

# Learning Compositional Neural Information Fusion for Human Parsing

Wenguan Wang<sup>\*1,3</sup>, Zhijie Zhang<sup>\*2,1</sup>, Siyuan Qi<sup>3</sup>, Jianbing Shen<sup>1</sup>, Yanwei Pang<sup>2†</sup>, Ling Shao<sup>1</sup>

<sup>1</sup>Inception Institute of Artificial Intelligence, UAE

<sup>2</sup>School of Electrical and Information Engineering, Tianjin University <sup>3</sup>University of California, Los Angeles, USA

{wenguanwang.ai, teesoloj}@gmail.com

<https://github.com/ZzzjzzZ/CompositionalHumanParsing>

## Abstract

This work proposes to combine neural networks with the compositional hierarchy of human bodies for efficient and complete human parsing. We formulate the approach as a neural information fusion framework. Our model assembles the information from three inference processes over the hierarchy: direct inference (directly predicting each part of a human body using image information), bottom-up inference (assembling knowledge from constituent parts), and top-down inference (leveraging context from parent nodes). The bottom-up and top-down inferences explicitly model the compositional and decompositional relations in human bodies, respectively. In addition, the fusion of multi-source information is conditioned on the inputs, i.e., by estimating and considering the confidence of the sources. The whole model is end-to-end differentiable, explicitly modeling information flows and structures. Our approach is extensively evaluated on four popular datasets, outperforming the state-of-the-arts in all cases, with a fast processing speed of 23fps. Our code and results have been released to help ease future research in this direction.

## 1. Introduction

Human parsing, which aims to decompose humans into semantic parts (e.g., arms, legs, etc.), is a crucial yet challenging task for detailed human body configuration analysis in 2D monocular images. It has gained increasing attention owing to its essential role in many areas of application, such as surveillance analysis [34], and fashion synthesis [82], to name a couple.

Recent human parsing approaches have made remarkable progress. Some representative ones [6, 78, 49] are built upon well-designed deep learning architectures for semantic segmentation (e.g., fully convolutional networks (FCNs) [46], DeepLab [4], etc.). Though these achieve

\*Equal contribution.

†Corresponding author: Yanwei Pang.

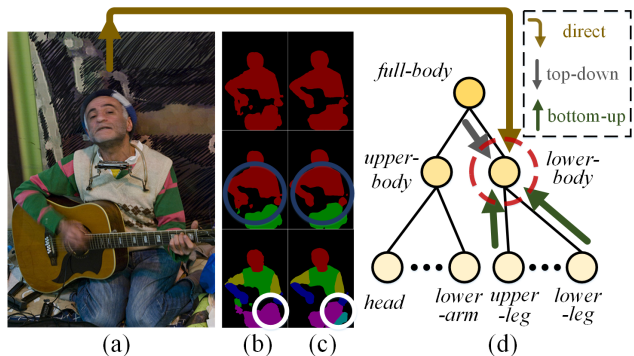


Figure 1: We represent the human body (a) as a hierarchy of multi-level semantic parts, and treat human parsing as a multi-source information fusion process. For each part, information from three sources (direct, bottom-up, and top-down processes) are fused to better capture the structures in this problem. For clarity, we only show the information fusion of the *lower-body* node in the red circle in (d). Compared to directly inferring the human semantics in (b), (c) shows better results after compositional neural information fusion (d).

promising results, they fail to make full use of the rich structures in this task. Some others use extra human joints to better constrain body configurations [22, 71, 54], requiring additional training data of human keypoints and ignoring the compositional relations within human bodies.

In this paper, we segment body parts at multiple levels (see Fig. 1), in contrast to most previous human parsers which only focus on atomic parts (represented as leaf nodes in the human hierarchy). The insight is that estimating the whole graph provides us cross-level information that can assist learning and inference for each body part. This is also evidenced by human perception studies [33, 62, 53]; a global shape can either precede or follow the recognition of its local parts, and both contribute to the final recognition.

We further specify this as a *multi-source information fusion* procedure, which integrates information from the following three processes (see Fig. 1 (d)). **1) Direct inference** (or unconscious inference) from the input image. For example, sometimes humans directly recognize objects by re-

lying on intuitive understanding [51, 69]. **2) Top-down inference**, which recognizes fine-grained components from a whole entity. For example, when recognizing small fine-grained parts, exploring contextual information of the entire object is essential [23, 24, 67] (see the regions in the white circles in Fig. 1). **3) Bottom-up inference**, which associates constituent parts to predict upper-level nodes. When objects are partially occluded or contain complex topologies, humans can assemble sub-parts to assist in recognizing the entities [13, 19] (see the regions in the blue circles in Fig. 1).

Employing the strong learning power of deep neural networks [28, 37], we build a compositional neural information fusion for these three inference processes in an end-to-end manner. This yields a hierarchical human parsing framework to better capture the compositional constraints and human part semantics. In addition, we design our model as a conditional fusion, *i.e.*, the assembly of different information is dependent on the confidence estimations for the sources, instead of simply assuming all the sources are reliable. This is achieved by a learnable gate mechanism, leading to more accurate parsing results.

This paper makes three contributions. 1) We formulate the human parsing problem as a neural information fusion process over a compositionally structured network. 2) We analyze three important sources of information, leading to a novel network architecture that conditionally incorporates direct, top-down, and bottom-up inferences. 3) Our model achieves state-of-the-art performances for comprehensive evaluations on four public datasets (LIP [22], PASCAL-Person-Part [71], ATR [39] and Fashion Clothing [49]). Testing with more than 20K images demonstrates the superiority over existing methods of exploiting compositional structural information for human parsing.

## 2. Related Work

**Hierarchical/Graphical Models in Computer Vision:** Hierarchical/graphical models are powerful for building structured representations, which can reflect task-specific relations and constraints. From early distributional semantic models, part-based models [16, 17], MRF/CRF [31], And-Or grammar model [59], to deep structural networks [30, 15], graph neural networks [20], trainable CRF [79], *etc.*, hierarchical/graphical models have found applications in a wide variety of core computer vision tasks, such as object recognition [55], human parsing [40, 41, 81], pose estimation [34, 66, 61, 68, 35], visual dialog *etc.*, to the extent that they are now ubiquitous in the field. Inspired by their general success, we leverage structural information to design our approach. In addition to directly inferring segments from the image features, we further derive two additional inference processes, *i.e.*, bottom-up and top-down inference, to better capture human structures. This encourages more reasonable results that are consistent with the human body

configuration.

**Information Fusion:** Our method is also inspired by the idea of fusing information from different sources to obtain a better prediction of the target. One typical application of this is sensor fusion, which is a broad field, discussed in more detail in [32]. Many machine learning models can be regarded as information fusion methods: *e.g.*, product of experts [26], Bayesian fusion, ensemble methods [10], and graphical models [64]. Motivated by this general idea, we learn to adaptively fuse the direct inference along with top-down and bottom-up predictions in the compositional human structure for our final prediction.

**Human Parsing Models:** Traditional human parsing models are typically built upon hand-crafted visual features (*e.g.*, color, HoG) [73, 43, 65, 74, 48, 58, 75], low-level image decompositions (*e.g.*, super-pixel) [43, 74, 75], and heuristic hypotheses (*e.g.*, grammars for human body configuration) [3, 12, 8, 11]. Though impressive results have been achieved, these pioneering works require a lot of carefully hand-designed pipelines, and suffer the limited representability of the hand-crafted features.

With the renaissance of connectionism in the computer vision community, recent research efforts take deep neural networks as their main building blocks [70, 54, 78, 50, 49, 77, 45]. More specifically, some efforts address the task as an active template regression problem [39], propagate semantic information from a retrieved, annotated image corpus [44], merge multi-level image context in a unified convolutional neural network [42], or use Graph LSTMs to model human configurations [40, 41]. Some others leverage extra pose information to assist the task [72, 22, 71, 14, 54]. In contrast to the above approaches addressing category-level understanding of human semantics, a few methods operate at an instance level [36, 80, 56].

The aforementioned deep human parsers generally achieve promising results, due to the strong learning power of neural networks [46, 4] and the plentiful availability of annotated data [22, 71]. However, they typically need to pre-segment images into superpixels [40, 41], which breaks the end-to-end story and is time-consuming, or rely on extra human landmarks [72, 22, 71, 14, 54], requiring additional annotations or pre-trained pose estimators. Though [81] also performs multi-level, fine-grained parsing, it neither explores different information flows within human hierarchies nor models the problem from the view of multi-source information fusion.

In contrast, we elaborately design a compositional neural information fusion framework, which explicitly captures human compositional structures and dynamically combines direct, bottom-up and top-down inference modes over the hierarchy. The overall model inherits the complementary advantages of FCNs and hierarchical models, yielding a unified, end-to-end trainable human parsing frame-

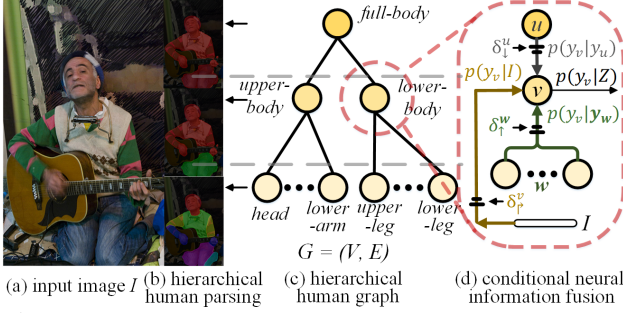


Figure 2: Given an input image (a), our compositional and conditional neural information fusion is performed over the human graph (c) to produce hierarchical parsing results.

work with a strong learning ability, improved representational power, as well as high processing speed.

### 3. Our Approach

Formally, we represent the hierarchical human body structure as a graph  $\mathcal{G} = (\mathcal{V}, \mathcal{E}, \mathcal{Y})$ , where nodes  $v \in \mathcal{V}$  represent human parts in different levels, and edges  $e \in \mathcal{E}$  are two-tuples  $e = (u, v)$  representing the compositional relation that node  $v$  is a part of node  $u$ . As shown in Fig. 2 (c), the nodes are further grouped into  $L (= 3)$  levels:  $\mathcal{V} = \mathcal{V}^1 \cup \dots \cup \mathcal{V}^L$ , where  $\mathcal{V}^1$  are the leaf nodes (the most fine-grained semantic parts typically considered in common human parsers),  $\mathcal{V}^2 = \{\text{upper-body}, \text{lower-body}\}$ , and  $\mathcal{V}^3 = \{\text{full-body}\}$ . For each node  $v$ , we want to infer a segmentation map  $y_v \in \mathcal{Y}$  that is a probability map of its label. Please note that such a problem setting does not introduce any additional annotation requirement, since higher-level annotations can be obtained by simply combining the lower-level labels.

There are three different sources of information when inferring  $y_v$  for  $v$ : 1) the raw input image, 2)  $y_u$  for the parent node  $u$ , and 3)  $y_w$  for all the child nodes  $w$ . We treat the final prediction of  $y_v$  as a fusion of the information from these three sources. Next, we briefly review different methods to modeling this information fusion problem that motivate our solution and network design for human parsing.

#### 3.1. Information Fusion

Information fusion refers to the process of combining information from several sources  $Z = \{z_1, z_2, \dots, z_n\}$  in order to form a unified picture of the measured/predicted target  $y$ . Each source provides an estimation of the target. These sources can be the raw data  $x$  or some other quantities that can be inferred from  $x$ . Several approaches have been proposed to tackle this problem.

- Product of experts (PoE) [26] treats each source as an “expert”. It multiplies the probabilities and then renormalizes:

$$p(y|Z) = \frac{\prod_{i=1}^n p(y|z_i)}{\sum_y \prod_{i=1}^n p(y|z_i)}. \quad (1)$$

- Bayesian fusion. Denoting  $Z_s = \{z_1, z_2, \dots, z_s\}$  as the set of the first  $s$  sources, it factorizes the posterior probability:

$$p(y|Z) = \frac{p(Z_n|y)p(y)}{p(Z_n)} = \frac{p(y)p(z_1|y) \prod_{s=2}^n p(Z_s|Z_{s-1}, y)}{p(z_1) \prod_{s=2}^n p(Z_s|Z_{s-1})}. \quad (2)$$

However, it is too difficult to learn all the conditional distributions. By assuming the independence of different information sources, we have the Naive Bayes:

$$p(y|Z) \propto p(y) \prod_i p(z_i|y), \quad (3)$$

which serves as an approximation of the true distribution.

- Ensemble methods. In this approach, each  $z_i$  is a classifier that predicts  $y$ . A typical ensemble method is Bayesian voting [10], which weights the prediction of each classifier to get the final prediction:

$$p(y|Z) = \sum_{z_i} p(y|z_i)p(z_i|x). \quad (4)$$

The AdaBoost [18] algorithm also falls into this category.

- Graphical models (e.g., conditional random fields). In such models, each  $z_i$  can be viewed as a node that contributes to the conditional probability:

$$p_{\theta}(y|Z) = \exp\left\{\sum_i \phi_{\theta_i}(y, z_i) - A(\theta)\right\}, \quad (5)$$

where  $A(\theta)$  is the log-partition function that normalizes the distribution. Computing  $A(\theta)$  is often intractable, hence the solution is usually given by approximation methods, such as Monte Carlo methods or (loopy) belief propagation [60].

#### 3.2. Compositional Neural Information Fusion

The above methods can all be viewed as ways to approximate the true underlying distribution  $p(y|Z)$ , which can be written as a function of predictions from different information sources  $Z$ :

$$p(y|Z) = f(p(y|z_1), p(y|z_2), \dots, p(y|z_n)). \quad (6)$$

There are potential drawbacks to following the exact solution of one of the above methods. First, they are not entirely consistent with each other. For example, the PoE multiplies all  $p(y|z_i)$  together, whereas ensemble methods compute their weighted sum. Each method approximates the true distribution in a different way and has its own tradeoff. Second, exact inference is difficult and solutions are often approximative (e.g., contrastive divergence [27] is used for PoE and Monte Carlo methods for graphical models).

Therefore, instead of exactly following the computation of one of the above methods, we leverage neural networks to directly model this fusion function, due to their strong ability for flexible feature learning and function approximation [28, 37]. The hope is that we can directly learn to fuse multi-source information for a specific task.

However, the fusion network should not be learned arbitrarily without inductive biases [9, 52, 2], which is the preference for structural explanations exhibited in human reasoning processes. Here, we exploit the compositional nature of the problem and design the network with the following observations:

- In the compositional structure  $\mathcal{G}$ , the final prediction  $p(y_v|Z)$  for each node  $v$  combines information from three

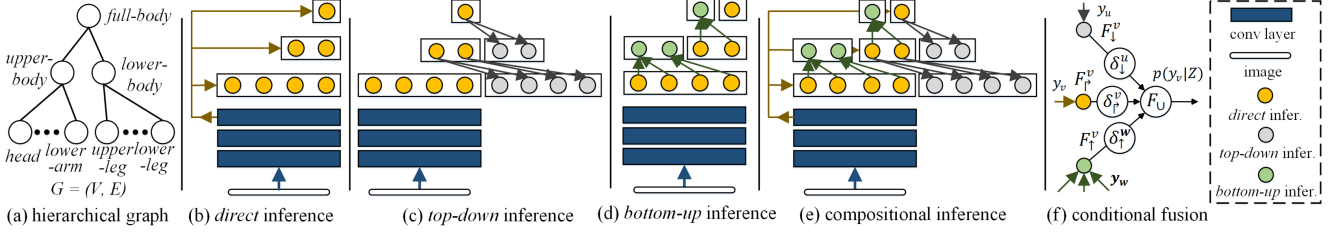


Figure 3: Illustration of our conditional neural information fusion network for hierarchical human parsing. See text for details.

different sources: 1) the direct inference  $p(y_v|x)$  from the raw image input, 2) the top-down inference  $p(y_v|y_u)$  from the parent node  $u$ , which utilizes the **decompositional** relation, and 3) the bottom-up inference  $p(y_v|y_w)$ , which assembles predictions  $y_w$  for all the child nodes  $w$  to leverage the **compositional** relation.

- In many cases, simply fusing different estimations could be problematic. The final decision should be conditioned on the **confidence** of each information source.

Based on the above observations, we design our parser network to learn a *compositional neural information fusion*:

$$p(y_v|Z) = f(\delta_r^v p(y_v|x), \delta_d^u p(y_v|y_u), \delta_b^w p(y_v|y_w)), \quad (7)$$

where the confidence  $\delta$  is a learnable continuous function with outputs from 0 to 1. The symbols  $\uparrow$ ,  $\downarrow$ , and  $\uparrow$  denote direct, top-down, and bottom-up inference, respectively. As shown in Fig. 2 (d), this function fuses information from the three sources in the compositional structure, taking into account the confidence of each source. For neural network realizations of this function, the probability terms can be relaxed to logits, which are essentially log-probabilities.

When carrying out such a prediction, there is one computational issue. Notice that the top-down/bottom-up inferences rely on an estimation of the parent/child node(s). This forms a circular dependency between a parent and its children. To solve this, we treat the direct inference result from the raw data as an initial estimation, and the top-down/bottom-up inferences rely on this initial estimation<sup>1</sup>. Therefore, we decompose the algorithm into three consecutive steps:

- 1. Direct inference.** Given the raw data as input, we assign an estimation  $\tilde{y}_v$  for each node  $v \in \mathcal{V}$ .
- 2. Top-down/bottom-up inference.** We estimate  $p(y_v|\tilde{y}_u)$  and  $p(y_v|\tilde{y}_w)$  based on the estimated  $\tilde{y}_u$  and  $\tilde{y}_w$  in step 1.
- 3. Conditional information fusion.** Based on the above results, we obtain a final prediction for each node  $v$  by  $y_v^* = \text{argmax}_y f(\delta_r^v p(y_v|x), \delta_d^u p(y_v|\tilde{y}_u), \delta_b^w p(y_v|\tilde{y}_w))$ .

This procedure motivates the overall network architecture, where each step above can be learned as a module by a neural network. Next, we discuss our network design.

<sup>1</sup>For some nodes, bottom-up or top-down inference might not exist. The terminal leaf nodes  $\mathcal{V}^1$  do not have bottom-up inference, while the root node  $\mathcal{V}^3$  only has direct and bottom-up inference. For clarity of the method description, we discuss the general case with all three sources.

### 3.3. Network Architecture

Our model stacks the following parts to form an end-to-end system for hierarchical human parsing. The system does not require any preprocessing and the modules are FCNs, so it achieves high efficiency.

**Direct Inference Network.** This directly predicts a segmentation map  $\tilde{y}_v$  for each node  $v$  (a human part), using information from the image (see Fig. 3 (b)). Formally, given an input image  $I \in \mathbb{R}^{K \times K \times 3}$ , a backbone network  $B$  (i.e., a DeepLabV3-like network, parameterized by  $\mathbf{W}_B$ ) is first employed to obtain a new effective image representation  $h_I$ :

$$\text{image embedding: } h_I = F_B(I; \mathbf{W}_B) \in \mathbb{R}^{k \times k \times c}. \quad (8)$$

As the nodes  $\mathcal{V}$  capture explicit semantics, a specific feature  $h_v$  for each node  $v$  is desired for more efficient representation. However, using several different, node-specific embedding networks will lead to a high computational cost. To remedy this, for each  $l$ -th level, we first apply a *level-specific* FCN (LSF) to describe the level-wise semantics and contextual relations:

$$\text{level-specific embedding: } h_{\text{LSF}}^l = F_{\text{LSF}}^l(h_I; \mathbf{W}_{\text{LSF}}^l) \in \mathbb{R}^{k \times k \times c}, \quad (9)$$

where  $l \in \{1, 2, 3\}$ . More specifically, three LSFs ( $F_{\text{LSF}}^1$ ,  $F_{\text{LSF}}^2$ , and  $F_{\text{LSF}}^3$ ) are learned to extract three level-specific embeddings ( $h_{\text{LSF}}^1$ ,  $h_{\text{LSF}}^2$ , and  $h_{\text{LSF}}^3$ ). Further, for each node  $v$ , an independent channel-attention block, Squeeze-and-Excitation (SE) [29], is applied to obtain its specific feature:

$$\text{node-specific embedding: } h_v = F_{\text{SE}}^v(h_{\text{LSF}}^l; \mathbf{W}_{\text{SE}}^v) \in \mathbb{R}^{k \times k \times c}, \quad (10)$$

where  $v \in \mathcal{V}^l$  (i.e.,  $v$  is located in the  $l$ -th level). By explicitly modelling the interdependencies between channels,  $F_{\text{SE}}^v$  allows us to adaptively recalibrate the channel-wise features of  $h_{\text{LSF}}^l$  to generate node-wise representations. Meanwhile, due to its light-weight nature, we can achieve our goal with minimal computational overhead. Then, the direct inference network  $F_{\uparrow}$  reads the feature and predicts the segmentation map  $\tilde{y}_v$ :

$$\text{logit}(\tilde{y}_v|I) = F_{\uparrow}(h_v; \mathbf{W}_{\uparrow}) \in \mathbb{R}_{\geq 0}^{k \times k}. \quad (11)$$

**Top-down Inference Network.** Based on the outputs from the direct inference network, the top-down inference predicts segmentation maps by considering human decompositional structures. Specifically, for node  $v$ , the top-down network  $F_{\downarrow}$  leverages the initial estimation  $\tilde{y}_u$  of its parent

node  $u$  as high-level contextual information for prediction (see Fig. 3 (c)):

$$\text{logit}(y_v|\tilde{y}_u) = F_{\downarrow}(y_v|\tilde{y}_u; h_v, \mathbf{W}_{\downarrow}) = F_{\downarrow}([\tilde{y}_u, h_v]) \in \mathbb{R}_{\geq 0}^{k \times k}. \quad (12)$$

Here, the concatenated feature  $[\tilde{y}_u, h_v]$  is fed into the FCN-based  $F_{\downarrow}$ , parameterized by  $\mathbf{W}_{\downarrow}$ , for top-down inference.

**Bottom-up Inference Network.** One major difference to the top-down network is that, for each node  $v$ , the bottom-up network needs to gather information (*i.e.*,  $\tilde{y}_w \in \mathbb{R}_{\geq 0}^{k \times k \times |w|}$ ) from multiple descendants  $w$ . Thanks to the compositional relations between  $w$  and  $v$ , we can transform  $\tilde{y}_w$  to a fixed one-channel representation  $\tilde{y}_w$  through *position-wise max-pooling* PMP (across channels):

$$\tilde{y}_w = \text{PMP}([\tilde{y}_w]_{w \in \omega}) \in \mathbb{R}_{\geq 0}^{k \times k \times 1}, \quad (13)$$

where  $[\cdot]$  is a concatenation operation. Then, the bottom-up network  $F_{\uparrow}$  gives a prediction according to compositional relations (see Fig. 3 (d)):

$$\text{logit}(y_v|\tilde{y}_w) = F_{\uparrow}(y_v|\tilde{y}_w; h_v, \mathbf{W}_{\uparrow}) = F_{\uparrow}([\tilde{y}_w, h_v]) \in \mathbb{R}_{\geq 0}^{k \times k}. \quad (14)$$

**Conditional Fusion Network.** Before making the final prediction, we estimate the confidence  $\delta$  of each information source using a neural gate function. For the direct inference of a node  $v$ , we estimate the confidence by:

$$\delta_v^v = \sigma(\mathbf{C}_v^v \cdot \text{CAP}(h_v)) \in [0, 1], \quad (15)$$

where  $\sigma$  is the *sigmoid* function. Here, CAP stands for *channel-wise average pooling*, which has been proved a simple yet effective way for capturing the global statistics of convolutional features [29, 63].  $\mathbf{C}_v^u \in \mathbb{R}^{1 \times C}$  indicates a small fully connected layer that maps the  $C$ -dimensional statistic vector  $\text{CAP}(h_v) \in \mathbb{R}^C$  of  $h_v$  into a confidence score.

The confidence scores for the top-down and bottom-up processes follow a similar computational framework:

$$\begin{aligned} \delta_{\downarrow}^u &= \sigma(\mathbf{C}_{\downarrow}^u \cdot \text{CAP}(h_u)) \in [0, 1], \\ \delta_{\uparrow}^w &= \sigma(\mathbf{C}_{\uparrow}^w \cdot \text{CAP}([h_w]_{w \in \omega})) \in [0, 1], \end{aligned} \quad (16)$$

where  $\mathbf{C}_{\downarrow}^u \in \mathbb{R}^{1 \times C}$  and  $\mathbf{C}_{\uparrow}^w \in \mathbb{R}^{1 \times C|w|}$ . Specifically, for the bottom-up process, we concatenate all the child node embeddings  $[h_w]_{w \in \omega} \in \mathbb{R}^{k \times k \times C|w|}$ . This means our decision is made upon the confidence of the union of the child nodes. Here, the confidence of a source can be viewed as a global score or statistic for interpreting the quality of the feature, which is learnt in an implicit manner.

Finally, for each node  $v$ , the fusion network  $F_{\cup}$  combines the results from the three inference networks above for final prediction (see Fig. 3 (e)):

$$\text{logit}(y_v|Z) = F_{\cup}(\delta_v^v F_v^v, \delta_{\downarrow}^u F_{\downarrow}^u, \delta_{\uparrow}^w F_{\uparrow}^w; \mathbf{W}_{\cup}) \in \mathbb{R}_{\geq 0}^{k \times k \times 1}, \quad (17)$$

where  $F_{\cup} : \mathbb{R}_{\geq 0}^{k \times k \times 3} \rightarrow \mathbb{R}_{\geq 0}^{k \times k \times 1}$  is implemented by a small FCN, parameterized by  $\mathbf{W}_{\cup}$ . Fig. 4 provides an illustration of our conditional fusion process. As can be seen,  $\delta$  provides a learnable gate mechanism that suggests how much

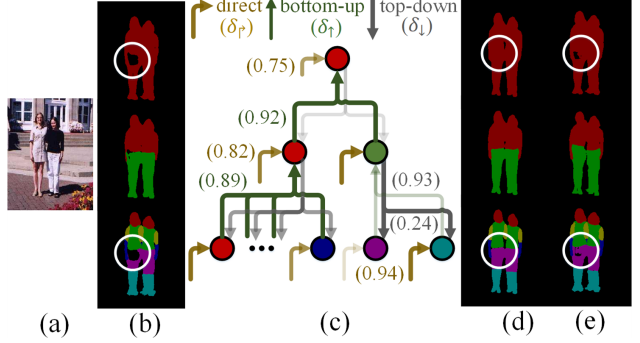


Figure 4: Illustration of our compositional inference and conditional fusion. (a) Input image. (b) Parsing results of direct inference. (c) Conditional information fusion, where the arrows with darker colors indicate higher values of gates  $\delta$ . For clarity, in (·) we only show the gate values for a few inference processes. (d) Parsing results w/ compositional inference and conditional fusion. (e) Parsing results of compositional inference only. The improved regions are highlighted in white circles.

information can be used from a source (see Fig. 4 (c)). It is able to dynamically change the amount of information for different inference processes, *i.e.*, condition on the sources (see Fig. 4 (d)). Thus, it yields better results than statically fusing the information with a weight-fixed fusion function (see Fig. 4 (e)). More detailed studies of our conditional and compositional fusion can be found in § 4.4 and Fig. 6.

**Loss Function.** To obtain the final segmentation map from  $\text{logit}(y_v|Z)$ , we apply a *softmax* function over the logits of nodes in the same level. Thus, for each level, all the inference networks;  $F_{\uparrow}$ ,  $F_{\downarrow}$ ,  $F_{\cup}$ , and the fusion network  $F_{\cup}$  are trained by the standard cross-entropy loss:

$$\mathcal{L}^l = \mathcal{L}_r^{\text{CE}} + \mathcal{L}_{\downarrow}^{\text{CE}} + \mathcal{L}_{\uparrow}^{\text{CE}} + \mathcal{L}_{\cup}^{\text{CE}}. \quad (18)$$

### 3.4. Implementation Details

**Backbone Network.** Our feature extraction network  $F_B$  in Eq. 8 uses the convolutional blocks of ResNet101 [25]. The stride is set to 16, *i.e.*, the resolution of the output is 1/16 of that of the input, for high computational efficiency. In addition, the ASPP module [5] is applied for extracting more effective features with multi-scale context. The ASPP-enhanced feature is compressed by a  $1 \times 1$  convolutional layer with *ReLU* activation. The compressed 512- $d$  feature is further  $\times 2$  upsampled and element-wisely added with the feature from the second convolutional block of ResNet101, to encode more spatial details. Thus, given an input image  $I$  with a size of  $K \times K$ , the feature extraction network  $B$  produces a new image representation  $h_I \in \mathbb{R}^{\frac{K}{8} \times \frac{K}{8} \times 512}$ .

**Direct Inference Network.** We implement  $F_{\text{LSF}}^l$  (Eq. 9) using a  $3 \times 3$  convolutional layer with Batch Normalization (BN) and *ReLU* activation, whose parameters are shared by all the nodes located in the  $l$ -th level. This is used for extracting specific features  $\{h_{\text{LSF}}^1, h_{\text{LSF}}^2, h_{\text{LSF}}^3\}$  for the three semantic-levels. For each node  $v$ , an independent SE [29] block,  $F_{\text{SE}}^v$  in Eq. 10, is further applied to extract its spe-

cific embedding  $\mathbf{h}_v \in \mathbb{R}^{\frac{K}{8} \times \frac{K}{8} \times 512}$  with an extremely light-weight architecture. Then,  $F_\uparrow$  in Eq. 11 is implemented by a stack of three  $1 \times 1$  convolutional layers.

**Top-down/Bottom-up Inference Network.** The architectures of the top-down  $F_\downarrow$  (Eq. 12) and bottom-up  $F_\uparrow$  (Eq. 14) inference networks are very similar, and only differ in their strategies of processing the input features (see Eq. 13). Both are achieved by three cascaded convolutional layers, with convolution sizes of  $3 \times 3$ ,  $3 \times 3$  and  $1 \times 1$ , respectively.

**Information Fusion Network.**  $F_\cup$  in Eq. 17 consists of three  $1 \times 1$  convolutional layers with *ReLU* activations for non-linear mapping.

## 4. Experiments

### 4.1. Experimental Settings

**Datasets:** We perform extensive experiments on the following four widely-tested datasets:

- **LIP [22]** has 50,462 single-person images with elaborate pixel-wise annotations of 19 part categories (e.g., *hair*, *face*, *left-/right-arms*, *left-/right-legs*, *left-/right-shoes*, etc.). The images are divided into 30,462 samples for training, 10,000 for validation and 10,000 for testing.
- **PASCAL-Person-Part [71]** contains multiple humans per image in unconstrained poses and occlusions (1,716 for training and 1,817 for testing). It provides careful pixel-wise annotations for six body parts (i.e., *head*, *torso*, *upper-/lower-arms*, and *upper-/lower-legs*).
- **ATR [39]** includes 7,700 images (6,000 for training, 700 for validation and 1,000 for testing), annotated at pixel-level with 17 categories, e.g., *hat*, *sunglass*, *face*, *upper-clothes*, *pants*, *left-/right-arms*, *left-/right-legs*, etc.
- **Fashion Clothing [49]** consists of Colorful Fashion Parsing [43], Fashionista [74], and Clothing Co-Parsing [75]. It is more concerned with human clothing details, including 17 categories (e.g., *glass*, *hair*, *pants*, *shoes*, *shirt*, *upper-clothes*, *skirt*, *scarf*, *socks*, etc.). It has 4,371 images in total (3,934 for training, and 437 for testing).

**Evaluation Metrics:** For LIP, following its standard protocol [78], we report pixel accuracy, mean accuracy and mean Intersection-over-Union (mIoU). For PASCAL-Person-Part, following conventions [70, 71, 50], the performance is evaluated in terms of mIoU. For ATR and Fashion Clothing, we report five metrics as [49] does, including pixel accuracy, foreground accuracy, average precision, average recall, and average F1-score.

**Training Settings:** During training, the weights of the backbone network are loaded from ResNet101 [25] pre-trained on ImageNet [57], and the remaining layers are randomly initialized. For data preparation, following [56, 21], we apply data augmentation techniques for all the training data, including randomly scaling, cropping and left-right flipping. The random scale is set from 0.5 to 2.0, while the crop size is set to  $473 \times 473$ . For optimization, we adopt

Methods	pixAcc.	Mean Acc.	Mean IoU
SegNet [1]	69.04	24.00	18.17
FCN-8s [46]	76.06	36.75	28.29
DeepLabV2 [4]	82.66	51.64	41.64
Attention [6]	83.43	54.39	42.92
Attention+SSL [22]	84.36	54.94	44.73
ASN [47]	-	-	45.41
SSL [22]	-	-	46.19
MMAN [50]	-	-	46.81
SS-NAN [78]	87.59	56.03	47.92
MuLA [54]	<b>88.5</b>	60.5	49.3
CE2P [56]	87.37	63.20	53.10
Ours	88.03	<b>68.80</b>	<b>57.74</b>

Table 1: **Comparison of pixel accuracy, mean accuracy and mIoU on LIP val [22].** (Higher values are better. The best score is marked in **bold**. These notes are the same for other tables.)

SGD with a momentum of 0.9, and weight\_decay of 0.0005. For the learning rate, we use the ‘poly’ learning rate schedule [4, 76],  $lr = base\_lr \times (1 - \frac{iters}{total\_iters})^{power}$ , in which  $power=0.9$  and  $base\_lr=0.007$ . The  $total\_iters$  is  $epochs \times batch\_size$ , where  $batch\_size=40$  and  $epochs=150$ . We use multiple GPUs for the consumption of the large  $batch\_size$ , and implement Synchronized Cross-GPU BN.

**Testing Phase:** Following general protocol [76, 54], we average the per-pixel classification scores at multiple scales with flipping, i.e., the scale is 0.5 to 1.5 (in increments of 0.25) times the original size. Our model does not require any other pre-/post-processing steps (i.e., over-segmentation [40, 38], human pose [71], CRF [71]), and thus achieves a processing speed of 23.0fps, averaged on PASCAL-Person-Part, which is faster than previous deep human parsers, such as Joint [71] (0.1fps), Attention+SSL [22] (2.0fps), MMAN [50] (3.5fps) and MuLA [54] (15fps).

**Reproducibility:** Our method is implemented on PyTorch and trained on four NVIDIA Tesla V100 GPUs with a 32GB memory per-card. All the testing procedures are carried out on a single NVIDIA TITAN Xp GPU with 12GB memory for a fair speed comparison. To provide full details of our training and testing processes, we release our code in <https://github.com/ZzzjzzZ/CompositionalHumanParsing>.

### 4.2. Quantitative Results

We compare the proposed method with several strong baselines on the four aforementioned challenging datasets.

**LIP [22]:** We compare our method with 11 state-of-the-arts on LIP val set in Table 1. Our method achieves a huge boost in average IoU (4.64% better than the second best method, CE2P [56] and 8.4% better than the third best, MuLA [54]). To verify its effectiveness in detail, we report per-class IoU in Table 2. Our model improves the performance over almost all classes, especially for the ones typically associated with small regions (e.g., *gloves*, *sunglasses*, *socks*, *shoes*), due to our top-down inference strategy. The

Methods	Hat	Hair	Glov	Sung	Clot	Dress	Coat	Sock	Pant	Suit	Scarf	Skirt	Face	L-Arm	R-Arm	L-Leg	R-Leg	L-Sh	R-Sh	B.G.	Ave.
SegNet [1]	26.60	44.01	0.01	0.00	34.46	0.00	15.97	3.59	33.56	0.01	0.00	0.00	52.38	15.30	24.23	13.82	13.17	9.26	6.47	70.62	18.17
FCN-8s [46]	39.79	58.96	5.32	3.08	49.08	12.36	26.82	15.66	49.41	6.48	0.00	2.16	62.65	29.78	36.63	28.12	26.05	17.76	17.70	78.02	28.29
DeepLabV2 [4]	56.48	65.33	29.98	19.67	62.44	30.33	51.03	40.51	69.00	22.38	11.29	20.56	70.11	49.25	52.88	42.37	35.78	33.81	32.89	84.53	41.64
Attention [6]	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
Attention+SSL [22]	59.75	67.25	28.95	21.57	65.30	29.49	51.92	38.52	68.02	24.48	14.92	24.32	71.01	52.64	55.79	40.23	38.80	28.08	29.03	84.56	44.73
ASN [47]	56.92	64.34	28.07	17.78	64.90	30.85	51.90	39.75	71.78	25.57	7.97	17.63	70.77	53.53	56.70	49.58	48.21	34.57	33.31	84.01	45.41
SSL [22]	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 [50]	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
SS-NAN [78]	63.86	70.12	30.63	23.92	70.27	33.51	56.75	40.18	72.19	27.68	16.98	26.41	75.33	55.24	58.93	44.01	41.87	29.15	32.64	<b>88.67</b>	47.92
CE2P [56]	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
Ours	<b>69.55</b>	<b>73.45</b>	<b>45.17</b>	<b>41.45</b>	<b>70.57</b>	<b>38.52</b>	<b>57.94</b>	<b>54.02</b>	<b>75.07</b>	<b>28.00</b>	<b>31.92</b>	<b>30.20</b>	<b>76.38</b>	<b>68.28</b>	<b>69.49</b>	<b>65.52</b>	<b>65.51</b>	<b>52.67</b>	<b>53.38</b>	87.99	<b>57.74</b>

Table 2: Per-class comparison of mIoU with state-of-the-art methods on LIP v1 [22].

Methods	Head	Torso	U-Arm	L-Arm	U-Leg	L-Leg	B.G.	Ave.
HAZN [70]	80.79	59.11	43.05	42.76	38.99	34.46	93.59	56.11
Attention [6]	81.47	59.06	44.15	42.50	38.28	35.62	93.65	56.39
LG-LSTM [41]	82.72	60.99	45.40	47.76	42.33	37.96	88.63	57.97
Attention+SSL [22]	83.26	62.40	47.80	45.58	42.32	39.48	94.68	59.36
Attention+MMAN [50]	82.58	62.83	48.49	47.37	42.80	40.40	94.92	59.91
Graph LSTM [40]	82.69	62.68	46.88	47.71	45.66	40.93	94.59	60.16
SS-NAN [78]	86.43	67.28	51.09	48.07	44.82	42.15	97.23	62.44
Structure LSTM [38]	82.89	67.15	51.42	48.72	51.72	45.91	97.18	63.57
Joint [71]	85.50	67.87	54.72	54.30	48.25	44.76	95.32	64.39
DeepLabV2 [4]	-	-	-	-	-	-	-	64.94
MuLA [54]	-	-	-	-	-	-	-	65.1
PCNet [81]	86.81	69.06	55.35	55.27	50.21	48.54	96.07	65.90
Holistic [36]	-	-	-	-	-	-	-	66.3
WSHP [14]	87.15	72.28	57.07	56.21	52.43	50.36	<b>97.72</b>	67.60
PGN [21]	<b>90.89</b>	<b>75.12</b>	55.83	<b>64.61</b>	55.42	41.57	95.33	68.40
Ours	88.02	72.91	<b>64.31</b>	63.52	<b>55.61</b>	<b>54.96</b>	96.02	<b>70.76</b>

Table 3: Per-class comparison of mIoU with state-of-the-art methods on PASCAL-Person-Part test [71].

results are also impressive for *arms*, *legs*, and *shoes*, demonstrating our model’s ability to distinguish between “left” and “right” with the help of composition relations.

**PASCAL-Person-Part [71]:** On its *test* set, we compare our method with 15 state-of-the-arts using IoU score. As shown in Table 3, our model outperforms previous methods across the vast majority of classes and on average.

**ATR [39]:** Table 4 gives evaluation on ATR *test* set. Our model again outperforms other competitors across most metrics. In particular, it achieves an average F-1 score of 85.51%, which is 3.45% better than TGPNet [49] and 5.37% better than Co-CNN [42].

**Fashion Clothing [49]:** We compare our method with five famous models on Fashion Clothing *test*, where we take the pre-computed evaluation from [49]. From Table 5, we observe our model surpasses other competitors across all metrics by a large margin. Notably, it yields an F-1 score of 58.12%, significantly outperforming TGPNet [49] and Attention [6] by +6.20% and +9.44%, respectively.

Overall, our model consistently obtains promising results over different datasets, which clearly demonstrates its superior performance and strong generalizability. This also distinguishes our model from several previous state-of-the-art deep human parsers, such as [22, 71, 14, 54], since it does not use extra pose annotations during training.

Methods	pixAcc.	F.G. Acc.	Prec.	Recall	F-1
Yamaguchi [74]	84.38	55.59	37.54	51.05	41.80
Paperdoll [73]	88.96	62.18	52.75	49.43	44.76
M-CNN [44]	89.57	73.98	64.56	65.17	62.81
ATR [39]	91.11	71.04	71.69	60.25	64.38
DeepLabV2 [4]	94.42	82.93	78.48	69.24	73.53
PSPNet [76]	95.20	80.23	79.66	73.79	75.84
Attention [6]	95.41	85.71	81.30	73.55	77.23
DeepLabV3+ [7]	95.96	83.04	80.41	78.79	79.49
Co-CNN [42]	96.02	83.57	<b>84.95</b>	77.66	80.14
TGPNet [49]	<b>96.45</b>	<b>87.91</b>	83.36	80.22	81.76
Ours	96.26	<b>87.91</b>	84.62	<b>86.41</b>	<b>85.51</b>

Table 4: Comparison of accuracy, foreground accuracy, average precision, recall and F1-score on ATR *test* [39]. Please see the supplementary material for per-class performance.

Methods	pixAcc.	F.G. Acc.	Prec.	Recall	F-1
Yamaguchi [74]	81.32	32.24	23.74	23.68	22.67
Paperdoll [73]	87.17	50.59	45.80	34.20	35.13
DeepLabV2 [4]	87.68	56.08	35.35	39.00	37.09
Attention [6]	90.58	64.47	47.11	50.35	48.68
TGPNet [49]	91.25	66.37	50.71	53.18	51.92
Ours	<b>92.20</b>	<b>68.59</b>	<b>56.84</b>	<b>59.47</b>	<b>58.12</b>

Table 5: Comparison of pixel accuracy, foreground pixel accuracy, average precision, average recall and average f1-score on Fashion Clothing *test* [49].

### 4.3. Qualitative Results

In Fig. 5, we show some visual results on PASCAL-Person-Part *test* set. Our method yields more precise predictions compared to SS-NAN [78], DeepLabV2 [4] and PGN [21]. For example, in the last row, our method correctly labels the lower-legs of the rider, while other methods [78, 4, 21] face difficulties in this case. Our model also provides clearer details for small parts. Observed from the second row, the small lower-arm regions can be successfully segmented out with the constraint of top-down inference. In general, by effectively exploiting the human semantic hierarchy, our approach outputs reasonable results for confusing labels on the human parsing task.

### 4.4. Ablation Study

Table 6 shows an evaluation of our full model compared to ablated versions without certain key components. All the

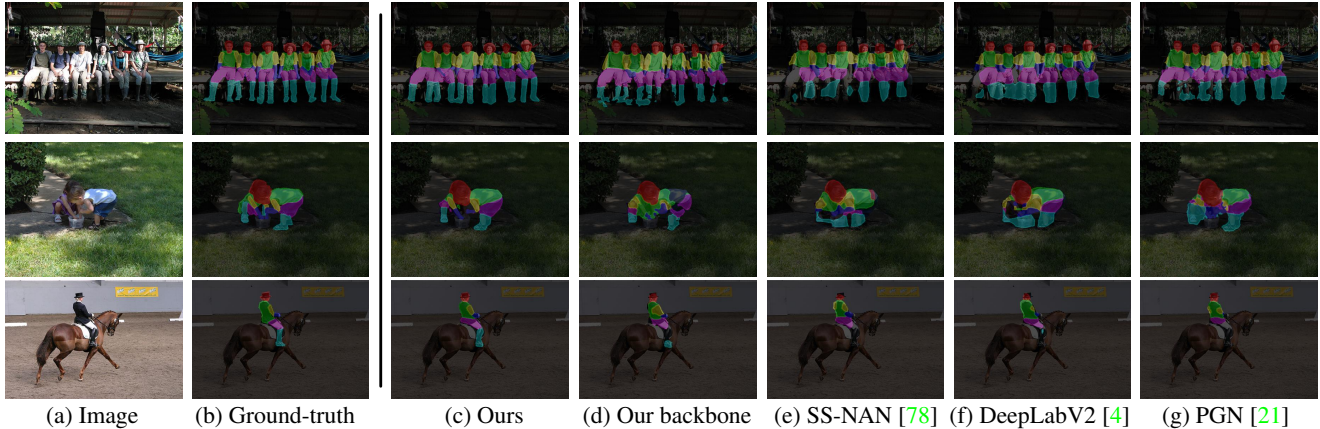


Figure 5: **Visual comparison results on PASCAL-Person-Part test set.** Our model (c) outputs more semantically meaningful and precise predictions, compared to our backbone network (d) and other famous competitors [78, 4, 21] (e-g). See § 4.3 for details.

Aspects	Methods	mIoU		
		1st Level	2nd Level	3rd Level
<b>Full model</b>	direct + bottom-up + top-down + conditional fusion	70.76	81.62	91.31
Backbone	direct infer. w/o hierarchy	64.14	-	-
Variant	direct	65.27	77.83	88.29
	direct + bottom-up	65.42	78.37	90.10
	direct + top-down	69.02	78.91	88.40
	direct + bottom-up + top-down	69.43	80.34	91.02

Table 6: **Ablation study on PASCAL-Person-Part test [71].**

variants are retrained independently with their specific network architectures. Here, 1st-Level denotes the automatic parts (e.g., head, leg, etc.) in  $\mathcal{V}^1$ , 2nd-Level  $\mathcal{V}^2$  (lower/upper body), and 3rd-Level  $\mathcal{V}^3$  (full body). The experiments are performed on PASCAL-Person-Part [71] test set using mIoU metric. Three essential conclusions can be drawn from our results. First, instead of only modeling the fine-grained parts in  $\mathcal{V}^1$  (i.e., backbone), even directly learning to parse the whole human hierarchy (i.e., direct) can bring a performance gain (64.14→65.27). This suggests that modeling the human hierarchy leads to a comprehensive understanding of human semantics. Second, further considering bottom-up and top-down inference provides substantial performance gain, demonstrating the benefit of exploiting human structures and efficient information fusion strategies in this problem. Note that in (direct vs. direct+bottom-up) and (direct+top-down vs. direct+bottom-up+top-down), even for the 1st-level nodes that do not have bottom-up inference, the training itself brings performance gain. The reason is that the bottom-up inference explicitly captures compositional relations and thus improves the quality of the learnt features. Similar observations can also be found in (direct vs. direct+top-down) and (direct+bottom-up vs. direct+bottom-up+top-down) for the 3rd-level node. These observations suggest the compositional information fusion not only improves the predictions during inference but also boosts the learning ability of our human parser model. Third, conditionally fusing information boosts performance, as the information

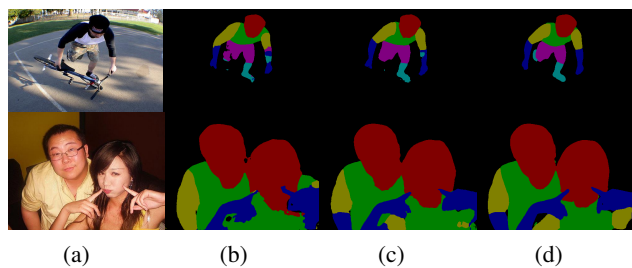


Figure 6: **Visual comparison results** from our (b) backbone, (c) compositional information fusion, and (d) full model.

from low-quality sources can be suppressed. This also provides a new glimpse into the information fusion mechanism over hierarchical models. A visual comparison between the results from our backbone network, our model only using compositional fusion and our full model can be found in Fig. 6 (b-d), which intuitively shows the improvements from our conditional and compositional information fusion.

## 5. Conclusion

In this work, we parse human parts in a hierarchical form, enabling us to capture human semantics from a more comprehensive view. We tackle this hierarchical human parsing problem through a neural information fusion framework that explores the compositional relations within human structures. It efficiently combines the information from the direct, top-down, and bottom-up inference processes while considering the reliability of each process. Extensive quantitative and qualitative comparisons performed on five datasets demonstrate that our method outperforms the current alternatives by a large margin.

**Acknowledgements** The authors thank Prof. Song-Chun Zhu and Prof. Ying Nian Wu from UCLA Statistics Department for helpful comments on this work. This work reported herein was supported in part by DARPA XAI grant N66001-17-2-4029, ARO grant W911NF-18-1-0296, CCF-Tencent Open Fund, and the National Natural Science Foundation of China (No. 61632018).



## References

- [1] Vijay Badrinarayanan, Alex Kendall, and Roberto Cipolla. Segnet: A deep convolutional encoder-decoder architecture for image segmentation. *IEEE TPAMI*, 39(12):2481–2495, 2017.
- [2] Peter W Battaglia, Jessica B Hamrick, Victor Bapst, Alvaro Sanchez-Gonzalez, Vinicius Zambaldi, Mateusz Malinowski, Andrea Tacchetti, David Raposo, Adam Santoro, Ryan Faulkner, et al. Relational inductive biases, deep learning, and graph networks. *arXiv preprint arXiv:1806.01261*, 2018.
- [3] Hong Chen, Zi Jian Xu, Zi Qiang Liu, and Song Chun Zhu. Composite templates for cloth modeling and sketching. In *CVPR*, 2006.
- [4] 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 TPAMI*, 40(4):834–848, 2018.
- [5] Liang-Chieh Chen, George Papandreou, Florian Schroff, and Hartwig Adam. Rethinking atrous convolution for semantic image segmentation. *CoRR*, abs/1706.05587, 2017.
- [6] Liang-Chieh Chen, Yi Yang, Jiang Wang, Wei Xu, and Alan L Yuille. Attention to scale: Scale-aware semantic image segmentation. In *CVPR*, 2016.
- [7] Liang-Chieh Chen, Yukun Zhu, George Papandreou, Florian Schroff, and Hartwig Adam. Encoder-decoder with atrous separable convolution for semantic image segmentation. In *ECCV*, 2018.
- [8] 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 *CVPR*, 2014.
- [9] Kenneth James Williams Craik. *The nature of explanation*, volume 445. CUP Archive, 1967.
- [10] Thomas G Dietterich. Ensemble methods in machine learning. In *International workshop on multiple classifier systems*, 2000.
- [11] Jian Dong, Qiang Chen, Xiaohui Shen, Jianchao Yang, and Shuicheng Yan. Towards unified human parsing and pose estimation. In *CVPR*, 2014.
- [12] Jian Dong, Qiang Chen, Wei Xia, Zhongyang Huang, and Shuicheng Yan. A deformable mixture parsing model with parselets. In *ICCV*, 2013.
- [13] Boris Epshtein, Ita Lifshitz, and Shimon Ullman. Image interpretation by a single bottom-up top-down cycle. *Proceedings of the National Academy of Sciences*, 105(38):14298–14303, 2008.
- [14] 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. In *CVPR*, 2018.
- [15] Hao-Shu Fang, Yuanlu Xu, Wenguan Wang, Xiaobai Liu, and Song-Chun Zhu. Learning pose grammar to encode human body configuration for 3d pose estimation. In *AAAI*, 2018.
- [16] Pedro F Felzenszwalb, Ross B Girshick, David McAllester, and Deva Ramanan. Object detection with discriminatively trained part-based models. *IEEE TPAMI*, 32(9):1627–1645, 2010.
- [17] Pedro F Felzenszwalb and Daniel P Huttenlocher. Pictorial structures for object recognition. *IJCV*, 61(1):55–79, 2005.
- [18] Yoav Freund and Robert E Schapire. A decision-theoretic generalization of on-line learning and an application to boosting. *Journal of Computer and System Sciences*, pages 119–139, 1997.
- [19] James Jerome Gibson. The senses considered as perceptual systems. 1966.
- [20] Justin Gilmer, Samuel S Schoenholz, Patrick F Riley, Oriol Vinyals, and George E Dahl. Neural message passing for quantum chemistry. In *ICML*, 2017.
- [21] Ke Gong, Xiaodan Liang, Yicheng Li, Yimin Chen, Ming Yang, and Liang Lin. Instance-level human parsing via part grouping network. In *ECCV*, 2018.
- [22] 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 *CVPR*, 2017.
- [23] Richard Langton Gregory. The intelligent eye. 1970.
- [24] Kalanit Grill-Spector and Kevin S Weiner. The functional architecture of the ventral temporal cortex and its role in categorization. *Nature Reviews Neuroscience*, 15(8):536, 2014.
- [25] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *CVPR*, 2016.
- [26] Geoffrey E Hinton. Products of experts. In *International Conference on Artificial Neural Networks*, 1999.
- [27] Geoffrey E Hinton. Training products of experts by minimizing contrastive divergence. *Neural computation*, pages 1771–1800, 2002.
- [28] Kurt Hornik, Maxwell Stinchcombe, and Halbert White. Multilayer feedforward networks are universal approximators. *Neural networks*, pages 359–366, 1989.
- [29] Jie Hu, Li Shen, and Gang Sun. Squeeze-and-excitation networks. In *CVPR*, 2018.
- [30] Ashesh Jain, Amir R Zamir, Silvio Savarese, and Ashutosh Saxena. Structural-RNN: Deep learning on spatio-temporal graphs. In *CVPR*, 2016.
- [31] Thorsten Joachims, Thomas Finley, and Chun-Nam John Yu. Cutting-plane training of structural svms. *Machine Learning*, 77(1):27–59, 2009.
- [32] Bahador Khaleghi, Alaa Khamis, Fakhreddine O Karray, and Saiedeh N Razavi. Multisensor data fusion: A review of the state-of-the-art. *Information fusion*, 14(1):28–44, 2013.
- [33] Ruth Kimchi. Primacy of wholistic processing and global/local paradigm: a critical review. *Psychological bulletin*, page 24, 1992.
- [34] Lubor Ladicky, Philip HS Torr, and Andrew Zisserman. Human pose estimation using a joint pixel-wise and part-wise formulation. In *CVPR*, 2013.
- [35] Jiefeng Li, Can Wang, Hao Zhu, Yihuan Mao, Hao-Shu Fang, and Cewu Lu. Crowdpose: Efficient crowded scenes pose estimation and a new benchmark. In *CVPR*, 2019.
- [36] Qizhu Li, Anurag Arnab, and Philip HS Torr. Holistic, instance-level human parsing. *arXiv preprint arXiv:1709.03612*, 2017.
- [37] Shiyu Liang and R. Srikant. Why deep neural networks for function approximation? In *JCLR*, 2017.
- [38] Xiaodan Liang, Liang Lin, Xiaohui Shen, Jiashi Feng,

- Shuicheng Yan, and Eric P Xing. Interpretable structure-evolving lstm. In *CVPR*, 2017.
- [39] Xiaodan Liang, Si Liu, Xiaohui Shen, Jianchao Yang, Luoqi Liu, Jian Dong, Liang Lin, and Shuicheng Yan. Deep human parsing with active template regression. *IEEE TPAMI*, 37(12):2402–2414, 2015.
- [40] Xiaodan Liang, Xiaohui Shen, Jiashi Feng, Liang Lin, and Shuicheng Yan. Semantic object parsing with graph lstm. In *ECCV*, 2016.
- [41] Xiaodan Liang, Xiaohui Shen, Donglai Xiang, Jiashi Feng, Liang Lin, and Shuicheng Yan. Semantic object parsing with local-global long short-term memory. In *CVPR*, 2016.
- [42] Xiaodan Liang, Chunyan Xu, Xiaohui Shen, Jianchao Yang, Si Liu, Jinhui Tang, Liang Lin, and Shuicheng Yan. Human parsing with contextualized convolutional neural network. In *ICCV*, 2015.
- [43] Si Liu, Jiashi Feng, Csaba Domokos, Hui Xu, Junshi Huang, Zhenzhen Hu, and Shuicheng Yan. Fashion parsing with weak color-category labels. *TMM*, 16(1):253–265, 2014.
- [44] Si Liu, Xiaodan Liang, Luoqi Liu, Xiaohui Shen, Jianchao Yang, Changsheng Xu, Liang Lin, Xiaochun Cao, and Shuicheng Yan. Matching-cnn meets knn: Quasi-parametric human parsing. In *CVPR*, 2015.
- [45] Si Liu, Yao Sun, Defa Zhu, Guanghui Ren, Yu Chen, Jiashi Feng, and Jizhong Han. Cross-domain human parsing via adversarial feature and label adaptation. In *AAAI*, 2018.
- [46] Jonathan Long, Evan Shelhamer, and Trevor Darrell. Fully convolutional networks for semantic segmentation. In *CVPR*, 2015.
- [47] Pauline Luc, Camille Couprie, Soumith Chintala, and Jakob Verbeek. Semantic segmentation using adversarial networks. In *NIPS-workshop*, 2016.
- [48] Ping Luo, Xiaogang Wang, and Xiaoou Tang. Pedestrian parsing via deep decomposition network. In *ICCV*, 2013.
- [49] Xianghui Luo, Zhuo Su, Jiaming Guo, Gengwei Zhang, and Xiangjian He. Trusted guidance pyramid network for human parsing. In *ACMMM*, 2018.
- [50] Yawei Luo, Zhedong Zheng, Liang Zheng, Tao Guan, Junqing Yu, and Yi Yang. Macro-micro adversarial network for human parsing. In *ECCV*, 2018.
- [51] Anthony J Marcel. Conscious and unconscious perception: An approach to the relations between phenomenal experience and perceptual processes. *Cognitive psychology*, 15(2):238–300, 1983.
- [52] Tom M Mitchell. *The need for biases in learning generalizations*. Department of Computer Science, Laboratory for Computer Science Research, 1980.
- [53] David Navon. Forest before trees: The precedence of global features in visual perception. *Cognitive psychology*, pages 353–383, 1977.
- [54] Xuecheng Nie, Jiashi Feng, and Shuicheng Yan. Mutual learning to adapt for joint human parsing and pose estimation. In *ECCV*, 2018.
- [55] Siyuan Qi, Wenguan Wang, Baoxiong Jia, Jianbing Shen, and Song-Chun Zhu. Learning human-object interactions by graph parsing neural networks. In *ECCV*, 2018.
- [56] 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*, 2019.
- [57] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael S. Bernstein, Alexander C. Berg, and Fei-Fei Li. Imagenet large scale visual recognition challenge. *IJCV*, 115(3):211–252, 2015.
- [58] Edgar Simo-Serra, Sanja Fidler, Francesc Moreno-Noguer, and Raquel Urtasun. A high performance crf model for clothes parsing. In *ACCV*, 2014.
- [59] Xi Song, Tianfu Wu, Yunde Jia, and Song-Chun Zhu. Discriminatively trained and-or tree models for object detection. In *CVPR*, 2013.
- [60] Charles Sutton, Andrew McCallum, et al. An introduction to conditional random fields. *Foundations and Trends® in Machine Learning*, pages 267–373, 2012.
- [61] Wei Tang, Pei Yu, and Ying Wu. Deeply learned compositional models for human pose estimation. In *ECCV*, 2018.
- [62] Michael J Tarr and Heinrich H Bülthoff. Image-based object recognition in man, monkey and machine. *Cognition*, pages 1–20, 1998.
- [63] Andreas Veit and Serge Belongie. Convolutional networks with adaptive inference graphs. In *ECCV*, 2018.
- [64] Martin J Wainwright, Michael I Jordan, et al. Graphical models, exponential families, and variational inference. *Foundations and Trends® in Machine Learning*, pages 1–305, 2008.
- [65] Nan Wang and Haizhou Ai. Who blocks who: Simultaneous clothing segmentation for grouping images. In *ICCV*, 2011.
- [66] Peng Wang, Xiaohui Shen, Zhe Lin, Scott Cohen, Brian Price, and Alan L Yuille. Joint object and part segmentation using deep learned potentials. In *ICCV*, 2015.
- [67] Wenguan Wang, Jianbing Shen, Ming-Ming Cheng, and Ling Shao. An iterative and cooperative top-down and bottom-up inference network for salient object detection. In *CVPR*, 2019.
- [68] Wenguan Wang, Yuanlu Xu, Jianbing Shen, and Song-Chun Zhu. Attentive fashion grammar network for fashion landmark detection and clothing category classification. In *CVPR*, 2018.
- [69] Tianfu Wu and Song-Chun Zhu. A numerical study of the bottom-up and top-down inference processes in and-or graphs. *IJCV*, 93(2):226–252, 2011.
- [70] Fangting Xia, Peng Wang, Liang-Chieh Chen, and Alan L Yuille. Zoom better to see clearer: Human and object parsing with hierarchical auto-zoom net. In *ECCV*, 2016.
- [71] Fangting Xia, Peng Wang, Xianjie Chen, and Alan L Yuille. Joint multi-person pose estimation and semantic part segmentation. In *CVPR*, 2017.
- [72] Fangting Xia, Jun Zhu, Peng Wang, and Alan L Yuille. Pose-guided human parsing by an and/or graph using pose-context features. In *AAAI*, 2016.
- [73] Kota Yamaguchi, M Hadi Kiapour, and Tamara L Berg. Paper doll parsing: Retrieving similar styles to parse clothing items. In *ICCV*, 2013.
- [74] Kota Yamaguchi, M Hadi Kiapour, Luis E Ortiz, and Tamara L Berg. Parsing clothing in fashion photographs. In *CVPR*, 2012.
- [75] Wei Yang, Ping Luo, and Liang Lin. Clothing co-parsing by joint image segmentation and labeling. In *CVPR*, 2014.
- [76] Hengshuang Zhao, Jianping Shi, Xiaojuan Qi, Xiaogang Wang, and Jiaya Jia. Pyramid scene parsing network. In *CVPR*, 2017.

- [77] Jian Zhao, Jianshu Li, Yu Cheng, Terence Sim, Shuicheng Yan, and Jiashi Feng. Understanding humans in crowded scenes: Deep nested adversarial learning and a new benchmark for multi-human parsing. In *ACMMM*, 2018.
- [78] Jian Zhao, Jianshu Li, Xuecheng Nie, Fang Zhao, Yunpeng Chen, Zhecan Wang, Jiashi Feng, and Shuicheng Yan. Self-supervised neural aggregation networks for human parsing. In *CVPR-workshop*, 2017.
- [79] Shuai Zheng, Sadeep Jayasumana, Bernardino Romera-Paredes, Vibhav Vineet, Zhizhong Su, Dalong Du, Chang Huang, and Philip HS Torr. Conditional random fields as recurrent neural networks. In *ICCV*, 2015.
- [80] Qixian Zhou, Xiaodan Liang, Ke Gong, and Liang Lin. Adaptive temporal encoding network for video instance-level human parsing. In *ACMMM*, 2018.
- [81] Bingke Zhu, Yingying Chen, Ming Tang, and Jinqiao Wang. Progressive cognitive human parsing. In *AAAI*, 2018.
- [82] Shizhan Zhu, Raquel Urtasun, Sanja Fidler, Dahua Lin, and Chen Change Loy. Be your own prada: Fashion synthesis with structural coherence. In *ICCV*, 2017.