

Self-Correction for Human Parsing

Peike Li¹, Yunqiu Xu², Yunchao Wei¹, Yi Yang^{1,2}

¹ReLER Lab, Centre for Artificial Intelligence, University of Technology Sydney

²Baidu Research

peike.li@yahoo.com, imyunqiuXu@gmail.com, {yunchao.wei, yi.yang}@uts.edu.au

Abstract

Labeling pixel-level masks for fine-grained semantic segmentation tasks, e.g. human parsing, remains a challenging task. The ambiguous boundary between different semantic parts and those categories with similar appearance usually are confusing, leading to unexpected noises in ground truth masks. To tackle the problem of learning with label noises, this work introduces a purification strategy, called Self-Correction for Human Parsing (SCHP), to progressively promote the reliability of the supervised labels as well as the learned models. In particular, starting from a model trained with inaccurate annotations as initialization, we design a cyclically learning scheduler to infer more reliable pseudo-masks by iteratively aggregating the current learned model with the former optimal one in an online manner. Besides, those correspondingly corrected labels can in turn to further boost the model performance. In this way, the models and the labels will reciprocally become more robust and accurate during the self-correction learning cycles. Benefiting from the superiority of SCHP, we achieve the best performance on two popular single-person human parsing benchmarks, including LIP and Pascal-Person-Part datasets. Our overall system ranks 1st in CVPR2019 LIP Challenge. Code is available at [this url](#).

1. Introduction

Human parsing, as a fine-grained semantic segmentation task, aims to assign each image pixel from the human body to a semantic category, e.g. arm, leg, dress, skirt. Understanding the detailed semantic parts of human is crucial in several potential application scenarios, including image editing, human analysis, virtual try-on and virtual reality. Recent advances on fully convolutional neural networks [22, 3] achieves various of well-performing methods for the human parsing task [18, 27].

To learn reliable models for human parsing, a large amount of pixel-level masks are required for supervision. However, labeling pixel-level annotations for human pars-

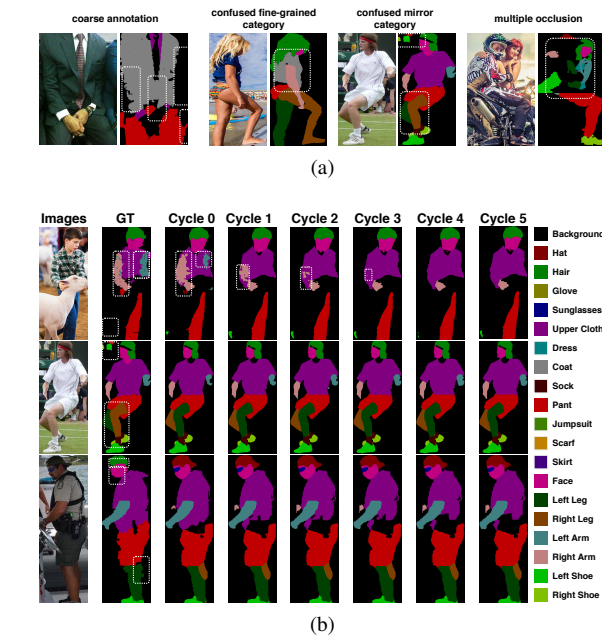


Figure 1: (a) Different types of label noises exist in the ground truth. (b) Our self-correction mechanism progressively promotes the reliability of the supervised labels. Label noises are emphasized by white dotted boxes. Better zoom in to see the details.

ing is much harder than those traditional pixel-level understanding tasks. In particular, for those traditional semantic segmentation tasks [22, 3], all the pixels belonging to one instance share the same semantic label, which is usually easy to be identified by annotators. Differently, the human parsing task requires annotators to carefully distinguish those semantic parts of one person. Moreover, the situation will become even more challenging when the annotator got confused by the ambiguous boundaries between different semantic parts.

Due to these factors, there inevitably exists different types of label noises (as illustrated in Figure 1a) caused by the careless observations by annotators. This incomplete

and low quality of the annotation labels will set a significant obstacle, which is usually ignored and prevents the performance of human parsing from increasing to a higher level. In this work, we investigate the problem of learning with noise in the human parsing task. Our target is to improve the model performance and generalization by progressively refining the noisy labels during the training stage.

In this paper, we introduce a purification strategy named Self-Correction for Human Parsing (SCHP), which can progressively promote the reliability of the supervised labels, as well as the learned models during the training process. Concretely, the whole SCHP pipeline can be divided into two sub-procedures, *i.e.* the model aggregation and the label refinement procedure. Starting from a model trained on inaccurate annotations as initialization, we design a cyclically learning scheduler to infer more reliable pseudo-masks by iteratively aggregating the current learned model with the former optimal one in an online manner. Besides, those corrected labels can in turn to boost the model performance, simultaneously. In this way, the self-correction mechanism will enable the model or labels to mutually promote its counterpart, leading to a more robust model and accurate label masks as training goes on.

Besides, to tackle the problem of ambiguous boundaries between different semantic parts, we introduce a network architecture called Augmented Context Embedding with Edge Perceiving (A-CE2P). In principle, our network architecture is an intuitive generalization and augmentation of the CE2P [27] framework. We introduce a consistency constraint term to augment the CE2P, so that the edge information is not only implicitly facilitated the parsing result by feature map level fusion, but also explicitly constrained between the parsing and edge prediction. Note that we do not claim any novelty of our architecture structure, but only the superiority of the performance.

On the whole, our major contributions can be summarized as follows,

- We propose a simple yet effective self-correction strategy SCHP for the human parsing task by online model aggregating and label refining which could mutually promote the model performance and label accuracy.
- We introduce a general architecture framework A-CE2P for human parsing that both implicitly and explicitly captures the boundary information along with the parsing information.
- We extensively investigate our SCHP on two popular human parsing benchmarks. Our method achieves the new state-of-the-art. In particular, we achieve the mIoU score of 59.36 on the large scale benchmark LIP, which outperforms the previous closest approach by 6.2 points.

2. Related Work

Human Parsing. Several different aspects of the human parsing task have been studied. Some early works [31, 18] utilized pose estimation together with the human parsing simultaneously as a multi-task learning problem. In [27], they cooperated the edge prediction with human parsing to accurately predict the boundary area. Most of the prior works assumed the fact that ground truth labels are correct and well-annotated. However, due to time and cost consuming, there inevitably exists lots of different label noises (as shown in Figure 1a). Meanwhile, it is impracticable to clean the pixel-level labels manually. Guided by this intuition, we try to tackle this problem via a novel self-correction mechanism in this paper.

Pseudo-Labeling. Pseudo-labeling [17, 26] is a typical technique used in semi-supervised learning. In semi-supervised learning setting, they assign pseudo-labels to the unlabeled data. However, in our fully supervised learning scheme, we are unable to locate the label noises, thus all ground truth labels are treated equally. From the perspective of distillation, the generated pseudo-label data contains much so-called *dark knowledge* [13] which could serve as a purification signal. Inspired by these findings, we design a cyclically learning scheduler to infer more reliable pseudo-masks by iteratively aggregating the current learned model with the former optimal one in an online manner. Also those corrected labels can in turn to boost the model performance, simultaneously.

Self-Ensembling. There is a line of researches [16, 29, 15] that exploit self-ensembling methods in various scenarios. For example, [29] averaged model weights as self-ensembling and adopted in the semi-supervised learning task. In [15], they averaged the model weight and led to better generalization. Different from their method, our proposed self-correction approach is to correct the noisy training label via a model and label mutually promoting process. By an online manner, we average both model weights and the predictions simultaneously. To the best of our knowledge, we make a first attempt to formulate the label noise problem as the mutual model and label optimization in fine-grained semantic segmenting to boost the performance. Furthermore, our proposed method is online training with a cyclical learning scheduler and only exhaust little extra computation.

3. Methodology

3.1. Revisiting CE2P

CE2P [27] is a well-performing framework for the human parsing task. In the CE2P network, they cooperate the edge prediction with human parsing to accurately predict the boundary area. Concretely, CE2P consists of three key branches, *i.e.* *parsing* branch, *edge* branch and *fusion*

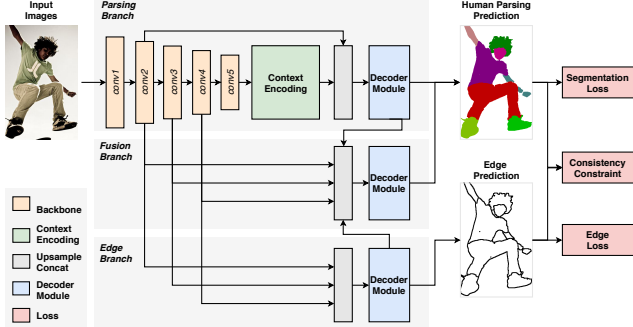


Figure 2: An overview of the Augmented-CE2P framework.

branch. In particular, the *edge* branch is employed to generate class-agnostic boundary maps. In the *fusion* branch, both semantic-aware feature representations from the *parsing* branch and the boundary-aware feature representations from the *edge* branch are concatenated to further produce a refined human parsing prediction.

Although CE2P is a framework that has already incorporated the most useful functions from the semantic segmentation community. However, there are still some aspects that could be further strengthened. First, the conventional cross-entropy loss indirectly optimizes mean intersection-over-union (mIoU) metric, which is a crucial metric to reveal the comprehensive performance of the model. Second, CE2P only implicitly facilitates the parsing results with the edge predictions by feature-level fusion. There is no explicit constraint to ensure the parsing results maintaining the same geometry shape of the boundary predictions.

Moreover, few efforts have been made to investigate the versatility of the CE2P framework *i.e.* the ability to accommodate other modules. Based on the key function, the *parsing* branch can be divided into three modules, *i.e.* backbone module, context encoding module and decoder module. Concretely, the backbone module could be plugged in with any fully-convolutional structure backbone such as ResNet-based [12] semantic segmentation network. The context encoding module utilizes the global context information to distinguish the fine-grained categories information. This module could be any effective context discovering module, *e.g.* feature pyramid based approaches like PSP [33], ASPP [4], or attention-based modules like OCNet [32]. More detailed network architecture could refer to our code.

3.2. Augmented-CE2P

The Augmented-CE2P (A-CE2P) is an intuitive generalization and augmentation of the CE2P, which can yield increased performance gain by augmenting additional powerful modules. In this work, our self-correction training employs the A-CE2P as the basic framework for conduct-

ing human parsing. We demonstrate the overview of the A-CE2P framework in Figure 2. Notably, several unique characteristics of A-CE2P are described as follows.

Targeted Learning Objectives. For an image I , suppose the human parsing ground truth label is \hat{y}_k^n and the parsing prediction is y_k^n , where n is the number of pixels for the k -th class. We define the pixel-level supervised objective using conventional cross-entropy loss as:

$$\mathcal{L}_{cls} = -\frac{1}{N} \sum_k \sum_n \hat{y}_k^n \log p(y_k^n). \quad (1)$$

here N is the number of pixels, K is the number of classes.

It is known that conventional cross-entropy loss is usually convenient to train a neural network, but it facilitates mean intersection-over-union (mIoU) indirectly. To tackle this issue, following by [1], we additionally introduce a tractable surrogate loss function for optimizing the mIoU directly. The final parsing loss function can be defined as a combination of the cross-entropy loss and the mIoU loss \mathcal{L}_{miou} ,

$$\mathcal{L}_{parsing} = \mathcal{L}_{cls} + \mathcal{L}_{miou}. \quad (2)$$

Consistency Constraint. In the CE2P, the balanced cross-entropy loss \mathcal{L}_{edge} is adopted to optimize the edge prediction, so that the learned edge-aware features can help distinguish human parts and facilitate human parsing via the fusion branch indirectly.

In the A-CE2P, we propose to further exploit the predicted boundary information by explicitly maintaining the consistency between the parsing prediction and the boundary prediction, *i.e.* ensure that the predicted parsing result matches the predicted edge as exact as possible. Intuitively, we add a constraint term to penalized the mismatch:

$$\mathcal{L}_{consistent} = \frac{1}{|N^+|} \sum_{n \in N^+} |\tilde{e}^n - e^n|, \quad (3)$$

where e^n is the edge maps predicted from the *edge* branch and \tilde{e}^n is the edge maps generated from the parsing result y_k^n . To prevent the non-edge pixels dominate the loss, we only allow the positive edge pixels $n \in N^+$ for contributing the consistency constraint term.

In brief, the overall learning objective of our framework is

$$\mathcal{L} = \lambda_1 \mathcal{L}_{edge} + \lambda_2 \mathcal{L}_{parsing} + \lambda_3 \mathcal{L}_{consistent}, \quad (4)$$

where λ_1, λ_2 and λ_3 are hyper-parameters to control the contribution among these three losses. We jointly train the model in an end-to-end fashion by minimizing \mathcal{L} .

3.3. Learning with Noise via Self-Correction

Based on the A-CE2P, we proposed the self-correction method that allows us to refine the label and get a robust

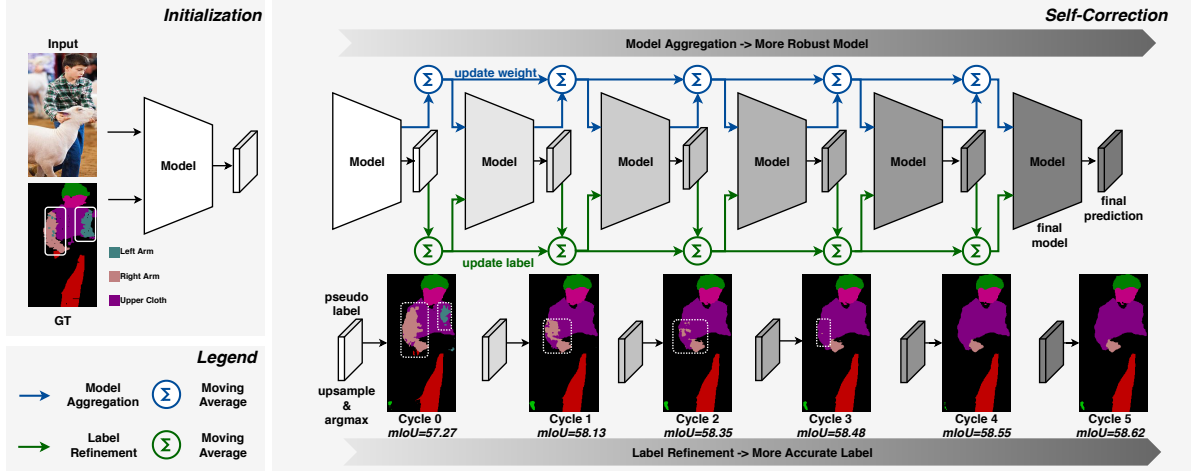


Figure 3: Illustration of the SCHP Pipeline. The self-correction mechanism will enable the model or labels to mutually promote its counterpart, leading to a more robust model and accurate label masks as training goes on. Label noises are specially marked with white boxes.

model via an online mutual improvement process, illustrated in Figure 3.

Training Strategy. Our proposed self-correction training strategy is a model and label aggregating process, which can promote the model performance and refine the ground truth labels iteratively. This promotion relies on the initial performance of the model. In other words, if intermediate results generated by the network are not accurate enough, they may potentially harm the iteration process. Therefore, we start to run our proposed self-correction algorithm after a good initialization, *i.e.* when the training loss starts to flatten with the original noisy labels. To make a fair comparison with other methods, we shorten the initial stage and keep the total training epochs as the same. After the initialization stage, we adopt a cyclically learning scheduler with warm restarts. Each cycle totally contains T epochs. In practice, we use a cosine annealing learning rate scheduler with cyclical restart [23]. Formally, η_{max} and η_{min} are set to the initial learning rate and final learning rate, while T_{cur} is the number of epochs since the last restart. Thus, the overall learning rate can be formulated as,

$$\eta = \eta_{min} + \frac{1}{2}(\eta_{max} - \eta_{min})(1 + \cos(\frac{T_{cur}}{T}\pi)). \quad (5)$$

Online Model Aggregation. We aim to discover all the potential information from the past optimal models to improve the performance of the future model. In our cyclical training strategy, intuitively, the model will converge to a local-minimum at the end of each cycle. And there has great model disparity among these sub-optimal models. Here we denote the set of all the sub-optimal model we get after each cycle as $\Omega = \{\hat{\omega}_0, \hat{\omega}_1, \dots, \hat{\omega}_M\}$ and M is the total number of cycles.

At the end of each cycle, we aggregate the current model weight $\hat{\omega}$ with the former sub-optimal one $\hat{\omega}_{m-1}$ to achieve a new model weight $\hat{\omega}_m$,

$$\hat{\omega}_m = \frac{m}{m+1}\hat{\omega}_{m-1} + \frac{1}{m+1}\hat{\omega}, \quad (6)$$

where m denotes the current cycle number and $0 \leq m \leq M$.

After updating the current model weight with the former optimal one from the last cycle, we forward all the training data for one epoch to re-estimate the statistics of the parameters (*i.e.* moving average and standard deviation) in all batch normalization [14] layers. During these successive cycles of model aggregation, the network leads to wider model optima as well as improved model’s generalization ability.

Online Label Refinement. It is known that soft, multi-class labels may contain dark information [13]. We aim to explore all these dark information to improve the performance and alleviate the label noises. After updating the model weight as mentioned in Eq. (6), we also update the ground truth of training labels. These generated pseudo-masks are more unambiguous, smooth and have the relational information between the fine-grain categories, which are taken as the supervised signal for the next cycle’s optimization. During successive cycles of pseudo-label refinement, this improves the network performance as well as the generalization ability of the model. Also, these pseudo-masks potentially alleviate or eliminate the noise in the original ground truth. Here we denote the predicted label after each cycles as $Y = \{\hat{y}_0, \hat{y}_1, \dots, \hat{y}_M\}$. Same as the model weight averaging process, we update the ground truth label

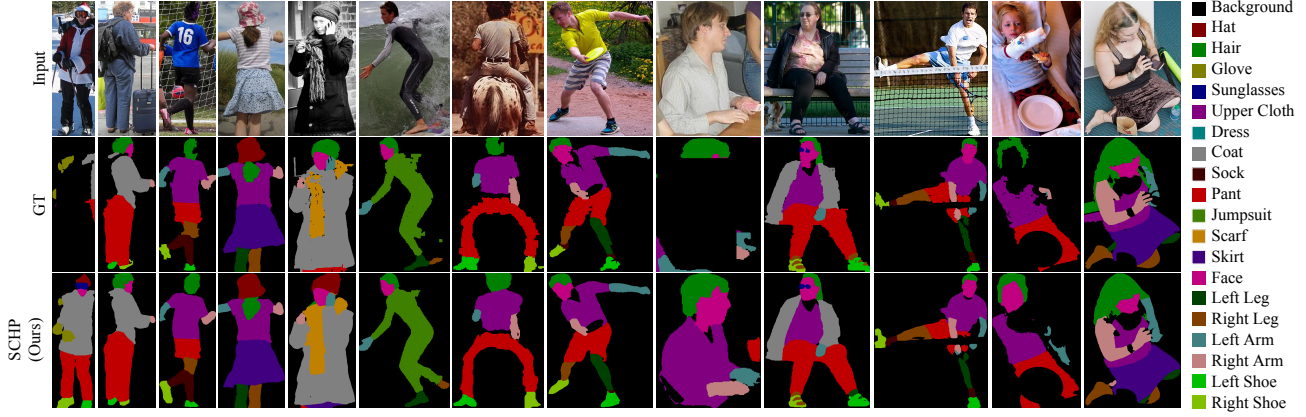


Figure 4: Visualization of SCHP results on LIP validation set. Note in most cases, our SCHP human parsing prediction is even better than the ground truth label. Zoom in to see details.

Algorithm 1: Self-Correction for Human Parsing

Input: Initialized model weight $\hat{\omega}_0$, original ground truth labels \hat{y}_0 , cycle epoch length T , total number of iterations M

Output: Final network model $\hat{\omega}_M$

Initialize the model weight $\hat{\omega} \leftarrow \hat{\omega}_0$;

for $m \leftarrow 1, 2, \dots, M$ **do**

for $t \leftarrow 1, 2, \dots, T$ **do**

 Update the learning rate η by Eq. (5);

for each batch in training set do

 Calculate loss \mathcal{L} by Eq. (4) using \hat{y}_{m-1} ;

 Gradient descending $\hat{\omega} \leftarrow \hat{\omega} - \eta \nabla \mathcal{L}$;

end

end

 Model aggregation by Eq. (6) to update $\hat{\omega}_m$;

 Re-calculate the BN layer parameters;

 Re-calculate the pseudo-mask \hat{y} using $\hat{\omega}_m$;

 Label refinement by Eq. (7) to update \hat{y}_m ;

end

as follows,

$$\hat{y}_m = \frac{m}{m+1} \hat{y}_{m-1} + \frac{1}{m+1} \hat{y}, \quad (7)$$

where \hat{y} is the generated pseudo-mask by model $\hat{\omega}_m$. The detail of our proposed self-correction procedure is summarized in Algorithm 1.

Note that the model and label are mutual improved step-by-step after each cyclical training process. The whole process is training in an online manner and does not need any extra training epochs. In addition, there is barely no extra computation required.

3.4. Discussion

Can SCHP generalize to other tasks? Our approach has no assumption for the data type. But the self-correction is based on the soft pseudo-label generated during the process. Thus our method could be applied in some other task such as classification and segmentation, but may be not applicable to regression tasks like detection.

Can SCHP still benefit with clean data? Although we could achieve more performance gain with our proposed self-correction process on noisy datasets. However, when the ground truth is relatively clean, the online model aggregation process could serve as a self-model ensembling, which could lead to better performance and generalization. Still, the online label refinement process benefits from discovering the *dark knowledge* using pseudo-mask instead of one-hot ground truth pixel-level label.

Transduction vs. Induction. In this work, we mainly focus on the supervised training scheme. Nevertheless, our approach can also be operated under semi-supervised learning manner, *i.e.* we assume that all the test images are available at once and utilize all the test samples for self-correction process together with the training images jointly.

4. Experiments

In this section, we perform a comprehensive comparison of our SCHP with other single-person human parsing state-of-the-art methods along with thorough ablation experiments to demonstrate the contribution of each component.

4.1. Datasets and Evaluation Metrics

Datasets. We evaluate our proposed method on two single human parsing benchmarks, including LIP [18] and PASCAL-Person-Part [6].

Method	hat	hair	glove	s-glass	u-clot	dress	coat	sock	pant	j-suit	scarf	skirt	face	l-arm	r-arm	l-leg	r-leg	l-shoe	r-shoe	bkg	mIoU
Attention [5]	58.87	66.78	23.32	19.48	63.20	29.63	49.70	35.23	66.04	24.73	12.84	20.41	70.58	50.17	54.03	38.35	37.70	26.20	27.09	84.00	42.92
DeepLab [3]	59.76	66.22	28.76	23.91	64.95	33.68	52.86	37.67	68.05	26.15	17.44	25.23	70.00	50.42	53.89	39.36	38.27	26.95	28.36	84.09	44.80
SSL [11]	58.21	67.17	31.20	23.65	63.66	28.31	52.35	39.58	69.40	28.61	13.70	22.52	74.84	52.83	55.67	48.22	47.49	31.80	29.97	84.64	46.19
MMAN [24]	57.66	65.63	30.07	20.02	64.15	28.39	51.98	41.46	71.03	23.61	9.65	23.20	69.54	55.30	58.13	51.90	52.17	38.58	39.05	84.75	46.81
MuLA [25]	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	49.30
JPPNet [18]	63.55	70.20	36.16	23.48	68.15	31.42	55.65	44.56	72.19	28.39	18.76	25.14	73.36	61.97	63.88	58.21	57.99	44.02	44.09	86.26	51.37
CE2P [27]	65.29	72.54	39.09	32.73	69.46	32.52	56.28	49.67	74.11	27.23	14.19	22.51	75.50	65.14	66.59	60.10	58.59	46.63	46.12	87.67	53.10
A-CE2P w/o SCHP	69.59	73.02	45.21	35.59	69.85	35.97	56.96	51.06	75.79	30.41	22.00	27.07	75.79	68.54	70.30	67.83	66.90	53.53	54.08	88.11	56.88
A-CE2P w/ SCHP	69.96	73.55	50.46	40.72	69.93	39.02	57.45	54.27	76.01	32.88	26.29	31.68	76.19	68.65	70.92	67.28	66.56	55.76	56.50	88.36	58.62
A-CE2P w/ SCHP [†]	70.63	74.09	51.40	41.70	70.56	40.06	58.17	55.17	76.57	33.78	26.63	32.83	76.63	69.33	71.76	67.93	67.42	56.56	57.55	88.40	59.36

Table 1: Comparison with state-of-the-arts on LIP validation set. [†] designates the test time augmentation.

Method	head	torso	u-arm	l-arm	u-leg	l-leg	bkg	mIoU
Attention [5]	81.47	59.06	44.15	42.50	38.28	35.62	93.65	56.39
HAZN [30]	80.76	60.50	45.65	43.11	41.21	37.74	93.78	57.54
LG-LSTM [20]	82.72	60.99	45.40	47.76	42.33	37.96	88.63	57.97
SS-JPPNet [18]	83.26	62.40	47.80	45.58	42.32	39.48	94.68	59.36
MMAN [24]	82.58	62.83	48.49	47.37	42.80	40.40	94.92	59.91
G-LSTM [19]	82.69	62.68	46.88	47.71	45.66	40.93	94.59	60.16
Part FCN [31]	85.50	67.87	54.72	54.30	48.25	44.76	95.32	64.39
DeepLab [3]	-	-	-	-	-	-	-	64.94
WSPH [9]	87.15	72.28	57.07	56.21	52.43	50.36	97.72	67.60
PGN [10]	90.89	75.12	55.83	64.61	55.42	41.57	95.33	68.40
DPC [2]	88.81	74.54	63.85	63.73	57.24	54.55	96.66	71.34
A-CE2P w/ SCHP	87.00	72.27	64.10	63.44	56.57	55.00	96.07	70.63
A-CE2P w/ SCHP [†]	87.41	73.80	64.98	64.70	57.43	55.62	96.26	71.46

Table 2: Comparison with state-of-the-arts on PASCAL-Person-Part validation set. [†] designates the test time augmentation.

LIP [18] is the largest human parsing dataset, which contains 50,462 images with elaborated pixel-wise annotations with 19 semantic human part labels. The images collected from the real-world scenarios contain human appearing with challenging poses and views, heavily occlusions, various appearances and low-resolutions. The datasets are divided images into 30,462 images for train set, 10,000 images for validation set and 10,000 for test set.

PASCAL-Person-Part [6] is a relatively small dataset annotated from PASCAL VOC 2010, which contains 1,716 train images, 1,817 validation images. The ground truth label consists of six semantic parts including head, torso, upper/lower arms, upper/lower legs and one background class. This dataset is challenging due to large variations in scale.

Metrics. We report three standard metrics for the human parsing task, including pixel accuracy, mean accuracy, mean intersection over union (mIoU). Note the mIoU metric generally represents the overall parsing performance of the method.

Implementation Details. We choose the ResNet-

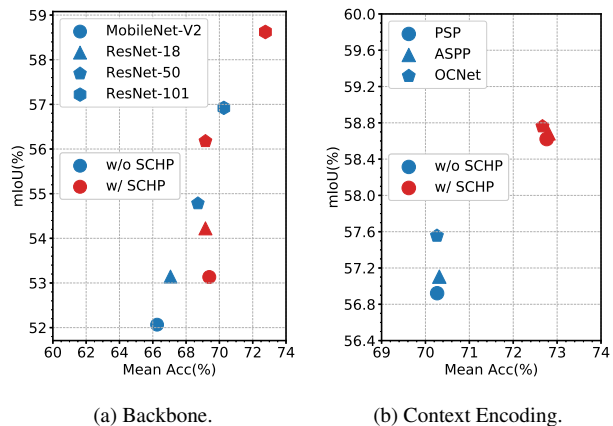


Figure 5: Effect of SCHP with different backbones and context encoding modules. All experiments are conducted on LIP validation set.

101 [12] as the backbone of the feature extractor and use an ImageNet [8] pre-trained weights. Specifically, we fix the first three residual layers and set the stride size of last residual layer to 1 with a dilation rate of 2. In this way, the final output is enlarged to 1/16 resolution size *w.r.t* the original image. We adopt pyramid scene parsing network [33] as the context encoding module. We use 473×473 as the input resolution. Training is done with a total batch size of 36. For our joint loss function, we set the weight of each term as $\lambda_1 = 1, \lambda_2 = 1, \lambda_3 = 0.1$. The initial learning rate is set as $7e-3$ with a linear increasing warm-up strategy for 10 epochs. We train our network for 150 epochs in total for a fair comparison, the first 100 epochs as initialization following 5 cycles each contains 10 epochs of the self-correction process.

4.2. Comparison with state-of-the-arts

In Table 1, we compare the performance of our network with other state-of-the-art methods on the LIP. It can be observed that even our baseline model outperforms all the

Loss			Pixel Acc	Mean Acc	mIoU
E	I	C			
-	-	-	86.93	65.12	53.75
✓	-	-	87.51	65.42	54.14
✓	✓	-	87.57	68.20	56.33
✓	✓	✓	87.68	68.79	56.88

Table 3: Each component of our loss function is evaluated on LIP validation set, including edge loss (E), IoU loss (I) and consistency constraint (C).

Method		Pixel Acc	Mean Acc	mIoU
MA	LR			
-	-	87.68	68.79	56.88
✓	-	87.90	71.27	57.94
-	✓	87.86	70.88	57.44
✓	✓	88.10	72.76	58.62

Table 4: The effect of our proposed model aggregation (MA) and label refinement (LR) strategy is evaluated on LIP validation set.

other state-of-the-art methods, which illustrates the effectiveness of the A-CE2P framework. In particular, we also apply test-time augmentation with multi-scale and horizontal flipping to make a fair comparison with others. Our SCHP outperforms the others with a large gain, achieving mIoU improvement by 6.26%, which is a significant boost considering the performance at this level. Our proposed approach achieves a large gain especially for some categories with few samples like *scarf*, *sunglasses* and some confusing categories such as *dress*, *skirt* and left-right confusion. The gains are mainly from using both model aggregation and label refinement for our self-correction process. Furthermore, the qualitative comparison between the predicted results of SCHP and ground truth annotations is shown in Figure 4. We can see that our SCHP can achieve even better parsing results than the ground truth ones.

To validate the generalization ability of our method, we further report the comparison with other state-of-the-arts on PASCAL-Person-Part in Table 2. It can be observed that our SCHP outperforms all the previous approaches. Particularly, our SCHP beats the DPC [2]. DPC is a network architecture search (NAS)-based method which easily achieves optimal results on the small dataset. Besides, instead of using more powerful backbone models such as Xception [7] in DPC, we only adopt ResNet-101 as the backbone. In addition, DPC leverages MS COCO [21] as additional data for pre-training, while our model is only pre-trained on ImageNet. All these results well demonstrate the superiority and generalization of our proposed approach.

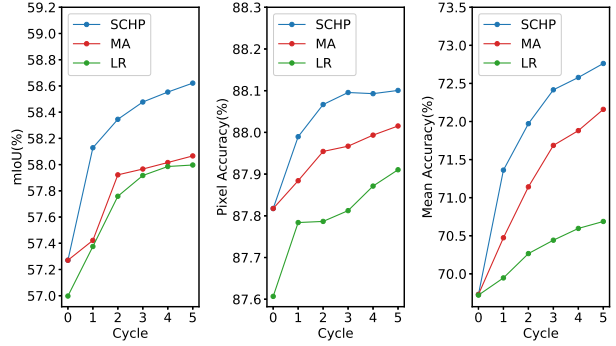


Figure 6: Performance curves w.r.t different training cycles. All experiments are conducted on LIP validation set. The mIoU, pixel accuracy, mean accuracy are reported, respectively.

4.3. Ablation Study

We perform extensive ablation experiments to illustrate the effect of each component in our SCHP. All experiments are conducted on LIP benchmark.

Alternatives Architectures. Since our proposed SHCP is a generic framework, we could merely plug-and-play with various backbones and context encoding modules. Figure 5a shows SHCP with different backbones from lightweight model MobileNet-V2 [28] to relatively heavy backbone ResNet-101 [12]. Interestingly, the lightweight MobileNet-V2 achieves the mIoU score of 52.06, with the benefit of SCHP leads to 53.13. This result is even better than some previous results [27] achieved by ResNet-101. We note that deeper network (18 vs. 50 vs. 101) tends to perform better. Regardless of different backbones, our SCHP consistently brings consistent positive gains, 1.08, 1.40 and 1.70, respectively in terms of mIoU. It suggests the robustness of our proposed method. As shown in Figure 5b, we further examine the robustness of our SCHP by varying the context encoding module. In particular, we choose three different types of modules, including multi-level global average pooling based module pyramid scene parsing network (PSP) [33], multi-level atrous convolutional pooling based module atrous spatial pyramid pooling (ASPP) [4] and attention-mechanism based module OCNNet [32]. Despite the similar basic performance of these three modules, our SCHP unfailingly obtains mIoU increased by 1.70, 1.30, 1.00 points for PSP, ASPP and OCNNet respectively. This further highlights the effectiveness of self-correction mechanism in our approach. Note that although we could achieve even better results with these advanced modules, the network structure modification is not the focus of this study. In our baseline model, we choose to use ResNet-101 as backbone and PSP as context encoding module.

Influence of Learning Objectives. Our network is

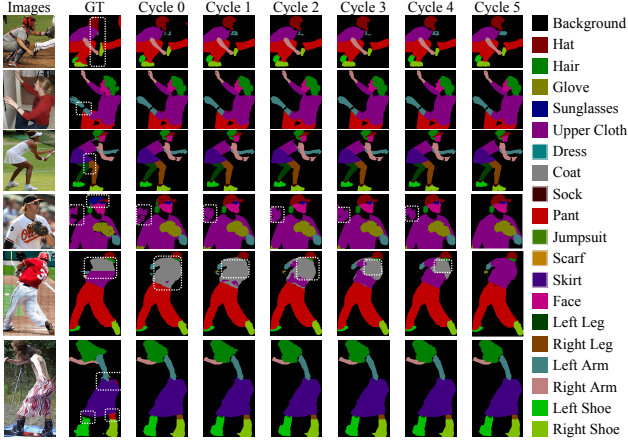


Figure 7: Visualization of our self-correction process in LIP train set. Label noises are emphasized with white dotted boxes. Better zoom in to see the details.

trained in an end-to-end manner with composite learning objectives describes as Eq. (4). An evaluation of different learning objectives is shown in Table 3. In this table, E denotes the binary cross-entropy loss to optimize the boundary prediction. I denotes the tractable surrogate function optimizing the mIoU metric. C denotes the consistency constraint term for maintain the consistency between parsing result and boundary result. Without all these three terms, only the basic cross-entropy loss function for parsing takes effect. By introducing the edge information, the performance improves mIoU by about 0.4. This gain is mainly due to the accurate prediction at the boundary area between semantic parts. Additionally, we compare the result further adding the IoU loss. As can be seen, the IoU loss significantly boosts the mean accuracy by 2.8 points and mIoU by 2.2 points, but the pixel accuracy almost remains the same level. This highlights that IoU loss largely resolves the challenge of prediction accuracy especially at some small area and infrequent categories. Finally, the result shows a gain of 0.59 mean accuracy and 0.55 mIoU when applying the consistency between parsing segmentation and edge prediction.

Effect of Self-Correction. In Table 4, we validate the effect of each component in our proposed SCHP, including the model aggregation (MA) process and the label refinement (LR) process. All experiments are conducted upon with A-CE2P framework. When there are no MA and LR involved, our method reduces to the conventional training process. By employing the MA process, the result shows a gain of mIoU by 1.1 points. While only benefit from the LR process, we achieve 0.6 point improvement. We achieve the best performance by simultaneously introducing these two processes. We can observe that the model aggregation and label refinement mutually promote each other in our SCHP.

To better qualitatively describe our SCHP, Figure 1b

shows the visualization of the generated pseudo-masks during the self-correction cycles. Note that all these pseudo-masks are up-sampled to the original size and applied argmax operation for better illustration. Label noises like inaccurate boundary, confused fine-grained categories, confused mirror categories, multi-person occlusion are alleviated and partly resolved during the self-correction cycles. Unsurprisingly, some of the boundaries of our corrected labels are prone to be more smooth than the ground truth labels. All these results demonstrate the effectiveness of our proposed self-correction method. Intuitively, our self-correction process is a mutual promoting process benefiting both model aggregation and label refinement. During the self-correction cycles, the model gets increasingly more robust, while by exploring the dark information from pseudo-masks produced by the enhanced model, the label noises are corrected in an implicit manner. The fact that corrected labels are smooth than the ground truth also illustrates the effectiveness of our model architecture design for combining the edge information.

Influence of Self-Correction Cycles. We achieve the goal of self-correction by a cyclically learning scheduler. The number of cycles is a virtual hyper-parameter for this process. To make a fair comparison with other methods [27], we maintain the entire training epoch unchanged. The performance curves are shown in Figure 6. It is evident that the performance consistently improves during the process, with the largest improvement after the first cycle and tendency saturates at the end. Our method may achieve even higher performance when extending more training epochs. It is noteworthy the performance of MA, LR and SCHP is not same at cycle 0. This small gap is caused due to the re-estimation of BatchNorm parameters. From the performance curve, We also intelligibly demonstrate the mutual benefit of the model aggregation and the label refinement process. More visualization of the self-correction process is illustrated in Figure 7.

5. Conclusion and Future Work

In this paper, we propose an effective self-correction strategy to deal with the label noises for the human parsing task. Our proposed method achieves the new state-of-the-art with a large margin gain. Moreover, the self-correction mechanism is a general strategy for training and can be incorporated into any frameworks to make further performance improvement. In the future, we would like to extend our method to multiple-person human parsing and video multiple-person human parsing tasks.

6. Acknowledgement

We thank Ting Liu and Tao Ruan for providing insights and expertise to improve this work.

References

- [1] Maxim Berman, Amal Rannen Triki, and Matthew B Blaschko. The Iovász-softmax loss: A tractable surrogate for the optimization of the intersection-over-union measure in neural networks. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 4413–4421, 2018. 3
- [2] Liang-Chieh Chen, Maxwell Collins, Yukun Zhu, George Papandreou, Barret Zoph, Florian Schroff, Hartwig Adam, and Jon Shlens. Searching for efficient multi-scale architectures for dense image prediction. In *Advances in Neural Information Processing Systems (NeurIPS)*, pages 8699–8710, 2018. 6, 7
- [3] Liang-Chieh Chen, George Papandreou, Iasonas Kokkinos, Kevin Murphy, and Alan L Yuille. Deeplab: Semantic image segmentation with deep convolutional nets, atrous convolution, and fully connected crfs. *IEEE Transactions on Pattern Recognition and Machine Intelligence*, 40(4):834–848, 2017. 1, 6
- [4] Liang-Chieh Chen, George Papandreou, Florian Schroff, and Hartwig Adam. Rethinking atrous convolution for semantic image segmentation. *arXiv preprint arXiv:1706.05587*, 2017. 3, 7
- [5] Liang-Chieh Chen, Yi Yang, Jiang Wang, Wei Xu, and Alan L Yuille. Attention to scale: Scale-aware semantic image segmentation. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 3640–3649, 2016. 6
- [6] Xianjie Chen, Roozbeh Mottaghi, Xiaobai Liu, Sanja Fidler, Raquel Urtasun, and Alan Yuille. Detect what you can: Detecting and representing objects using holistic models and body parts. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 1971–1978, 2014. 5, 6
- [7] François Chollet. Xception: Deep learning with depthwise separable convolutions. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 1251–1258, 2017. 7
- [8] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 248–255, 2009. 6
- [9] Hao-Shu Fang, Guansong Lu, Xiaolin Fang, Jianwen Xie, Yu-Wing Tai, and Cewu Lu. Weakly and semi supervised human body part parsing via pose-guided knowledge transfer. *arXiv preprint arXiv:1805.04310*, 2018. 6
- [10] Ke Gong, Xiaodan Liang, Yicheng Li, Yimin Chen, Ming Yang, and Liang Lin. Instance-level human parsing via part grouping network. In *European Conference on Computer Vision (ECCV)*, pages 770–785, 2018. 6
- [11] Ke Gong, Xiaodan Liang, Dongyu Zhang, Xiaohui Shen, and Liang Lin. Look into person: Self-supervised structure-sensitive learning and a new benchmark for human parsing. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 932–940, 2017. 6
- [12] Kaifeng He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 770–778, 2016. 3, 6, 7
- [13] Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. Distilling the knowledge in a neural network. *arXiv preprint arXiv:1503.02531*, 2015. 2, 4
- [14] Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. *arXiv preprint arXiv:1502.03167*, 2015. 4
- [15] Pavel Izmailov, Dmitrii Podoprikin, Timur Garipov, Dmitry Vetrov, and Andrew Gordon Wilson. Averaging weights leads to wider optima and better generalization. *arXiv preprint arXiv:1803.05407*, 2018. 2
- [16] Samuli Laine and Timo Aila. Temporal ensembling for semi-supervised learning. *arXiv preprint arXiv:1610.02242*, 2016. 2
- [17] Dong-Hyun Lee. Pseudo-label: The simple and efficient semi-supervised learning method for deep neural networks. In *Workshop on Challenges in Representation Learning, ICML*, volume 3, page 2, 2013. 2
- [18] Xiaodan Liang, Ke Gong, Xiaohui Shen, and Liang Lin. Look into person: Joint body parsing & pose estimation network and a new benchmark. *IEEE Transactions on Pattern Recognition and Machine Intelligence*, 41(4):871–885, 2018. 1, 2, 5, 6
- [19] Xiaodan Liang, Xiaohui Shen, Jiashi Feng, Liang Lin, and Shuicheng Yan. Semantic object parsing with graph lstm. In *European Conference on Computer Vision (ECCV)*, pages 125–143, 2016. 6
- [20] Xiaodan Liang, Xiaohui Shen, Donglai Xiang, Jiashi Feng, Liang Lin, and Shuicheng Yan. Semantic object parsing with local-global long short-term memory. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 3185–3193, 2016. 6
- [21] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *European Conference on Computer Vision (ECCV)*, pages 740–755, 2014. 7
- [22] Jonathan Long, Evan Shelhamer, and Trevor Darrell. Fully convolutional networks for semantic segmentation. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 3431–3440, 2015. 1
- [23] Ilya Loshchilov and Frank Hutter. Sgdr: Stochastic gradient descent with warm restarts. *arXiv preprint arXiv:1608.03983*, 2016. 4
- [24] Yawei Luo, Zhedong Zheng, Liang Zheng, Tao Guan, Junqing Yu, and Yi Yang. Macro-micro adversarial network for human parsing. In *European Conference on Computer Vision (ECCV)*, pages 418–434, 2018. 6
- [25] Xuecheng Nie, Jiashi Feng, and Shuicheng Yan. Mutual learning to adapt for joint human parsing and pose estimation. In *European Conference on Computer Vision (ECCV)*, pages 502–517, 2018. 6
- [26] Scott Reed, Honglak Lee, Dragomir Anguelov, Christian Szegedy, Dumitru Erhan, and Andrew Rabinovich. Training deep neural networks on noisy labels with bootstrapping. *arXiv preprint arXiv:1412.6596*, 2014. 2
- [27] Tao Ruan, Ting Liu, Zilong Huang, Yunchao Wei, Shikui Wei, and Yao Zhao. Devil in the details: Towards accurate

- single and multiple human parsing. In *AAAI Conference on Artificial Intelligence (AAAI)*, volume 33, pages 4814–4821, 2019. [1](#), [2](#), [6](#), [7](#), [8](#)
- [28] Mark Sandler, Andrew Howard, Menglong Zhu, Andrey Zhmoginov, and Liang-Chieh Chen. Mobilenetv2: Inverted residuals and linear bottlenecks. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 4510–4520, 2018. [7](#)
- [29] Antti Tarvainen and Harri Valpola. Mean teachers are better role models: Weight-averaged consistency targets improve semi-supervised deep learning results. In *Advances in Neural Information Processing Systems (NeurIPS)*, pages 1195–1204, 2017. [2](#)
- [30] Fangting Xia, Peng Wang, Liang-Chieh Chen, and Alan L Yuille. Zoom better to see clearer: Human part segmentation with auto zoom net. In *European Conference on Computer Vision (ECCV)*, pages 648–663, 2016. [6](#)
- [31] Fangting Xia, Peng Wang, Xianjie Chen, and Alan L Yuille. Joint multi-person pose estimation and semantic part segmentation. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 6769–6778, 2017. [2](#), [6](#)
- [32] Yuhui Yuan and Jingdong Wang. Ocnet: Object context network for scene parsing. *arXiv preprint arXiv:1809.00916*, 2018. [3](#), [7](#)
- [33] Hengshuang Zhao, Jianping Shi, Xiaojuan Qi, Xiaogang Wang, and Jiaya Jia. Pyramid scene parsing network. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2881–2890, 2017. [3](#), [6](#), [7](#)