

# Human Pose Estimation with Spatial Contextual Information

Hong Zhang<sup>1</sup> Hao Ouyang<sup>2</sup> Shu Liu<sup>3</sup> Xiaojuan Qi<sup>5</sup>  
Xiaoyong Shen<sup>3</sup> Ruigang Yang<sup>1,4</sup> Jiaya Jia<sup>3,5</sup>

<sup>1</sup>Baidu Research, Baidu Inc. <sup>2</sup>Hong Kong University of Science and Technology  
<sup>3</sup>YouTu Lab, Tencent <sup>4</sup>University of Kentucky <sup>5</sup>The Chinese University of Hong Kong

{fykalviny, ououkenneth, liushuhust, qxj0125, goodshenxy}@gmail.com  
yangruigang@baidu.com leojia@cse.cuhk.edu.hk

## Abstract

*We explore the importance of spatial contextual information in human pose estimation. Most state-of-the-art pose networks are trained in a multi-stage manner and produce several auxiliary predictions for deep supervision. With this principle, we present two conceptually simple and yet computational efficient modules, namely Cascade Prediction Fusion (CPF) and Pose Graph Neural Network (PGNN), to exploit underlying contextual information. Cascade prediction fusion accumulates prediction maps from previous stages to extract informative signals. The resulting maps also function as a prior to guide prediction at following stages. To promote spatial correlation among joints, our PGNN learns a structured representation of human pose as a graph. Direct message passing between different joints is enabled and spatial relation is captured. These two modules require very limited computational complexity. Experimental results demonstrate that our method consistently outperforms previous methods on MPII and LSP benchmark.*

## 1. Introduction

Human pose estimation refers to the problem of determining precise pixel location of important keypoints of human body. It serves as a fundamental tool to solve other high level tasks, such as human action recognition [48, 28], tracking [9, 51] and human-computer interaction [41]. There are already a variety of solutions where remaining challenges include large change in appearance, uncommon body postures and occlusion.

Recent successful human pose estimation methods are based on Convolutional Neural Networks (CNNs). State-of-the-art methods [35, 50, 52] train pose networks in a multi-stage fashion. These networks produce several auxiliary prediction maps. Then the predictions are refined iteratively in different stages until the final result is produced. It

needs to learn semantically strong appearance features and prevent gradient vanishing during training.

Spatial contextual correlation among different joints plays an important role in human pose estimation [12, 45]. Fig. 1 shows the prediction maps of different stages. Rough locations of *head* and *left knee* are easy to identify in the first stage. However, joints like {*right knee*, *left ankle*} in Fig. 1 are with large deformation and occlusion, which are hard to determine only based on the local regions. Fortunately, location of joints like *left knee* is associated with *left ankle*. So the prediction result of *left knee* in the first stage could be indicated as a prior to help infer the location of *left ankle* in the following stage.

Moreover, since human pose estimation is related to structure, it is important to design appropriate guideline to choose directions of information propagation for the joints that are unclear or occluded. Probabilistic Graphical Models (PGMs) are used to facilitate message passing among joints. In [46, 12], MRF or CRF is utilized to describe the distribution of human body. Nevertheless, the status of each joint needs to be sequentially updated, which means before updating the status of current joints, status of the previous joints is to be refreshed. The sequential nature of the updating scheme makes it easy to accumulate error. The multilevel compositional models [56, 43] considered the relations of joints. These methods all rely on hierarchy structures. Pose grammars are based on the prior knowledge of the human body.

To make good use of the underlying spatial contextual information, we propose two conceptually simple and computational efficient modules to estimate body joints.

**Our Contribution #1** To utilize the contextual information, we propose *Cascade Prediction Fusion* (CPF) to make use of auxiliary prediction maps. The prediction maps at previous stage could be deemed as a prior to support predictions in following stages. This procedure is different from that of [35, 50], where prediction maps were concatenated

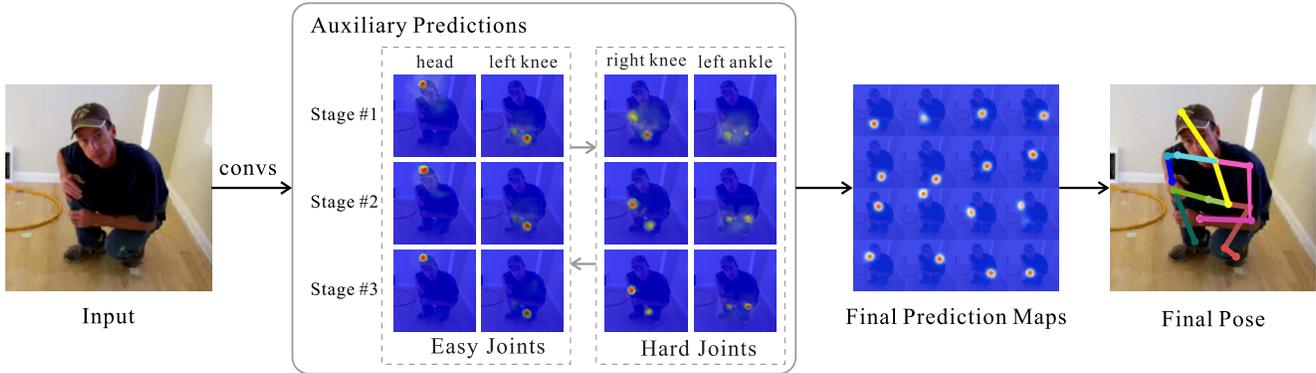


Figure 1. Pipeline of multi-stage prediction. A set of auxiliary predictions are generated. In the first stage, it is easy to identify easy joints while others with severe deformation are still confusing. Relative positions between joints help resolve ambiguity in the second stage. All joints converge to the final prediction in the third stage.

with or added to image feature maps and then fed to following huge CNN trunks. As shown in Fig. 2, we create a light-weight path to gradually accumulate auxiliary prediction maps for final accurate pose estimation. The predictions at different stages are with varied properties. Specifically, predictions from lower layers are with more accurate localization signals while those at higher layers are with stronger semantic information to distinguish among similar keypoints. Our network effectively fuses information from different stages by the shorter path created by CPF.

**Our Contribution #2** We introduce the *Pose Graph Neural Network* (PGNN), which is flexible and efficient to learn a structured representation of body joints. Our PGNN is built on a graph that can be integrated in various pose estimation networks. Each node in the graph is associated with neighboring body parts. Spatial relations are thus captured through edge construction. Direct message passing between different nodes is enabled for precise prediction.

Ours is different in modeling spatial relation. PGNN is a novel way to adaptively select the message passing directions in parallel. Instead of defining an explicit sequential order for a human body structure, it dynamically arranges the update sequences. Via simultaneous update, we manage the short- and long-term relation. Finally, PGNN learns a structured graph representation to boost performance.

Our system is also end-to-end trainable, which not only estimates body location but also configures the spatial structures. We evaluate the system on two representative human pose benchmark datasets, *i.e.*, MPII and LSP. It accomplishes new state-of-the-arts with high computational efficiency.

## 2. Related Work

**Human Pose Estimation** The key of human pose estimation lies in joint detection and spatial relation configuration. Previous human pose estimation methods can be divided

into two groups. The first is to learn feature representation using powerful CNN. These methods detect body joint location directly or predict the score maps for body joints. Early methods like DeepPose [47] regressed joint locations with multiple stages. Later, Fan *et al.* [17] combined local and global features to improve performance. To connect the input and output space, Carreira *et al.* [5] iteratively concatenated the input image with previous prediction in each step. Following the paradigm of semantic segmentation [33, 6], methods of [35, 50, 52] used Gaussian peaks to represent part locations. Then a fully convolutional neural network [33] is applied to estimate body joint location. These methods can produce high quality representation and do not predict structure among body joints, however.

The other group focuses on modeling spatial relationship between body joints. The pictorial structures [36] modeled spatial deformation by designing pairwise terms between different joints. To deal with human poses with large variation, a mixture model is learned for each joint. Yang *et al.* [53] used a part mixture model to infer spatial relation with a tree structure. This structure may not capture very complicated relation. Subsequent methods introduced other models, such as loopy structure [49] and poselet [36] to further improve the performance.

Later methods [45, 46] modeled structures via CNN. Tompson *et al.* [46] utilized the Markov Random Field (MRF) to model distribution of body parts. Convolutional priors were used in defining the pairwise terms of joints. The method of [11] utilized geometrical transform kernels to capture relation of joints on feature maps.

**Graph Neural Network** Previous work on feature learning for graph-structure can be divided into two categories. One direction is to apply CNN to graphs. Based on graph Laplacian, methods of [3, 15, 25] applied CNN to spectral domain. In order to operate CNN directly on graph, the method of [16] used a special hash function. The other

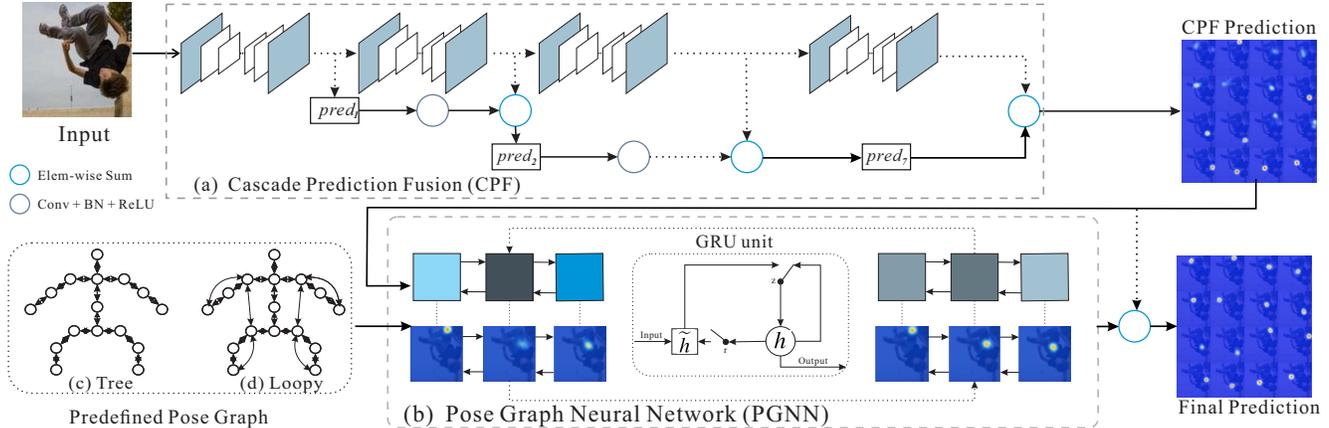


Figure 2. **Framework.** Our system takes an image as input, and generates the prediction maps. The architecture is with two components where CPF is for computing the prediction maps and the other PGNN is for refining these maps until final prediction.

line focuses on recurrently applying neural networks to each node of the graph. State of each node can be updated based on history and new information passing through edges. The Graph Neural Network (GNN) was first proposed in [40]. It utilized multi-layer perceptrons (MLP) to learn hidden state of nodes in graphs. However, the contraction map assumption is a restriction. As an extension, the Gated Graph Neural Network (GGNN) [27] adopted recurrent gating function [8] to update the hidden state and output sequences. Parameters of the final model can be effectively optimized by the back-propagation through time (BPTT) algorithm. Very recently, GGNNs was used in image classification [34], situation recognition [26] and RGBD semantic segmentation [38].

### 3. Our Method

In this section, we describe the two major components in our method. One is a cascaded multi-stage prediction module where previous-stage prediction serves as a prior to guide present prediction and accumulate auxiliary prediction as shown in Fig. 2(a). The other is to model different parts in a graph, augmented by Pose Graph Neural Network (PGNN) to learn representation, as shown in Fig. 2(b).

#### 3.1. Cascade Prediction Fusion (CPF)

For common pose estimation methods [35, 50], a set of prediction maps are iteratively refined for body parts. We propose CPF to take the underlying contextual information encoded in auxiliary prediction into consideration.

These prediction maps are in different semantic levels while all of them can be utilized for final predictions. As detailed in [55, 13], the lower-layer features focus on local appearance and describe details. It is crucial for accurate joints localization. Meanwhile, the global representations from higher layers help discriminate among different

body joints. Our CPF is designed to gradually integrate different semantic information from lower to higher layers. Fig. 2(a) shows the way to incorporate CPF into Hourglass [35] framework. It can be built on top of most multi-stage pose estimation frameworks and iterates from the first stage to the final predictions.

For stage  $i$ , instead of simply fusing the prediction map  $pred_{i-1}$  from last stage with  $pred_i$  from current stage directly, we provide  $pred_{i-1}$  as a prior, which is used for producing  $pred_i$ . Particularly, the coarse prediction map  $pred_{i-1}$  undergoes a  $1 \times 1$  convolution to increase channels and is then merged with image features from stage  $i$  by using element-wise addition.  $pred_i$  is generated by taking the fused feature map as input.

CPF is different from DenseNet [22] and DLA [54]. DenseNet emphasizes more on feature reuse and gradient vanish issues. DLA unifies semantic and spatial fusion in the feature level. In contrast, CPF focuses on exploring and aggregating the contextual information encoded in prediction maps.

#### 3.2. Graph Neural Network (GNN)

Graph neural network (GNN) is a general model handling graph structured data. GNN takes the graph  $G = \{K, E\}$  as input where  $K$  and  $E$  represent the nodes and edges of the graph respectively. Each node  $k \in K$  is associated with a hidden state vector  $h_k$ , which is recurrently updated. The hidden state vector at time step  $t$  is denoted as  $h_k^t$ . The hidden state is updated by taking as input its current state vector and the incoming messages  $x_k^t$  from its neighboring nodes  $\mathcal{N}_k$ .  $\mathcal{A}$  is a function to collect messages from neighboring nodes.  $\mathcal{T}$  is a function to update the hid-

den state. Formally, the hidden state is updated as

$$\begin{aligned} x_k^t &= \mathcal{A}(h_u^{t-1} | u \in \mathcal{N}_k), \\ h_k^t &= \mathcal{T}(h_k^{t-1}, x_k^t). \end{aligned} \quad (1)$$

In the following, we present our new GNN named PGNN for pose estimation.

**Graph Construction** Each node  $k$  in PGNN represents one body joint and each edge is defined as the connection between neighboring joints. Fig. 2(c) shows an example of how to construct a tree-like graph for human poses. The prediction maps are treated as unary maps learned from a backbone network, which will be detailed in Sec. 4. The hidden state of each node is initialized with its corresponding spatial prediction feature map derived from the original image. The status of node  $k$  is initialized as

$$h_k^0 = \mathcal{F}_k(\Theta, I), k \in \{1 \dots K\}, \quad (2)$$

where  $\mathcal{F}$  indicates the backbone network,  $\Theta$  is a set of parameters for the network, and  $I$  is the original input image.

**Information Propagation** We use the constructed graph to exploit the semantic spatial relation and refine the appearance representation for each joint in steps. Before updating the hidden state of each node, it first aggregates messages of the hidden state at time step  $t - 1$  from neighboring node  $k'$ . As demonstrated in [11], convolutional layers can be used as geometrical transform kernels. It advances message passing between feature maps. It is noted that the weights of convolution for different edges are not shared. So  $\mathcal{A}$  is expressed as

$$x_k^t = \sum_{k, k' \in \Omega} W_{p,k} h_{k'}^{t-1} + b_{p,k}, \quad (3)$$

where  $W_{p,k}$  is the convolution weights and  $b_{p,k}$  is the bias of the  $k^{th}$  node.  $\Omega$  is a set of connected edges.

Eq. (4) gives the formulation of  $\mathcal{T}$ . It updates the  $k^{th}$  node with the aggregated messages and the  $t - 1$  step of hidden state. We follow the same gating mechanism with GRU [27] and enjoy more computational efficiency and less memory consumption. Again, we utilize convolution operations and do not share weights.  $W_{z,k}$ ,  $U_{z,k}$ ,  $b_{z,k}$ ,  $W_{r,k}$ ,  $U_{r,k}$ ,  $b_{r,k}$ ,  $W_{h,k}$ ,  $U_{h,k}$  and  $b_{h,k}$  are the weights and biases for the  $k^{th}$  node in the update function. With this method, the aggregated information is softly combined with its own memory, which can be expressed as

$$\begin{aligned} z_k^t &= \sigma(W_{z,k} x_k^t + U_{z,k} h_k^{t-1} + b_{z,k}), \\ r_k^t &= \sigma(W_{r,k} x_k^t + U_{r,k} h_k^{t-1} + b_{r,k}), \\ \tilde{h}_k^t &= \tanh(W_{h,k} x_k^t + U_{h,k} (r_k^t \odot h_k^{t-1}) + b_{h,k}), \\ h_k^t &= (1 - z_k^t) \odot h_k^{t-1} + z_k^t \odot \tilde{h}_k^t. \end{aligned} \quad (4)$$

**Output and Learning** After  $T$ -time propagation, we get the final prediction

$$\tilde{P}_k = h_k^T + h_k^0, \quad (5)$$

where  $h_k^T$  is the final hidden state collected from the corresponding node.  $h_k^0$  is the initialization hidden state, which encodes the appearance information of a joint. We get the final prediction by adding these two prediction maps. The graph network is trained by minimizing the  $\ell_2$  loss of

$$L_2 = \frac{1}{K} \sum_{k=1}^K \sum_{x,y} \|\tilde{P}_k(x,y) - P_k(x,y)\|^2, \quad (6)$$

where  $(x, y)$  is the pixel location,  $P_k(x, y)$  is the ground truth label at pixel  $(x, y)$ .  $\tilde{P}_k$  is the prediction map obtained in Eq. (5). The model is trained with back-propagation through time (BPTT).

### 3.2.1 Graph Types

PGNN can handle a variety of graphs. It develops a novel message passing scheme so that each body receives message from specific neighboring joints. Intuitively, a fully connected graph is expected to be the ideal choice to collect information from all other joints. However, for some joints, such as *head* and *ankle*, it is hard to capture the relationship.

To address this problem, we utilize two types of structure, *i.e.*, tree and loopy structure. It is not known beforehand which one is better. A tree is a simple structure, which captures the relation of neighboring joints. Loopy structure is more complex, allowing message passing in a loop. The structure we use in this paper is illustrated in Fig. 2(c) and (d). Although many tree-like or loopy graphs can be derived, PGNN tackles them in the same way. The graphs are undirected and enable bidirectional message passing.

### 3.2.2 Relationship to Other Methods

Most current state-of-the-art methods focus more on appearance learning of body parts. They capture spatial relation by enlarging the receptive fields. However, poses are with large variation, making the structure information in prediction feature maps still have the potential to boost the performance. Other models like Recurrent Neural Network (RNN) and Probabilistic Graphical Model (PGM) can also model the relation. We will detail the difference between PGNN and these models in the following.

**PGNN vs. RNN** RNN can be deemed as a special case of PGNN. It is also able to pass information across nodes of a graph, where each body part is denoted as a node and the joints relations are propagated through edges. However, the

graph structure requires to be chains for RNN. For its construction, at each time step, the state of current node in RNN is updated by its current state and the hidden state of the parent node. It is different from our PGNN, which collects information from a set of neighboring nodes. Moreover, the order of RNN input is manually defined. A slightly inappropriate setting may destroy the naturally structured relationship of joints.

Tree-structured RNN [42] can handle tree-structured data, which propagates information through a tree sequentially. In addition, before updating the state of subsequent layers  $L_t$ , it must update the ancestors  $L_{t-1}$  at first. Contrarily, PGNN updates all states of the node simultaneously. In addition, RNN shares weights at different time steps. The transfer matrix between nodes in the graph is shared through  $T$ -time update. Note that in our model, each edge of the graph has different transformation weights.

**PGNN vs. PGM** PGNN is also closely related to probabilistic graphical model, which is widely used for pose estimation [46, 11] to integrate joint associations. In fact, our model can be viewed as generalization of these models by designing specific update. As detailed in [46], for a body part  $r$ , the final marginal likelihood  $\tilde{Q}_r$  is defined as

$$\tilde{Q}_r = \frac{1}{Z} \prod_{v \in \mathcal{V}} (q_{r|v} * q_v + b_{v \rightarrow r}), \quad (7)$$

where  $\mathcal{V}$  is a set of neighboring nodes of  $r$ .  $q_v$  is the joint probability,  $q_{r|v}$  is the conditional prior and  $b_{v \rightarrow r}$  is the bias, respectively.  $Z$  is the partition function. When the aggregation function is formulated as the product and update function is represented by Eq. (7), PGNN is degraded to the MRF model. With these derivations, it becomes clear that PGNN is a more general model to integrate joint associations by designing specific graph structure and making its own way to update and aggregate functions.

### 3.3. Backbone Networks

To verify the generality of our method, we use two backbone networks. One is our modified ResNet-50 [20] and the other is the widely used 8-stack Hourglass [35].

#### 3.3.1 ResNet-50

ResNet has demonstrated its power on many high-level tasks, including object detection [20] and instance segmentation [19, 31, 32]. To show the generalization ability, we first modify the ResNet-50 network with a few novel steps for human pose estimation. It achieves decent results, even comparable with using other much deeper networks.

Our strategy is to first convert the vanilla ResNet-50 into a fully convolutional network by removing the final classification and average pooling. We integrate CPF in ResNet-50

and further improve the results using PGNN. The following two techniques are also used to adapt ResNet-50.

**Feature Pyramid Network (FPN)** As introduced in Sec. 3.1, multi-stage prediction is very important for training a pose network. To this end, we adopt Feature Pyramid Network (FPN) [30] in the vanilla ResNet-50. FPN leverages the pyramid shape in networks for prediction at different feature levels. Similar to FPN, we also use a lateral connection ( $1 \times 1$  conv) to merge the information from both bottom-up and top-down pathways. Finally, we produce three auxiliary predictions at three different levels.

**Dilated Convolution** Dilated Convolution [6] is used to enlarge the receptive field without introducing extra parameters. An input image is down-sampled 32 times after fed into the vanilla ResNet-50. However, the feature map is too coarse to precisely localize the joints. To address this problem, we first decrease the stride of convolution layers of the last two blocks from 2px to 1px. This results in shrinking the receptive field. Since for human pose estimation as demonstrated in [35, 50], the spatial information needs to be captured by a large enough area, we replace the  $3 \times 3$  convolution layers of the last two blocks with the dilated convolution. Finally, we reduce the stride to 8px.

**Other Implementation Details** All models are implemented by Torch [14]. We use ImageNet pre-trained model as the base and adopt RMSProp [44] to optimize parameters. The network is trained in a total of 250 epochs with batch size 8. The initial learning rate is 0.001. It decreased by 10 times at the 200<sup>th</sup> epoch.

#### 3.3.2 Hourglass

The 8-stack Hourglass (Hg) is adopted as the other backbone network to verify our method. It is much deeper than ResNet-50 and is widely adopted by many pose estimation frameworks. With CPF and PGNN integrated in Hourglass, we achieve new state-of-the-art results.

**Implementation Details** The network is implemented using Torch and optimized with RMSProp. The parameters are randomly initialized. We train the network in 300 epochs with batch size 6. The learning rate starts at 0.00025 and decreases by 10 times at the 240<sup>th</sup> epoch.

## 4. Experiments

**Datasets** We evaluate our CPF and PGNN on two representative benchmark datasets, *i.e.*, MPII human pose dataset (MPII) [1] and extended Leeds Sports Poses (LSP) [24]. MPII contains about 25,000 images with

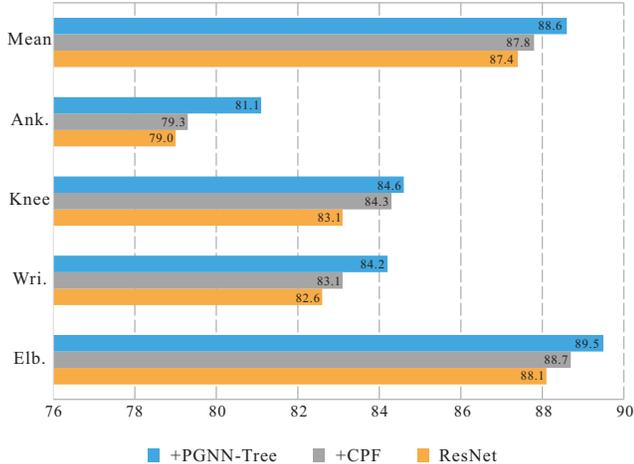


Figure 3. Prediction results of PCKh @0.5 on MPII validation set. We compare the results by adding CPF and further integrating PGNN. The backbone is ResNet-50 and PGNN passes message with a tree-like structure.

over 40,000 annotated poses. We use the same setting as that in [45] to split training and validation sets. The dataset is very challenging since it covers daily human activities with large pose variety. The LSP dataset consists of 11,000 training images and 1,000 testing ones from sport activities.

**Data Augmentation** During training, the input image is cropped and warped to size  $256 \times 256$  according to the annotated body position and scale. We augment the dataset by scaling it with a factor ( $[0.75, 1.25]$ ), rotation ( $\pm 30$ ), horizontal flipping and illumination adjustment to enhance data diversity, and further improve the robustness of the model for various cases. During testing, we crop the image with the given rough center location and scale of the person. For LSP dataset, we simply use the size and center of the image as rough scale and center.

**Unary Maps** For the two backbone networks, *i.e.*, ResNet-50 and Hourglass, we take the last prediction score maps as the unary maps. The reason is that the prediction at the final stage is made based on feature with strong semantics, which gathers previous prediction through CPF. It is generally with decent prediction accuracy. The score maps are of size  $H \times W \times C$  where  $H$  is the height,  $W$  is the width, and  $C$  is the channel size. In our experiments, the size of  $C$  depends on different datasets.  $W$  and  $H$  are all with size 64, which is 1/8 of the original input image.

**Evaluation Criteria** We use the Percentage Correct Key-points (PCK) to evaluate results on the LSP dataset. For MPII, we use PCKh [1], a modified version of PCK. It normalizes the distance errors with respect to the size of head.

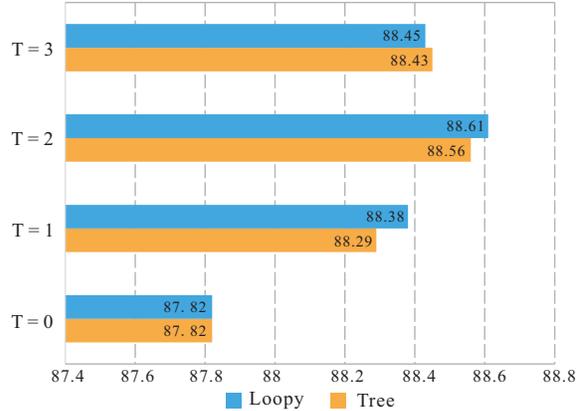


Figure 4. Results at different timesteps with tree-like and loopy-like structure of PCKh @0.5. The backbone is ResNet-50.

#### 4.1. Ablation Study

To investigate the effectiveness of our proposed CPF and PGNN modules, we conduct ablative analysis on the validation set of MPII Human Pose dataset. We set the modified ResNet-50 as our baseline network. To show the efficacy of our models, all results are tested without flipping or multi-scale testing.

**CPF** To evaluate the effectiveness of CPF, we compare results with and without CPF on modified ResNet-50. Fig. 3 shows results on MPII validation set. “ResNet” refers to our modified ResNet-50. For ResNet-50 with CPF, some difficult joints like *knee* is 1.2% higher and result of *elbow* is improved by 0.6%.

In order to clearly demonstrate accuracy change in each stage. Fig. 5(b) shows the accuracy produced at different stages. It is observed that the accuracy increases gradually in steps on the ResNet. This manifests that our CPF effectively gathers information from previous predictions.

**PGNN** Other than adding CPF, we integrate PGNN to further enhance the accuracy. Fig. 3 gives the experimental results where “PGNN” stands for the graph updated twice based on a tree-like structure. The accuracy of parts such as difficult joints of *elbow*, *wrist* is further improved. The reason is that the contextual information propagated from confident parts through graph helps reduce error.

We use two types of graphs in our experiments. They are tree- and loopy-like graphs in PGNN. Fig. 4 presents the results using different PGNN structures. They are comparable – connecting the parts including  $\{elbow, ankle\}$  with other easy parts consistently improves performance. In our experiments, a naive loopy structure, shown in Fig. 2(d), is used. We simply add extra connections, *i.e.* *shoulder-wrist*, *ankle-hip*, and *shoulder-hip*. It is notable that performance of these two types of graphs with the same number of

steps are consistent. We thus believe allowing information to propagate between neighboring joints is of great importance. More sophisticated structures may further improve the performance, which will be our future work.

We also conduct experiments to compare results when propagating different times (*i.e.* with varying propagation number  $T$ ) in the system. The results are shown in Fig. 4. The performance increases by a small amount when increasing  $T$ , and saturates quickly at  $T = 3$ . We also notice that propagation is important in the first 2 steps. For the tree-like graph, as revealed in the comparison when applying  $T = 0$ ,  $T = 1$  and  $T = 2$ , we obtain the improvement of around 0.5% and 0.3%, respectively. Similar results are observed when using the loopy-like graph. However, the performance begins to drop at  $T = 3$ . Since it is hard to capture the semantic information between too far away joints, but instead confuses prediction at current joint.

Methods	Head	Sho.	Elb.	Wri.	Hip	Knee	Ank.	Mean
Pishchulin <i>et al.</i> [36]	74.3	49.0	40.8	34.1	36.5	34.4	35.2	44.1
Tompson <i>et al.</i> [46]	95.8	90.3	80.5	74.3	77.6	69.7	62.8	79.6
Carreira <i>et al.</i> [5]	95.7	91.7	81.7	72.4	82.8	73.2	66.4	81.3
Tompson <i>et al.</i> [45]	96.1	91.9	83.9	77.8	80.9	72.3	64.8	82.0
Hu&Ramanan <i>et al.</i> [21]	95.0	91.6	83.0	76.6	81.9	74.5	69.5	82.4
Pishchulin <i>et al.</i> [37]	94.1	90.2	83.4	77.3	82.6	75.7	68.6	82.4
Lifshitz <i>et al.</i> [29]	97.8	93.3	85.7	80.4	85.3	76.6	70.2	85.0
Gkioxary <i>et al.</i> [18]	96.2	93.1	86.7	82.1	85.2	81.4	74.1	86.1
Rafi <i>et al.</i> [39]	97.2	93.9	86.4	81.3	86.8	80.6	73.4	86.3
Insafutdinov <i>et al.</i> [23]	96.8	95.2	89.3	84.4	88.4	83.4	78.0	88.5
Wei <i>et al.</i> [50]	97.8	95.0	88.7	84.0	88.4	82.8	79.4	88.5
Chu <i>et al.</i> [13]	98.5	96.3	91.9	88.1	90.6	88.0	85.0	91.5
Chou <i>et al.</i> [10]	98.2	96.8	92.2	88.0	91.3	89.1	84.9	91.8
Chen <i>et al.</i> [7]	98.1	96.5	92.5	88.5	90.2	89.6	86.0	91.9
Yang <i>et al.</i> [52]	98.5	96.7	92.5	88.7	91.1	88.6	86.0	92.0
Newell <i>et al.</i> [35]	98.2	96.3	91.2	87.1	90.1	87.4	83.6	90.9
ResNet-ours	98.2	96.4	91.6	87.1	91.2	88.0	83.6	91.2
Hg-ours	<b>98.6</b>	<b>97.0</b>	<b>92.8</b>	<b>88.8</b>	<b>91.7</b>	<b>89.8</b>	<b>86.6</b>	<b>92.5</b>

Table 1. Results of PCKh @0.5 on the MPII test set. Note that ResNet is our modified ResNet-50. ResNet-50 and Hg are all trained with CPF and PGNN.

## 4.2. Experimental Results on MPII

**Accuracy** Tab. 1 lists our results on MPII test set. “Hg-ours” is trained on MPII combined with LSP. The results are produced with five-scale input with horizontal flip testing. Our method trained based on Hourglass yields result 92.5% PCKh at threshold 0.5, which is the highest on this dataset at the time of paper submission. For the challenging parts such as *knee* and *ankle*, we obtain improvement of 2.4% and 3.0% compared to the baseline Hourglass, respectively. Particularly, our method outperforms the method [13] with CRF as well. It is noteworthy that accuracy of our method (ResNet-ours) is also higher than the baseline ResNet-50, which proves the generalization ability.

Methods	Head	Sho.	Elb.	Wri.	Hip	Knee	Ank.	Mean
Belagiannis&Zisserman [2]	95.2	89.0	81.5	77.0	83.7	87.0	82.8	85.2
Lifshitz <i>et al.</i> [29]	96.8	89.0	82.7	79.1	90.9	86.0	82.5	86.7
Pishchulin <i>et al.</i> [37]	97.0	91.0	83.8	78.1	91.0	86.7	82.0	87.1
Insafutdinov <i>et al.</i> [23]	97.4	92.7	87.5	84.4	91.5	89.9	87.2	90.1
Wei <i>et al.</i> [50]	97.8	92.5	87.0	83.9	91.5	90.8	89.9	90.5
Bulat&Tzimiropoulos [4]	97.2	92.1	88.1	85.2	92.2	91.4	88.7	90.7
Yang <i>et al.</i> [52]	98.3	94.5	<b>92.2</b>	88.9	<b>94.4</b>	<b>95.0</b>	93.7	93.9
ResNet-ours	<b>98.5</b>	94.0	89.9	86.9	92.3	93.5	92.7	92.5
Hg-ours	98.4	<b>94.8</b>	92.0	<b>89.4</b>	<b>94.4</b>	94.8	<b>93.8</b>	<b>94.0</b>

Table 2. Comparison of PCK @0.2 on the LSP dataset. ResNet is short for ResNet-50. Both backbones are trained with CPF and PGNN.

**Complexity** In Fig. 5(a), we compare the number of parameters and computational complexity between Hourglass, previous method PRM [52] and our model. We note that the PRM adds 13.5% extra parameters compared with Hourglass, while our model only increases parameters by 0.8%. Additionally, we introduce very limited computation overhead (measured by GFLOPs) on Hourglass, contrary to much increased computational cost from PRM. Our results are with higher quality and decent computational efficiency.

**Visual Analysis** In Fig. 6, we visualize results of baseline and our model. The baseline model has difficulty in distinguishing among symmetric parts and uncommon body postures. For example, in  $\{col.1, row.2\}$  of Fig. 6, *ankle* with large deformation is hard to identify with inherent ambiguity. Our proposed CPF and PGNN provide an effective way to utilize contextual information to reduce the confusion. As a result, the associated joints *knee* help inferring the precise location of *ankle* in our model as shown in  $\{col.2 and row.2\}$  of Fig 6. More results on MPII and LSP generated by our method are shown in Fig. 7.

## 4.3. Experimental Results on LSP

Tab. 2 gives comparison with person-centric annotation. The results are evaluated with PCK scores at threshold 0.2. Following previous methods [13, 52], we add the MPII training set to the extended LSP training set. Our modified ResNet-50 outperforms most of the methods trained with deeper networks.

## 5. Concluding Remarks

We have presented effective *Cascade Prediction Fusion* (CPF) and *Pose Graph Neural Network* (PGNN) to explore contextual information for human pose estimation. CPF makes use of rich contextual information encoded in the auxiliary score maps to produce enhanced prediction. PGNN, differently, is adopted to provide an explicit information propagation scheme to refine prediction. These two components are independent while beneficial to each other

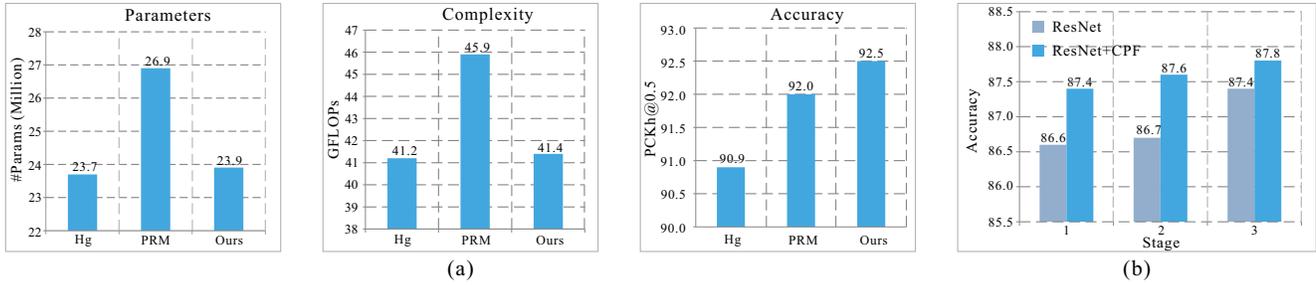


Figure 5. (a) Statistics of parameter numbers, GFLOPs and accuracy on three frameworks, *i.e.*, baseline Hourglass, PRM and our method respectively. (b) Prediction accuracy at different stages with modified ResNet-50.

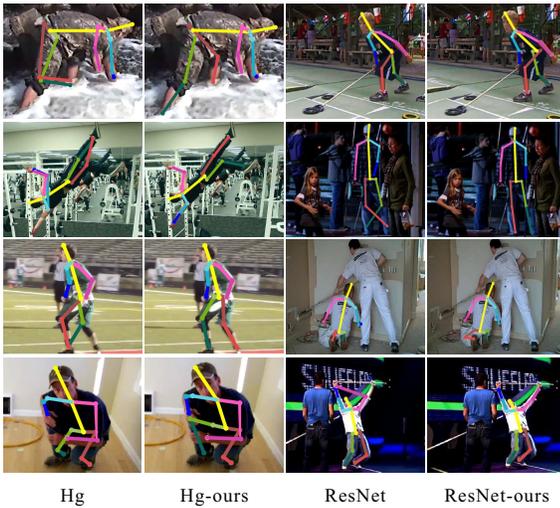


Figure 6. Results on MPII test set produced by different backbone networks, *i.e.* Hourglass and ResNet-50. Hg-ours and ResNet-50-ours are both trained with CPF and PGNN.

in human pose estimation. They are also general for most existing pose estimation networks to boost performance. Our future work will be to extend our framework to 3D and video data for deeper understanding of the temporal and spatial relationship.

## References

- [1] M. Andriluka, L. Pishchulin, P. Gehler, and B. Schiele. 2d human pose estimation: New benchmark and state of the art analysis. In *CVPR*, 2014. 5, 6
- [2] V. Belagiannis and A. Zisserman. Recurrent human pose estimation. In *FG*, 2017. 7
- [3] J. Bruna, W. Zaremba, A. Szlam, and Y. LeCun. Spectral networks and locally connected networks on graphs. *arXiv preprint arXiv:1312.6203*, 2013. 2
- [4] A. Bulat and G. Tzimiropoulos. Human pose estimation via convolutional part heatmap regression. In *ECCV*, 2016. 7
- [5] J. Carreira, P. Agrawal, K. Fragkiadaki, and J. Malik. Human pose estimation with iterative error feedback. In *CVPR*, 2016. 2, 7
- [6] L.-C. Chen, G. Papandreou, I. Kokkinos, K. Murphy, and A. L. Yuille. Deeplab: Semantic image segmentation with deep convolutional nets, atrous convolution, and fully connected crfs. *arXiv preprint arXiv:1606.00915*, 2016. 2, 5
- [7] Y. Chen, C. Shen, X.-S. Wei, L. Liu, and J. Yang. Adversarial posenet: A structure-aware convolutional network for human pose estimation. *arXiv preprint arXiv:1705.00389*, 2017. 7
- [8] K. Cho, B. Van Merriënboer, C. Gulcehre, D. Bahdanau, F. Bougares, H. Schwenk, and Y. Bengio. Learning phrase representations using rnn encoder-decoder for statistical machine translation. *arXiv preprint arXiv:1406.1078*, 2014. 3
- [9] N.-G. Cho, A. L. Yuille, and S.-W. Lee. Adaptive occlusion state estimation for human pose tracking under self-occlusions. *Pattern Recognition*, 46(3):649–661, 2013. 1
- [10] C.-J. Chou, J.-T. Chien, and H.-T. Chen. Self adversarial training for human pose estimation. *arXiv preprint arXiv:1707.02439*, 2017. 7
- [11] X. Chu, W. Ouyang, H. Li, and X. Wang. Structured feature learning for pose estimation. In *CVPR*, 2016. 2, 4, 5
- [12] X. Chu, W. Ouyang, X. Wang, et al. Crf-cnn: Modeling structured information in human pose estimation. In *NIPS*, 2016. 1
- [13] X. Chu, W. Yang, W. Ouyang, C. Ma, A. L. Yuille, and X. Wang. Multi-context attention for human pose estimation. *arXiv preprint arXiv:1702.07432*, 2017. 3, 7
- [14] R. Collobert, K. Kavukcuoglu, and C. Farabet. Torch7: A matlab-like environment for machine learning. In *BigLearn, NIPS Workshop*, 2011. 5
- [15] M. Defferrard, X. Bresson, and P. Vandergheynst. Convolutional neural networks on graphs with fast localized spectral filtering. In *NIPS*, 2016. 2
- [16] D. K. Duvenaud, D. Maclaurin, J. Iparraguirre, R. Bombarell, T. Hirzel, A. Aspuru-Guzik, and R. P. Adams. Convolutional networks on graphs for learning molecular fingerprints. In *NIPS*, 2015. 2
- [17] X. Fan, K. Zheng, Y. Lin, and S. Wang. Combining local appearance and holistic view: Dual-source deep neural networks for human pose estimation. In *CVPR*, 2015. 2
- [18] G. Gkioxari, A. Toshev, and N. Jaitly. Chained predictions using convolutional neural networks. In *ECCV*, 2016. 7
- [19] K. He, G. Gkioxari, P. Dollár, and R. Girshick. Mask r-cnn. In *ICCV*, 2017. 5

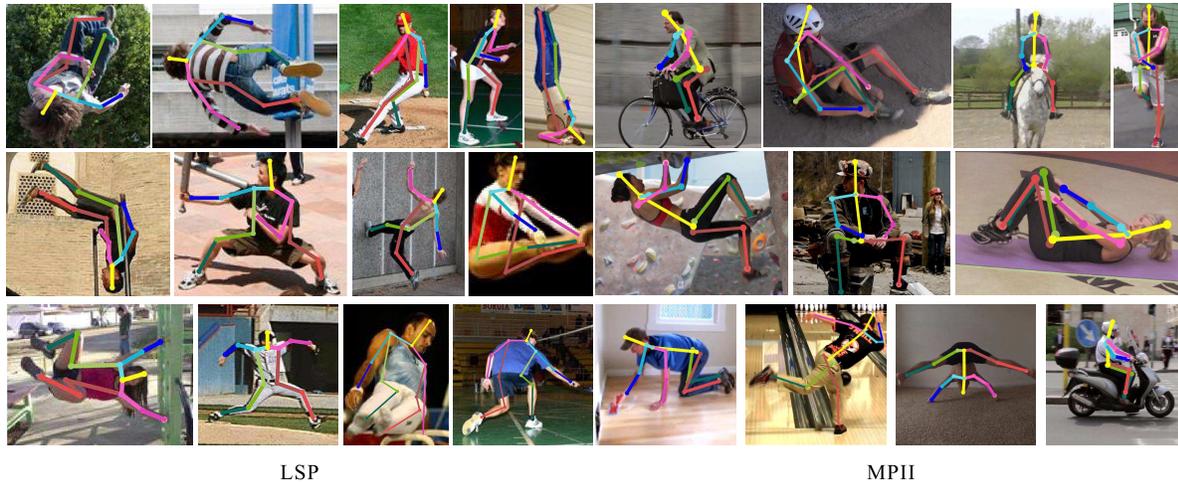


Figure 7. Example output on the LSP and MPII test data.

- [20] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *CVPR*, 2016. 5
- [21] P. Hu and D. Ramanan. Bottom-up and top-down reasoning with hierarchical rectified gaussians. In *CVPR*, 2016. 7
- [22] G. Huang, Z. Liu, K. Q. Weinberger, and L. van der Maaten. Densely connected convolutional networks. In *CVPR*, 2017. 3
- [23] E. Insafutdinov, L. Pishchulin, B. Andres, M. Andriluka, and B. Schiele. Deepcut: A deeper, stronger, and faster multi-person pose estimation model. In *ECCV*, 2016. 7
- [24] S. Johnson and M. Everingham. Clustered pose and non-linear appearance models for human pose estimation. In *BMVC*, 2010. 5
- [25] T. N. Kipf and M. Welling. Semi-supervised classification with graph convolutional networks. *arXiv preprint arXiv:1609.02907*, 2016. 2
- [26] R. Li, M. Tapaswi, R. Liao, J. Jia, R. Urtasun, and S. Fidler. Situation recognition with graph neural networks. In *ICCV*, 2017. 3
- [27] Y. Li, D. Tarlow, M. Brockschmidt, and R. Zemel. Gated graph sequence neural networks. In *ICLR*, 2016. 3, 4
- [28] Z. Liang, X. Wang, R. Huang, and L. Lin. An expressive deep model for human action parsing from a single image. In *ICME*, 2014. 1
- [29] I. Lifshitz, E. Fetaya, and S. Ullman. Human pose estimation using deep consensus voting. In *ECCV*, 2016. 7
- [30] T.-Y. Lin, P. Dollár, R. Girshick, K. He, B. Hariharan, and S. Belongie. Feature pyramid networks for object detection. In *CVPR*, 2017. 5
- [31] S. Liu, J. Jia, S. Fidler, and R. Urtasun. Sgn: Sequential grouping networks for instance segmentation. In *ICCV*, 2017. 5
- [32] S. Liu, L. Qi, H. Qin, J. Shi, and J. Jia. Path aggregation network for instance segmentation. In *CVPR*, 2018. 5
- [33] J. Long, E. Shelhamer, and T. Darrell. Fully convolutional networks for semantic segmentation. In *CVPR*, 2015. 2
- [34] K. Marino, R. Salakhutdinov, and A. Gupta. The more you know: Using knowledge graphs for image classification. In *CVPR*, 2017. 3
- [35] A. Newell, K. Yang, and J. Deng. Stacked hourglass networks for human pose estimation. In *ECCV*, 2016. 1, 2, 3, 5, 7
- [36] L. Pishchulin, M. Andriluka, P. Gehler, and B. Schiele. Poselet conditioned pictorial structures. In *CVPR*, 2013. 2, 7
- [37] L. Pishchulin, E. Insafutdinov, S. Tang, B. Andres, M. Andriluka, P. Gehler, and B. Schiele. Deepcut: Joint subset partition and labeling for multi person pose estimation. In *CVPR*, 2016. 7
- [38] X. Qi, R. Liao, J. Jia, S. Fidler, and R. Urtasun. 3d graph neural networks for rgbd semantic segmentation. In *ICCV*, 2017. 3
- [39] U. Rafi, B. Leibe, J. Gall, and I. Kostrikov. An efficient convolutional network for human pose estimation. In *BMVC*, 2016. 7
- [40] F. Scarselli, M. Gori, A. C. Tsoi, M. Hagenbuchner, and G. Monfardini. The graph neural network model. *IEEE Transactions on Neural Networks*, 20(1):61–80, 2009. 3
- [41] J. Shotton, T. Sharp, A. Kipman, A. Fitzgibbon, M. Finocchio, A. Blake, M. Cook, and R. Moore. Real-time human pose recognition in parts from single depth images. *Communications of the ACM*, 56(1):116–124, 2013. 1
- [42] K. S. Tai, R. Socher, and C. D. Manning. Improved semantic representations from tree-structured long short-term memory networks. *arXiv preprint arXiv:1503.00075*, 2015. 5
- [43] Y. Tian, C. L. Zitnick, and S. G. Narasimhan. Exploring the spatial hierarchy of mixture models for human pose estimation. In *ECCV*, 2012. 1
- [44] T. Tieleman and G. Hinton. Lecture 6.5-rmsprop: Divide the gradient by a running average of its recent magnitude. *COURSERA: Neural networks for machine learning*, 4(2):26–31, 2012. 5
- [45] J. Tompson, R. Goroshin, A. Jain, Y. LeCun, and C. Bregler. Efficient object localization using convolutional networks. In *CVPR*, 2015. 1, 2, 5, 7

- [46] J. J. Tompson, A. Jain, Y. LeCun, and C. Bregler. Joint training of a convolutional network and a graphical model for human pose estimation. In *NIPS*, 2014. [1](#), [2](#), [5](#), [7](#)
- [47] A. Toshev and C. Szegedy. Deeppose: Human pose estimation via deep neural networks. In *CVPR*, 2014. [2](#)
- [48] C. Wang, Y. Wang, and A. L. Yuille. An approach to pose-based action recognition. In *CVPR*, 2013. [1](#)
- [49] Y. Wang, D. Tran, and Z. Liao. Learning hierarchical poselets for human parsing. In *CVPR*, 2011. [2](#)
- [50] S.-E. Wei, V. Ramakrishna, T. Kanade, and Y. Sheikh. Convolutional pose machines. In *CVPR*, 2016. [1](#), [2](#), [3](#), [5](#), [7](#)
- [51] B. Xiaohan Nie, C. Xiong, and S.-C. Zhu. Joint action recognition and pose estimation from video. In *CVPR*, 2015. [1](#)
- [52] W. Yang, S. Li, W. Ouyang, H. Li, and X. Wang. Learning feature pyramids for human pose estimation. In *ICCV*, 2017. [1](#), [2](#), [7](#)
- [53] Y. Yang and D. Ramanan. Articulated pose estimation with flexible mixtures-of-parts. In *CVPR*, 2011. [2](#)
- [54] F. Yu, D. Wang, and T. Darrell. Deep layer aggregation. *arXiv preprint arXiv:1707.06484*, 2017. [3](#)
- [55] M. D. Zeiler and R. Fergus. Visualizing and understanding convolutional networks. In *ECCV*, 2014. [3](#)
- [56] L. L. Zhu, Y. Chen, and A. Yuille. Recursive compositional models for vision: Description and review of recent work. *Journal of Mathematical Imaging and Vision*, 41(1-2):122, 2011. [1](#)