



LONG BEACH  
CALIFORNIA  
June 16-20, 2019

# Team ByteDance-SEU



# 1st place

in the **Single-person Human Pose Estimation** Track of CVPR-19 LIP Challenge

Speaker: Kai Su

Kai Su<sup>\*,1,2</sup>, Dongdong Yu<sup>\*,1</sup>, Xin Geng<sup>2</sup>, Changhu Wang<sup>1</sup>

<sup>1</sup>ByteDance AI Lab <sup>2</sup>Southeast University

(\* Equal Contribution)

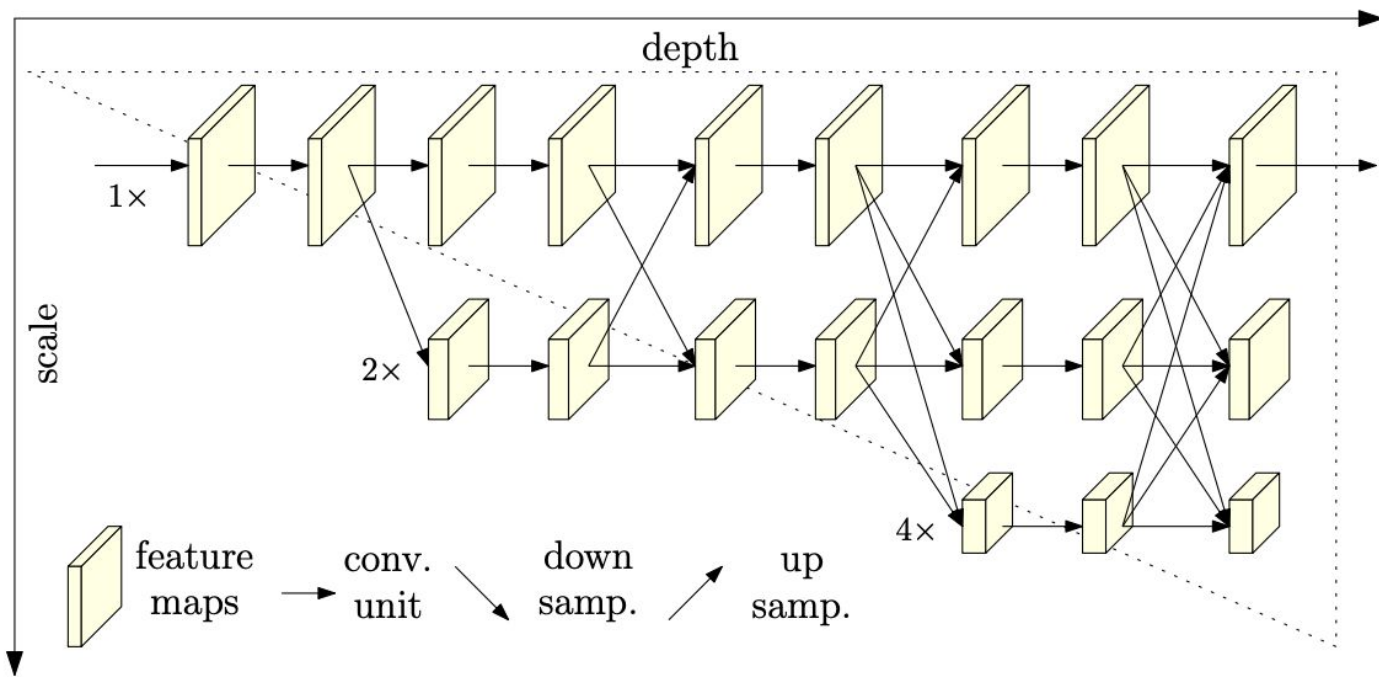


# Outline

- Human Pose Estimation Networks
- Datasets
- Experiments
- Summary



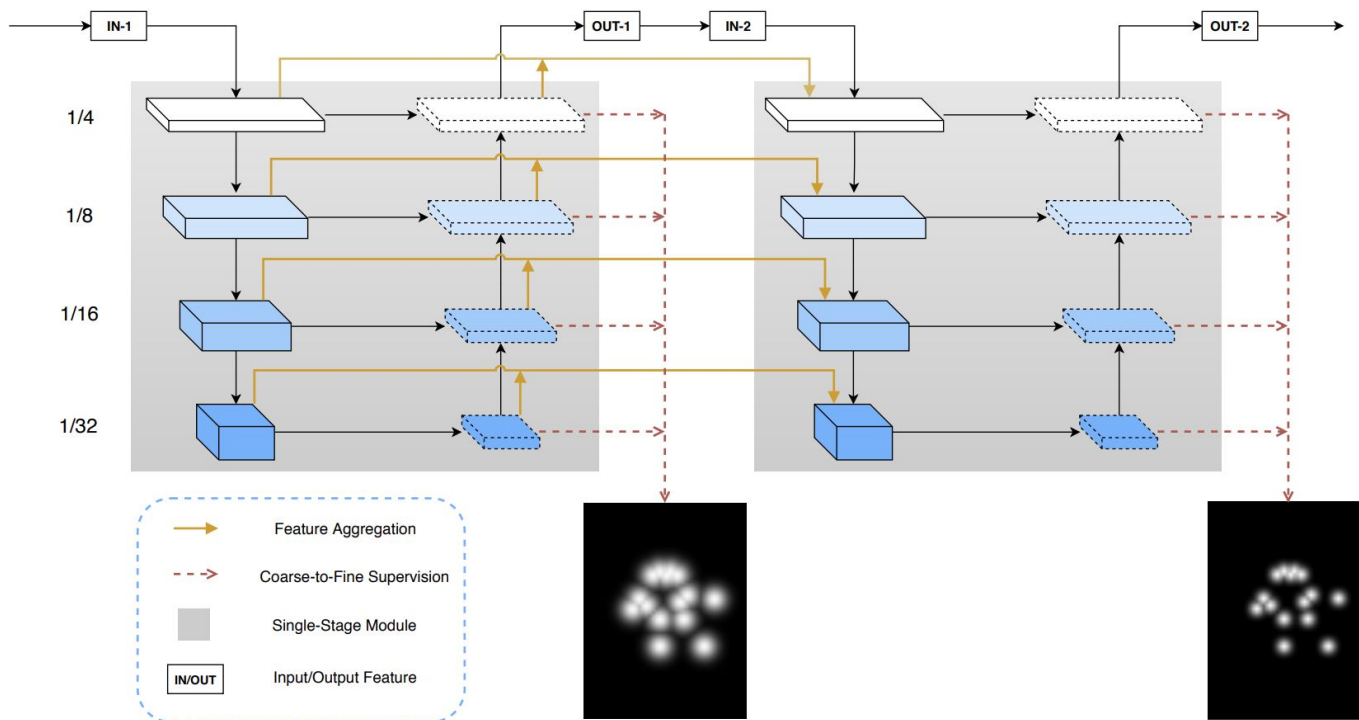
# High-Resolution Net (HRNet)



Sun, Ke, et al. "Deep High-Resolution Representation Learning for Human Pose Estimation." arXiv preprint arXiv:1902.09212 (2019).



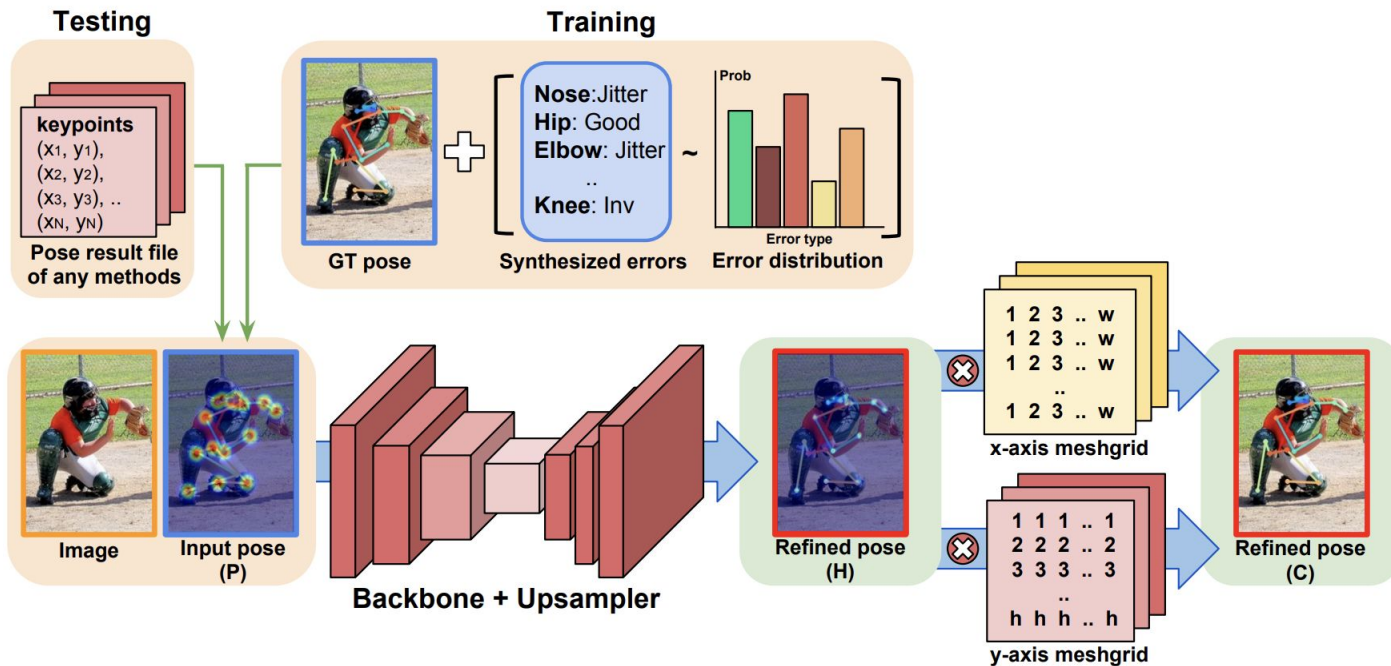
# Multi-Stage Networks (MSPN)



Li, Wenbo, et al. "Rethinking on Multi-Stage Networks for Human Pose Estimation." *arXiv preprint arXiv:1901.00148* (2019).



# Human Pose Refinement Network (PoseFix)



# Outline

- Human Pose Estimation
- Datasets
- Experiments
- Summary



# Datasets

Dataset	Number of images	Number of Keypoints
LIP	3w+ for training, 1w for validation, 1w for test	16
COCO	14w+ for training	17
AI Challenge	20w+ for training	14





# Outline

- Human Pose Estimation
- Datasets
- Experiments
- Summary



# Baseline

Method	Input Size	Training Datasets	PCKh on LIP test set
HRNet-w32	384*288	LIP train	<b>88.2</b>



# External Dataset

Method (HRNet-w32)	Input Size	PCKh on LIP test set
with LIP train (baseline)	384*288	88.2
+ COCO	384*288	91.2 <span style="color: red;">↑3.0</span>
+ AI Challenge	384*288	91.8 <span style="color: red;">↑0.6</span>
+ LIP validation	384*288	<b>91.8</b> <span style="color: blue;">↑0.0</span>



# External Dataset

Method (MSPN)	Input Size	PCKh on LIP test set
with COCO + LIP train (baseline)	384*288	91.4
+ AI Challenge	384*288	91.8 <span style="color: red;">↑0.4</span>
+ LIP validation	384*288	<b>91.9</b> <span style="color: red;">↑0.1</span>



# Input Size

Method (HRNet-32)	Training Datasets	PCKh on LIP test set
384*288	All Available Datasets	91.8
+ 384*288 $\rightarrow$ 512*384	All Available Datasets	<b>92.0</b> $\uparrow$ 0.2

Method (MSPN)	Training Datasets	PCKh on LIP test set
384*288	All Available Datasets	91.9
+ 384*288 $\rightarrow$ 512*384	All Available Datasets	<b>92.0</b> $\uparrow$ 0.1



# Ensemble

Method	Training Datasets	PCKh on LIP test set
HRNet-w32	All Available Datasets	92.0
+ MSPN	All Available Datasets	92.4 <span style="color: red;">↑0.4</span>
+ PoseFix	LIP train/val + COCO	<b>92.6</b> <span style="color: red;">↑0.2</span>

# Outline

- Human Parsing Networks
- Datasets
- Experiments
- Summary



# Summary

Baseline Result: 88.2

- + external dataset 88.2 -> 91.8 (+3.6)
- + enlarge the input size 91.8 -> 92.0 (+0.2)
- + ensemble with average heatmaps 92.0 -> 92.6 (+0.6)



# Bad Cases



# 2nd place

in the **Single-person Human Parsing** Track of CVPR-19 LIP Challenge

Speaker: Dongdong Yu

Dongdong Yu<sup>\*,1</sup>, Kai Su<sup>\*,1,2</sup>, Jian Wang<sup>1</sup>, Kaihui Zhou<sup>1</sup>, Xin Geng<sup>2</sup>, Changhu Wang<sup>1</sup>

<sup>1</sup>ByteDance AI Lab <sup>2</sup>Southeast University

(\* Equal Contribution)

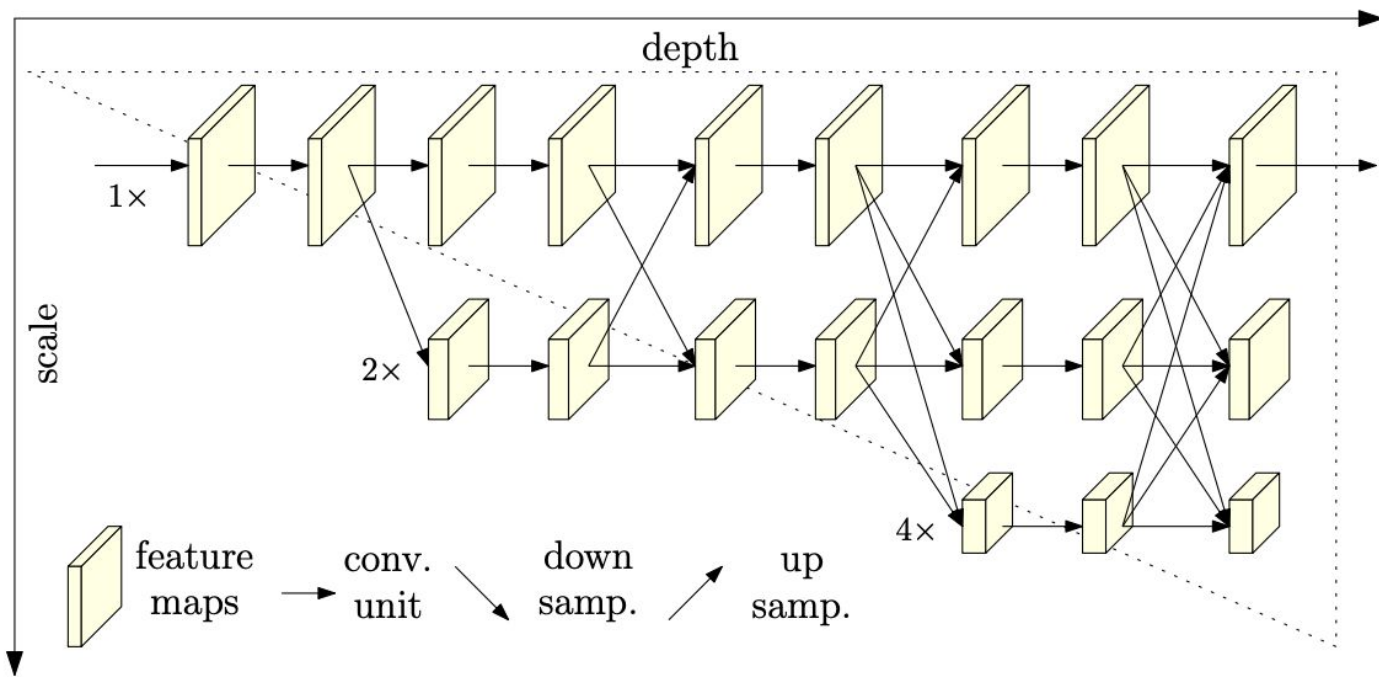


# Outline

- Human Parsing Networks
- Datasets
- Experiments
- Summary



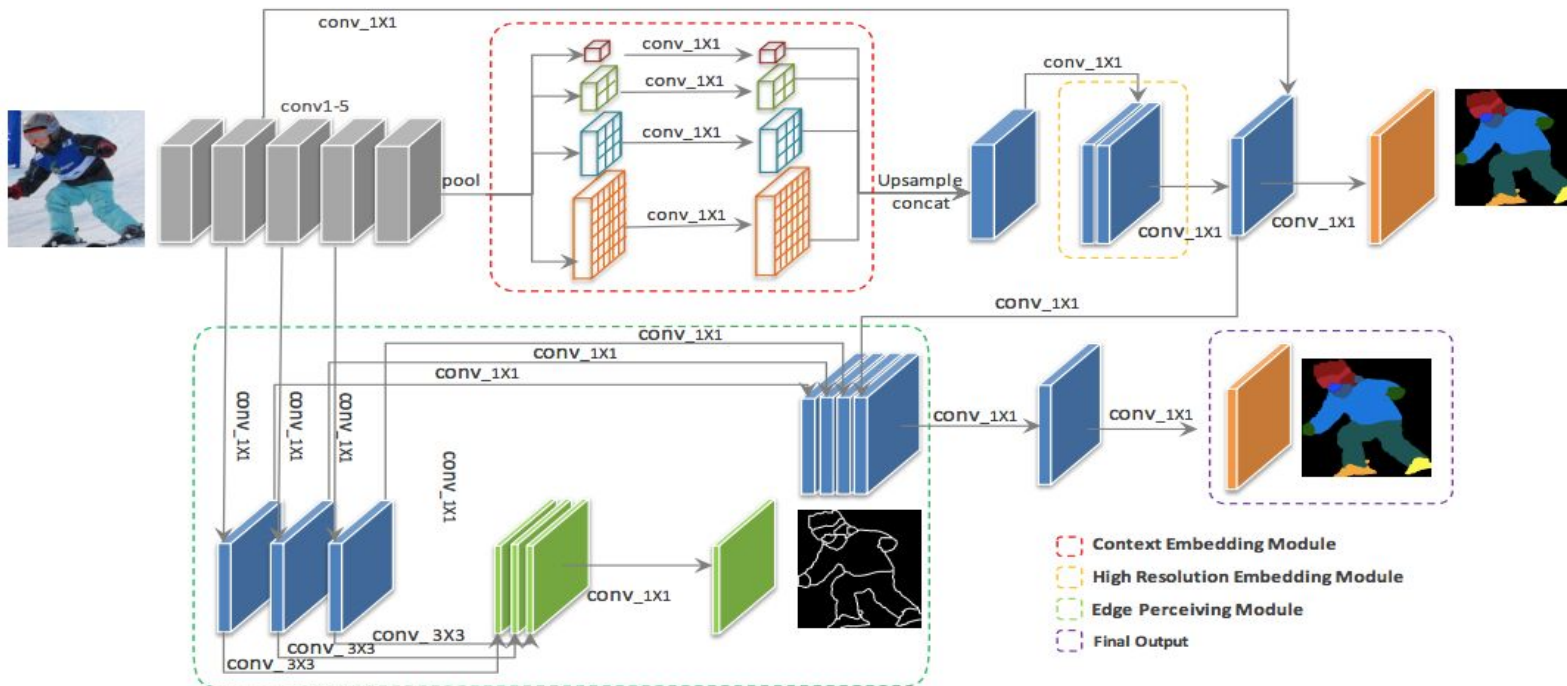
# High-Resolution Net (HRNet)



Sun, Ke, et al. "High-Resolution Representations for Labeling Pixels and Regions." arXiv preprint arXiv:1904.04514 (2019).

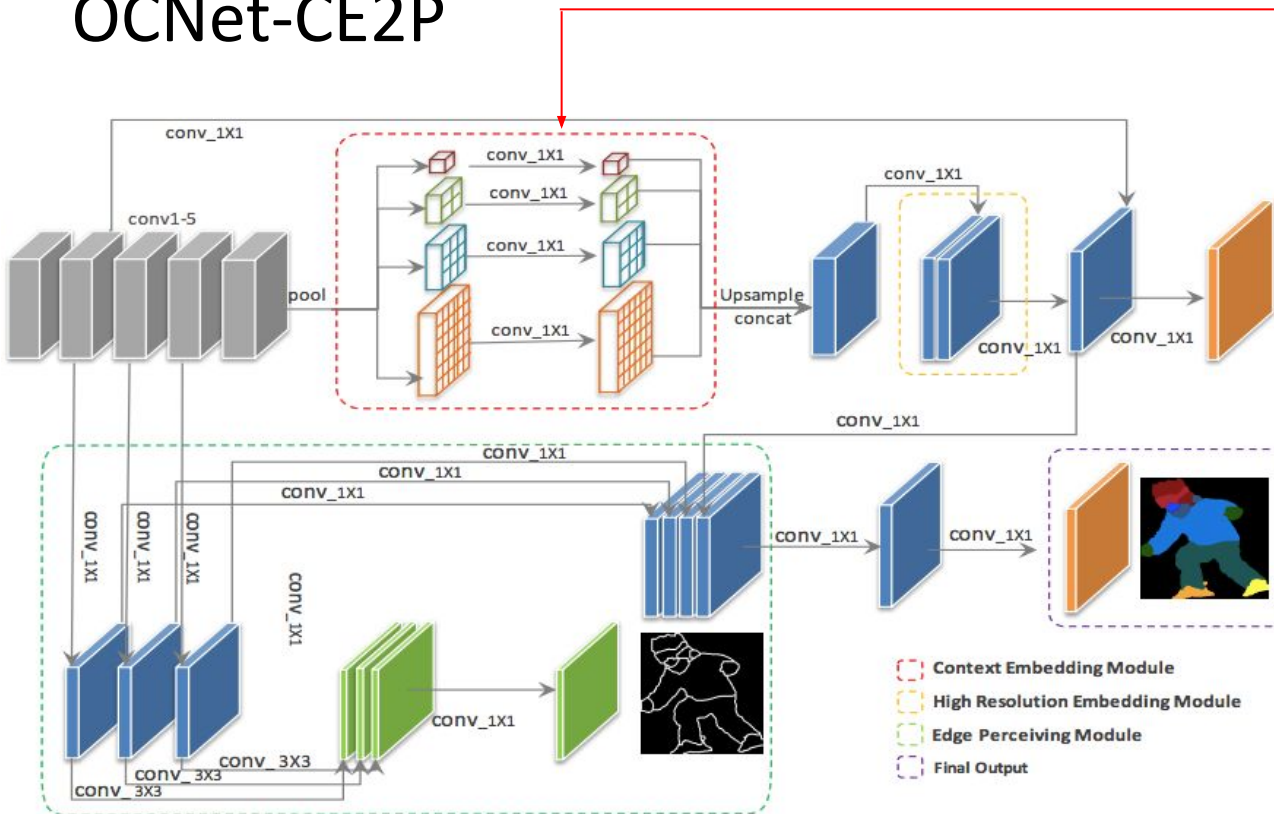


# Context Embedding with Edge Perceiving (CE2P)

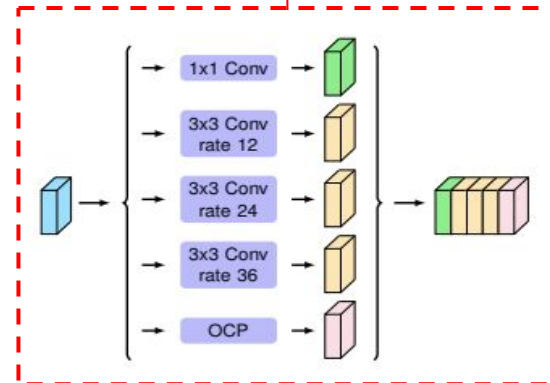




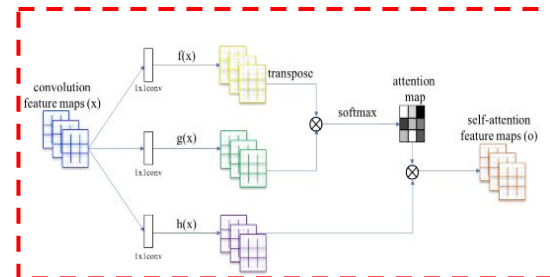
# OCNet-CE2P



## OCP ASPP:



## OCP Module:



Yuan, Yuhui, et al. "OCNet: Object Context Network for Scene Parsing." arXiv preprint arXiv:1809.00916 (2018).



# Outline

- Human Parsing Networks
- Datasets
- Experiments
- Summary



# Datasets

Dataset	Number of instances	Class Definition
LIP	3w+ for training 1w for validation 1w for test	Background(0), Hat(1), Hair(2), Glove(3), Sunglasses(4), Upper-clothes(5), Dress(6),Coat(7),Socks(8),Pants(9), <b>Jumpsuits(10)</b> , Scarf(11), Skirt(12), Face(13), Left-arm(14), Right-arm(15), Left-leg(16), Right-leg(17), Left-shoe(18), Right-shoe(19)
Multi-Person Human Parsing	8W+ human instance	Background(0), Hat(1), Hair(2), Glove(3), Sunglasses(4), Upper-clothes(5), Dress(6),Coat(7),Socks(8),Pants(9), <b>tosor-skin(10)</b> , Scarf(11), Skirt(12), Face(13), Left-arm(14), Right-arm(15), Left-leg(16), Right-leg(17), Left-shoe(18), Right-shoe(19)





# Outline

- Human Parsing Networks
- Datasets
- Experiments
- Summary



# Baseline

Method	Input Size	Training Datasets	Miou on LIP test set
HRNetv2-w48	480*480	LIP training	<b>55.64</b>
CE2P-OCNet (ResNet101)	480*480	LIP training	54.20



# External Dataset

Training Datasets	Method	Input Size	Miou on LIP test set
LIP training	CE2P-OCNet (ResNet101)	480*480	54.20
LIP training + val	CE2P-OCNet (ResNet101)	480*480	56.59 <span style="color: red;">↑2.39</span>
LIP training + val Multi-Human	CE2P-OCNet (ResNet101)	480*480	<b>60.45</b> <span style="color: red;">↑6.25</span>



# Input Size

Input Size	Training Datasets	Method	Miou on LIP test set
480*480	LIP training + val Multi-Human	CE2P-OCNet (ResNet101)	60.45
576*576	LIP training + val Multi-Human	CE2P-OCNet (ResNet101)	<b>60.65</b> <span style="color: red;">↑0.20</span>



# Different Backbone

Method	Training Datasets	Input Size	Miou on LIP test set
CE2P-OCNet (ResNet101)	LIP training + val Multi-Human	576*576	60.65
CE2P-OCNet (SENet152)	LIP training + val Multi-Human	576*576	61.00 <span style="color: red;">↑0.35</span>
CE2P-OCNet (SCNet101)	LIP training + val Multi-Human	576*576	61.23 <span style="color: red;">↑0.58</span>
CE2P-OCNet (DeResNet101)	LIP training + val Multi-Human	576*576	<b>62.17</b> <span style="color: red;">↑1.52</span>
HRNetv2-w48	LIP training + val Multi-Human	576*576	61.96 <span style="color: red;">↑1.31</span>



# Ensemble

Method	Training Datasets	Input Size	Miou on LIP test set
CE2P-OCNet (DeResNet101)	LIP training + val Multi-Human	576*576	62.17
Ensemble with Average Heatmaps	LIP training + val Multi-Human	576*576	64.00 <span style="color: red;">↑1.83</span>
Ensemble with Average Softmax-Heat maps	LIP training + val Multi-Human	576*576	<b>64.13</b> <span style="color: red;">↑1.96</span>



# Outline

- Human Parsing Networks
- Datasets
- Experiments
- Summary



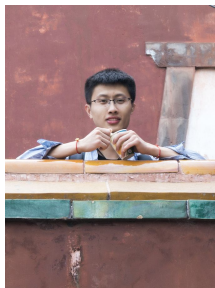
# Summary

Baseline Result: 54.20

- + external dataset 54.20 -> 60.45 (+6.25)
- + enlarge the input size 60.45 -> 60.65 (+0.20)
- + strong backbone 60.65 -> 62.17 (+1.52)
- + ensemble with average heatmaps 62.17 -> 64.00 (+1.83)
- + ensemble with average softmax-heatmaps 64.00 -> 64.13 (+0.13)



# Our Team



Kai Su



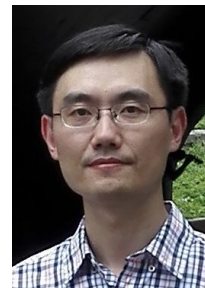
Dongdong Yu



Jian Wang



Kaihui Zhou



Xin Geng



Changhu Wang

# Thanks & Questions

Slides are available at

<https://7color94.github.io/files/CVPR-19-LIP.pdf>

