agnostic meta-learning for fast adaptation of deep networks. In: International Conference on Machine Learning. pp. 1126– 1135. PMLR (2017)
- Fischler, M.A., Elschlager, R.A.: The representation and matching of pictorial structures. IEEE Transactions on computers 100(1), 67–92 (1973)
- de Geus, D., Meletis, P., Lu, C., Wen, X., Dubbelman, G.: Part-aware panoptic segmentation. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 5485–5494 (2021)
- Girshick, R., Felzenszwalb, P., McAllester, D.: Object detection with grammar models. Advances in Neural Information Processing Systems 24, 442–450 (2011)
- Gong, K., Liang, X., Li, Y., Chen, Y., Yang, M., Lin, L.: Instance-level human parsing via part grouping network. In: Proceedings of the European Conference on Computer Vision (ECCV). pp. 770–785 (2018)
- He, J., Chen, J.N., Liu, S., Kortylewski, A., Yang, C., Bai, Y., Wang, C., Yuille, A.: Transfg: A transformer architecture for fine-grained recognition. arXiv preprint arXiv:2103.07976 (2021)

- 16 J. He et al.
- He, J., Kortylewski, A., Yuille, A.: Compas: Representation learning with compositional part sharing for few-shot classification. arXiv preprint arXiv:2101.11878 (2021)
- Kirillov, A., Girshick, R., He, K., Dollár, P.: Panoptic feature pyramid networks. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 6399–6408 (2019)
- B.M., Salakhutdinov, Tenenbaum, J.B.: Human-level 20. Lake, R., con- cept learning through probabilistic program induction. Science https://doi.org/10.1126/science.aab3050, **350**(6266), 1332 - 1338(2015).https://www.science.org/doi/abs/10.1126/science.aab3050
- Liang, X., Gong, K., Shen, X., Lin, L.: Look into person: Joint body parsing & pose estimation network and a new benchmark. IEEE transactions on pattern analysis and machine intelligence 41(4), 871–885 (2018)
- Lin, T.Y., Maire, M., Belongie, S., Hays, J., Perona, P., Ramanan, D., Dollár, P., Zitnick, C.L.: Microsoft coco: Common objects in context. In: European conference on computer vision. pp. 740–755. Springer (2014)
- Liu, Q., Kortylewski, A., Zhang, Z., Li, Z., Guo, M., Liu, Q., Yuan, X., Mu, J., Qiu, W., Yuille, A.: Cgpart: A part segmentation dataset based on 3d computer graphics models. arXiv preprint arXiv:2103.14098 (2021)
- Liu, Y., Zhang, X., Zhang, S., He, X.: Part-aware prototype network for few-shot semantic segmentation. In: European Conference on Computer Vision. pp. 142– 158. Springer (2020)
- 25. Lorenz, D., Bereska, L., Milbich, T., Ommer, B.: Unsupervised part-based disentangling of object shape and appearance. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) (June 2019)
- Mo, K., Zhu, S., Chang, A.X., Yi, L., Tripathi, S., Guibas, L.J., Su, H.: Partnet: A large-scale benchmark for fine-grained and hierarchical part-level 3d object understanding. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 909–918 (2019)
- Reddy, N.D., Vo, M., Narasimhan, S.G.: Carfusion: Combining point tracking and part detection for dynamic 3d reconstruction of vehicles. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 1906–1915 (2018)
- Ren, M., Triantafillou, E., Ravi, S., Snell, J., Swersky, K., Tenenbaum, J.B., Larochelle, H., Zemel, R.S.: Meta-learning for semi-supervised few-shot classification. arXiv preprint arXiv:1803.00676 (2018)
- Shao, S., Li, Z., Zhang, T., Peng, C., Yu, G., Zhang, X., Li, J., Sun, J.: Objects365: A large-scale, high-quality dataset for object detection. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 8430–8439 (2019)
- Snell, J., Swersky, K., Zemel, R.S.: Prototypical networks for few-shot learning. arXiv preprint arXiv:1703.05175 (2017)
- Song, X., Wang, P., Zhou, D., Zhu, R., Guan, C., Dai, Y., Su, H., Li, H., Yang, R.: Apollocar3d: A large 3d car instance understanding benchmark for autonomous driving. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 5452–5462 (2019)
- 32. Sun, Y., Zheng, L., Yang, Y., Tian, Q., Wang, S.: Beyond part models: Person retrieval with refined part pooling (and a strong convolutional baseline). In: Proceedings of the European conference on computer vision (ECCV). pp. 480–496 (2018)
- Thewlis, J., Bilen, H., Vedaldi, A.: Unsupervised learning of object landmarks by factorized spatial embeddings. In: Proceedings of the IEEE international conference on computer vision. pp. 5916–5925 (2017)

- Tian, Y., Wang, Y., Krishnan, D., Tenenbaum, J.B., Isola, P.: Rethinking few-shot image classification: a good embedding is all you need? In: Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XIV 16. pp. 266–282. Springer (2020)
- Wah, C., Branson, S., Welinder, P., Perona, P., Belongie, S.: The Caltech-UCSD Birds-200-2011 Dataset. Tech. Rep. CNS-TR-2011-001, California Institute of Technology (2011)
- Wang, J., Zhang, Z., Xie, C., Premachandran, V., Yuille, A.: Unsupervised learning of object semantic parts from internal states of cnns by population encoding. arXiv preprint arXiv:1511.06855 (2015)
- 37. Wang, P., Shen, X., Lin, Z., Cohen, S., Price, B., Yuille, A.L.: Joint object and part segmentation using deep learned potentials. In: Proceedings of the IEEE International Conference on Computer Vision. pp. 1573–1581 (2015)
- 38. Weber, M., Welling, M., Perona, P.: Unsupervised learning of models for recognition. In: European conference on computer vision. pp. 18–32. Springer (2000)
- Wu, J., Zhang, T., Zhang, Y., Wu, F.: Task-aware part mining network for few-shot learning. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 8433–8442 (2021)
- Xia, F., Wang, P., Chen, X., Yuille, A.L.: Joint multi-person pose estimation and semantic part segmentation. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 6769–6778 (2017)
- Xiang, Y., Mottaghi, R., Savarese, S.: Beyond pascal: A benchmark for 3d object detection in the wild. In: IEEE winter conference on applications of computer vision. pp. 75–82. IEEE (2014)
- 42. Xie, E., Wang, W., Yu, Z., Anandkumar, A., Alvarez, J.M., Luo, P.: Segformer: Simple and efficient design for semantic segmentation with transformers. arXiv preprint arXiv:2105.15203 (2021)
- Xu, W., Wang, H., Tu, Z., et al.: Attentional constellation nets for few-shot learning. In: International Conference on Learning Representations (2020)
- 44. Yi, L., Kim, V.G., Ceylan, D., Shen, I.C., Yan, M., Su, H., Lu, C., Huang, Q., Sheffer, A., Guibas, L.: A scalable active framework for region annotation in 3d shape collections. ACM Transactions on Graphics (ToG) 35(6), 1–12 (2016)
- 45. Yuille, A.L., Hallinan, P.W., Cohen, D.S.: Feature extraction from faces using deformable templates. International journal of computer vision 8(2), 99–111 (1992)
- Zhang, C., Cai, Y., Lin, G., Shen, C.: Deepemd: Differentiable earth mover's distance for few-shot learning (2020)
- Zhang, C., Cai, Y., Lin, G., Shen, C.: Deepemd: Few-shot image classification with differentiable earth mover's distance and structured classifiers. In: Proceedings of the IEEE/CVF conference on computer vision and pattern recognition. pp. 12203– 12213 (2020)
- Zhao, J., Li, J., Cheng, Y., Sim, T., Yan, S., Feng, J.: Understanding humans in crowded scenes: Deep nested adversarial learning and a new benchmark for multi-human parsing. In: Proceedings of the 26th ACM international conference on Multimedia. pp. 792–800 (2018)
- Zhou, B., Zhao, H., Puig, X., Fidler, S., Barriuso, A., Torralba, A.: Scene parsing through ade20k dataset. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 633–641 (2017)
- Zhu, S.C., Mumford, D.: A stochastic grammar of images. Now Publishers Inc (2007)