

AutoPose: Searching Multi-Scale Branch Aggregation for Pose Estimation

Xinyu Gong^{*1} Wuyang Chen^{*1} Yifan Jiang¹ Ye Yuan²
Xianming Liu³ Qian Zhang³ Yuan Li³ Zhangyang Wang¹

¹ Department of Electrical and Computer Engineering, The University of Texas at Austin

²Department of Computer Science & Engineering, Texas A&M University ³Horizon Robotics Inc.

{xinyu.gong, wuyang.chen, yifanjiang97, atlaswang}@utexas.edu
ye.yuan@tamu.edu, {xianming.liu, qian01.zhang, yuan.li}@horizon.ai

Abstract

We present AutoPose, a novel neural architecture search (NAS) framework that is capable of automatically discovering multiple parallel branches of cross-scale connections towards accurate and high-resolution 2D human pose estimation. Recently, high-performance hand-crafted convolutional networks for pose estimation show growing demands on multi-scale fusion and high-resolution representations. However, current NAS works exhibit limited flexibility on scale searching, they dominantly adopt simplified search spaces of single-branch architectures. Such simplification limits the fusion of information at different scales and fails to maintain high-resolution representations. The presented AutoPose framework is able to search for multi-branch scales and network depth, in addition to the cell-level micro structure. Motivated by the search space, a novel bi-level optimization method is presented, where the network-level architecture is searched via reinforcement learning, and the cell-level search is conducted by the gradient-based method. Within 2.5 GPU days, AutoPose is able to find very competitive architectures on the MS COCO dataset, that are also transferable to the MPII dataset. Our code is available at <https://github.com/VITA-Group/AutoPose>.

1. Introduction

Human 2D pose estimation targets at localizing anatomical keypoints of individuals from images. Accurate pose estimation plays an essential role in understanding human behaviors. This paper focuses on pose estimation for single person, which can be a cornerstone for downstream tasks including action recognition [34, 35, 36], gaming [19], multi-person pose estimation [3], and video tracking [32, 44], etc.

Recently, great endeavors have been made to exploit convolutional neural networks (CNN) to improve pose estimation.

Especially, multi-scale branches and context aggregation modules are increasingly leveraged to capture multi-grained features. For example, Newell *et al.* proposes the stacked Hourglass network [27], which repeatedly down-samples and restores resolutions through the network. Cascaded pyramid network [7] constructs a multi-scale network, keeping multi-scale information through the network and aggregating them at the end. HRNet [38] achieves state-of-the-art performance by constructing a multi-branch network, where high-resolution representations are maintained through the whole pipeline, in addition to multi-scale aggregation at each stage. While multi-scale features, as well as high resolution, are shown to be vital for accurate pose estimation, how to advance the design of pose estimation networks further is not immediately clear, since the arising complex branches and connectivity patterns can cost heavy human efforts to explore thoroughly.

Designing modern CNN architectures manually can result in tedious trial-and-errors. To this end, the neural architecture search (NAS) has recently drawn a booming interest. Aiming to discover an optimal network architecture from data, NAS has been successfully applied primarily on image classification [10, 22, 25, 30, 33, 43, 49], and lately also on object detection [8, 13], semantic segmentation [21, 26, 6], person re-identification [31], speech recognition [9], super-resolution [37], medical image analysis [48], and even generative models [14, 12] or Bayesian deep networks [2]. These works surge heavy interests in designing sophisticated architectures automatically, reducing human efforts as well as revealing/verifying design insights.

Despite the promise of NAS, there remains to be a major gap between its current capability, and the real sophistication needed by advanced computer vision applications. For example, for the sake of simplicity, most NAS methods usually define their search spaces in a way that can only cover chain-like single backbones (with skip connections), which will exclude the discovery of multiple (and inter-connected) backbones that have proven effectiveness [7, 38]. Recently,

^{*}Part of this work was done while Xinyu and Wuyang were Research Interns at Horizon Robotics Applied AI Lab.

Table 1. Comparison of AutoPose against other NAS works. AutoPose is the first NAS framework that integrally searches for cell, scale, aggregation (“Agg.”) and branches. We also compares the size of search space (“SS”), optimization method, and search time (GPU days).

Model	SS	Optimization	Cell	Scale	Agg.	Branch	Dataset	Task	Days
NASNet [50]	10^{16}	RL	✓				CIFAR-10	Cls	2000
DARTS [22]	10^{28}	GRAD	✓				CIFAR-10	Cls	4
DPC [5]	10^{11}	RS			✓		Cityscapes	Seg	2600
Auto-DeepLab [21]	10^{19}	GRAD	✓	✓			Cityscapes	Seg	3
Ours	10^{61}	GRAD+RL	✓	✓	✓	✓	MSCOCO	Pose	2.5

Auto-DeepLab [21] proposes a two-level hierarchy search space for semantic segmentation, which is able to search for the network scales at different stages. However, considering only connections within a single cell, Auto-DeepLab cannot derive an architecture that contains connections over multiple sub-paths over different scales at the same time. NAS-FPN [13] designs cross-scale connections only for the feature fusion head, while the backbone is not optimized.

This paper proposes AutoPose, a novel NAS framework that automatically designs architectures towards accurate 2D human pose estimation. The core of AutoPose is a hierarchical multi-scale **search space**, consisting of a novel network-level search space in addition to the cell-level search space. Our search space can simultaneously maintain multiple parallel branches of diverse scales in an end-to-end fashion throughout the network, which leads to much flexible and finer-grained exploration. To effectively fuse features from branches, AutoPose further searches for novel aggregation cells, which support dense and sophisticated cross-scale aggregations, instead of naively summing up all features [38].

As a result, the AutoPose search space is unusually gigantic (10^{61} candidates) with multi-level selections (from cell to network), making either gradient method or reinforcement learning alone *not directly applicable*. For example, traditional gradient-based search is unable to select more than one branches simultaneously [21]. On the other hand, reinforcement learning would be too time-consuming for such a huge space. Motivated by that, we present a novel bi-level **optimization algorithm** combining the two, for different aspects of search space. Gradient-based method is leveraged to search for the micro cell architectures (cell-level). An RNN controller trained via reinforcement learning then guides the dynamic searching over multiple scales and network depth (network level), which boosts the flexibility in discovering networks of capacity variations.

We summarize our contributions as follows:

- **High-resolution Multi-branch Search Space.** We show that: i) the searchable parallel branches of multiple scales; and ii) the cross-scale context aggregations, are two key components towards searching for accurate human pose estimation. Extensive ablation studies support our design philosophy.

- **Bi-level Optimization.** To fulfill the effective discovery over the gigantic fined-grained search space, we integrate the reinforcement learning and the gradient-based search method for the first time, to jointly optimize the network-level and the cell-level architectures.

- Experiments demonstrate that AutoPose for the first time automatically designs architectures of sophisticated connections, achieving competitive performance as hand-crafted networks for human pose estimation.

We noticed another recently proposed NAS framework of pose estimation called PoseNFS [45]. Unlike [21], PoseNFS follows a one-together search strategy directly over macro and micro levels. To simplify their search, the authors referred to the body structure prior knowledge to decompose the whole search space into multiple part-specific ones, with global pose constraint relationship enforced. However, its performance falls behind current state-of-the-art hand-crafted pose estimation models with a notable margin. In comparison, we decouple our search at two different granularity levels in one unified space, and show the feasibility and effectiveness of purely end-to-end data-driven search without drawing any explicit body part knowledge. We compare our results with PoseNFS in Section 4.

2. Related Work

2.1. Human Pose Estimation

Existing pose estimation work can be categorized to either bottom-up or top-down. Bottom-up methods [3, 18] first detect all parts in the image (i.e. parts of every person), then associate or group parts belonging to each individual. Bottom-up methods are more suitable for efficient pose estimation. Top-down methods [7, 11, 15, 17, 28, 38] first generate human proposals in the form of bounding boxes and then exploit non-maximum suppression (NMS) to remove the redundant proposals, after which a single-person pose estimation is conducted for every proposal.

The recent trend in pose estimation leverages high-resolution features and aggregating contexts from multiple scales. Chen *et al.* [7] proposes the Cascaded Pyramid Network which consists of a GlobalNet based on feature pyramid and a cascaded multi-scale RefineNet to focus more on fine-granularity hard keypoints. HRNet [38] achieves

state-of-the-art performance by utilizing a network of multiple branches across multiple scales, where information is aggregated at each stage of the network. In AutoPose, we follow this line of leveraging high-resolution features across multi-scale branches from end to end.

2.2. Neural Architecture Search Spaces

Most existing works in neural architecture search (NAS) focus on searching either stacked operators [30, 49] or repeated cell-structured directed acyclic graph [22] for the classification task. Search spaces in these works are mainly designed for chain-like single-path networks with fixed scale patterns. Recent NAS works explore more challenging high-level vision tasks like object detection [8, 13] and semantic segmentation [21, 26]. Although these dense prediction tasks are more sensitive to resolutions and require aggregation of multi-scale context for improving performance, these early NAS works for dense prediction tasks are still limited in their searchable scale-level connections. In contrast, our search space supports more complex cross-scale connections throughout the network. Different from the network-level search in Auto-Deeplab [21] which still only explore and derive only one branch, our AutoPose supports maintaining multiple parallel branches across scales. Our search space naturally provides the architecture exploration with more flexibility yet it is more challenging.

2.3. Neural Architecture Search Methods

Neural architecture optimization is first proposed with the reinforcement learning method [30, 49], which leverages a learned policy to control the selection of operators along with the network. Reinforcement learning is a popular methods for NAS works which involve non-differentiable optimization target [14, 39], as the policy-gradient method [41] can leverage discrete extrinsic rewards. The gradient-based method is proposed in [22] which provides an efficient way of searching the architecture using gradient descent [21]. Real *et al.* [33] successfully introduce an evolution algorithm into NAS, achieving competitive performance, also followed by many [47]. Egrinho *et al.* [25] takes the first step to incorporate the Monte Carlo tree search (MCTS) into NAS. [40] also adopted the search space divide-and-conquer strategy. NAO [23] embeds the complex network structures into a latent space to simplify the searching. All existing NAS works focus on using one search method coupled with their search spaces. However, a single search method may lack the capability of supporting complex search spaces. In our AutoPose, we for the first time integrate different methods for different parts of our search space.

3. Method

We illustrate the overview of our AutoPose framework in Fig. 1. The input image is first passed through a pre-

defined stem module (Sec. 4.2) to output feature maps of different scales*, which are served as inputs into our multi-branch search space. We split our search space into three stages, where the activated parallel branches (of different scales) in each stage are individually controlled by our RNN controller. The purpose of these separate stages is to support the different needs of scales from early to late stage along with the network. At the end, the outputs of different branches are upsampled via nearest-neighbor interpolation to the highest resolution, and are aggregated by summation.

3.1. High-resolution Multi-branch Search Space

Parallel branches of diverse scales and context aggregation modules are two key points for sufficiently capturing feature maps under different receptive fields, which motivates the design of our core functional unit “Multi-scale Module” (Fig. 2). In the following sections, we describe our network-level search space for maintaining multiple parallel branches throughout the network and cell-level structures for merging multi-scale features. We by default stack two identical multi-scale modules in each stage.

3.1.1 Network-level Search Space

Our core component “Multi-scale Module” (Fig. 2) is composed of B branches of different scales, where each branch includes two individual cells ($\mathcal{C}_{I1}, \mathcal{C}_{I2}$) and one aggregation cell (\mathcal{C}_M) (see Sec. 3.1.2). The activation status of each scale in each stage is individually controlled by our RNN controller. Specifically, each stage may activate one or more different scales. If a scale is activated, the input tensor of that scale will be processed by the two individual cells on that branch. Otherwise, the input tensor will go through a skip connection, which means the two individual cells and the aggregation cell on this deactivated branch will not participate in the network forward and backward process, and correspondingly the cell weights and architecture will not be updated by the gradient descent. The feature maps remain the same scale within each branch.

More formally, for the m -th multi-scale module, the RNN controller outputs a binary scale vector $[\beta_{2^2}^m, \beta_{2^3}^m, \dots, \beta_{2^{B+1}}^m]$. For the b -th branch (top to bottom in Fig. 1, $b \in N^+, 1 \leq b \leq B$) whose scale $s_b = 2^{b+1}$, its output, denoted as $H_{s_b}^m$, can be formulated as:

$$H_{s_b}^m = \begin{cases} \mathcal{C}_{M,s_b}^m \circ \mathcal{C}_{I2,s_b}^m \circ \mathcal{C}_{I1,s_b}^m (H_{s_b}^{m-1}) & \text{if } \beta_{s_b}^m = 1, \\ H_{s_b}^{m-1} & \text{if } \beta_{s_b}^m = 0. \end{cases} \quad (1)$$

Note that we allow the situation that all branches in a stage are deactivated, which means that our framework is able to search for the **network depth**. Fig. 3 shows an example when $B = 4$ and the controller generates a scale vector $[1, 0, 1, 0]$. In this scenario, only the branches of $4 \times$

*A feature map of scale s is of $\frac{1}{s^2}$ size of the input image by scaling both the height and width to $\frac{1}{s}$.

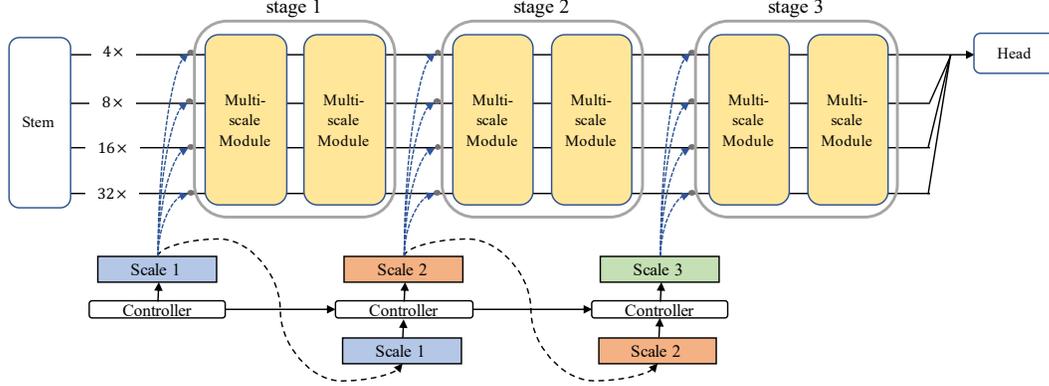


Figure 1. Our AutoPose framework. We support the multi-scale parallel branches and context aggregation via our stacked multi-scale modules. The activation of each scale ($\{4\times, 8\times, 16\times, 32\times\}$ when $B = 4$ shown here) in each stage is controlled by an RNN controller.

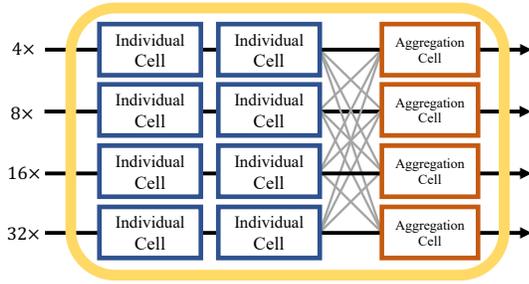


Figure 2. Our multi-scale module takes inputs of different scales. The inputs will first go through two individual cells and then will be aggregated together by the aggregation cells.

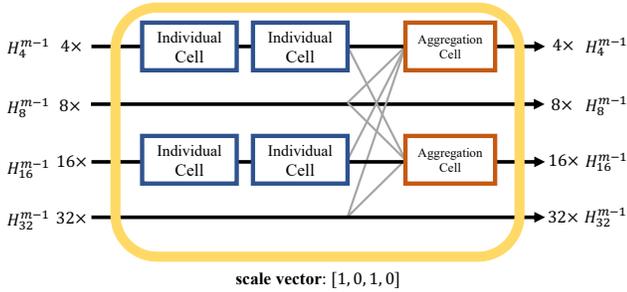


Figure 3. An example of a multi-scale module with a scale vector $[1, 0, 1, 0]$. Only the branches of scales $4\times$ and $16\times$ are activated.

and $16\times$ scales are activated and will process their input tensors, while the branches of $8\times$ and $32\times$ scales become skip connections. Note that the aggregation cells of $4\times$ and $16\times$ scales will take $\mathcal{C}_{I2,s_b}^m \circ \mathcal{C}_{I1,s_b}^m (H_{s_b}^{m-1})$ as inputs for $s_b = 4, 16$, and will take H_s^{m-1} as inputs for $s_b = 8, 32$.

3.1.2 Cell-level Search Space

We include two types of cells ((network basic units): the individual cell and the aggregation cell. As in Fig. 2, on each branch in the multi-scale module, two individual cells ($\mathcal{C}_{I1}, \mathcal{C}_{I2}$) are sequentially stacked, followed by a aggregation cell (\mathcal{C}_M). Both the individual cell \mathcal{C}_I and the aggregation cell \mathcal{C}_M are directed acyclic graphs (DAG), consist-

ing of an ordered sequence of K_I, K_M nodes, respectively. Each node $x^{(i)}$ is a feature map and each directed edge (i, j) is associated with a specific operation $O^{(i,j)}$ that transforms node. The output of a cell, denoted as H , is the concatenation of all nodes, followed by a 1×1 convolution, aiming at compressing the channel dimension.

Individual Cell An individual cell \mathcal{C}_I takes its predecessor (the previous cell's output) as its single input, and outputs one tensor. The i -th node $x_I^{(i)}$ is calculated as the sum of the transformed (processed by an operation) i previous nodes and the input of the cell. The computation of i -th intermediate node in the individual cell is formulated as:

$$x_I^{(i)} = \sum_{j < i} O^{(i,j)}(x_I^{(j)}). \quad (2)$$

$O^{(i,j)}$ denotes the operation of i -th node, which takes the j -th node $x_I^{(j)}$ as input.

Aggregation Cell A aggregation cell \mathcal{C}_M takes B outputs of different scales from B individual cells as its inputs. In particular, the B inputs from the individual cells will be first downsampled or upsampled accordingly to match the current aggregation cell's scale. The computation at i -th intermediate node in the aggregation cell is formulated as:

$$x_M^{(i)} = \sum_{j < i} O^{(i,j)}(x_M^{(j)}) + \sum_{b=1}^B O_b^i(H_{I2,s_b}), \quad (3)$$

where H_{I2,s_b} denotes the output of the second individual cell of scale $s_b = 2^{b+1}$ ($b \in N^+, 1 \leq b \leq B$), with associated operation O_b^i for the i -th node.

Continuous Relaxation of Cell-level Search Space We adopts the continuous relaxation [10, 21, 22] to represent the cell-level architecture. Specifically, the categorical choice of a particular operation $O^{(i,j)}$ is relaxed to be

continuous:

$$\bar{O}^{(i,j)}(x^{(j)}) = \sum_{k=1}^{|\mathcal{O}|} \alpha_k^{(i,j)} O_k^{(i,j)}(x^{(j)}) \quad (4)$$

where $\alpha^{(i,j)}$ is a differentiable $|\mathcal{O}|$ -dim vector with its k -th elements associated with operation $O^k \in \mathcal{O}$ for the edge from j -th node to i -th node. $\alpha_k^{(i,j)}$ is normalized by the ‘‘Gumbel-Softmax’’ with reparameterization trick, for the purpose of efficient search process (see supplementary materials). \mathcal{O} is the set of candidate operations, which contains six prevalent operations as below:

- 1×1 convolution
- 3×3 average pooling
- 3×3 convolution
- 3×3 max pooling
- 5×5 convolution
- skip connection

3.2. Bi-level Optimization

To enable the optimization over network-level search space, one ad-hoc idea may be also adopting a continuously relaxed network-level search space and applying gradient-based updates, similar to the approach proposed by Liu *et al.* [21]. However, this continuous relaxation cannot support the need for multiple scales. As demonstrated by the single-branch network in Auto-DeepLab, the gradient-based search method can only pick the top1 single branch, which is not naturally applicable to discovering parallel branches.

We instead unify the optimization of the network-level and cell-level architectures, and train a policy to control the selection of cross-scale connections. Specifically, we update the cell-level micro architecture via gradient descent, and our policy controller will sample the macro architecture of scales, where the network will activate the corresponding branches. The network-level macro search and cell-level micro search are integrated and alternatively proceed (see Sec. 4.2 for searching details). The search procedure is shown in the *search* function in Algorithm 1.

3.2.1 Network-level Macro Search

For the network-level macro search, we construct an RNN controller optimized via the REINFORCE gradient [41], to decide which branches should be activated. The optimization of controller is illustrated in *controller_update* function in Algo. 1. At each time step, the controller will sample a B -dimension vector from the policy $\pi(\theta)$, where the b -th element indicates the activation status of the branch of scale 2^{b+1} . We use average precision (AP) score on validation dataset as the reward \mathcal{R} .

3.2.2 Cell-level Micro Search

With continuous relaxation mentioned in Sec. 3.1.2, we are able to optimize the discrete architecture with gradient descent. We adopt the first-order approximation in [22]. The

Algorithm 1: Bi-level Optimization

Function *search* ($N_{epoch}, N_{iter}, N_{ctr}$):
 initialize network weights ω , cell architecture α , branch architecture β , policy controller $\pi(\theta)$
 $epoch = 0$
while $epoch < N_{epoch}$ **do**
 $iter = 0$
 while $iter < N_{iter}$ **do**
 $\beta \leftarrow \pi(\theta)$
 cell.update(β)
 $iter = iter + 1$
 end
 controller.update(N_{ctr})
 $epoch = epoch + 1$
end

Function *cell.update* (β):
 Continuously relax operation $\bar{o}^{(i,j)}$ weighted by the architecture parameter $\alpha^{(i,j)}$
 1. Update weights ω by descending $\nabla_{\omega} \mathcal{L}_{train}(\omega, \alpha, \beta)$
 2. Update cell architecture α by descending $\nabla_{\alpha} \mathcal{L}_{val}(\omega, \alpha, \beta)$
 3. Replace $\bar{o}^{(i,j)}$ with $o_{k^*}^{(i,j)}, k^* = \operatorname{argmax}_{1 \leq k \leq |\mathcal{O}|} \alpha_k^{(i,j)}$

Function *controller.update* (N_{ctr}):
 Fix network weights ω and cell architecture α
 $step = 0$
while $step < N_{ctr}$ **do**
 Update controller by maximize the expected reward
 $\mathbb{E}_{\beta \sim \pi(\theta)} [\mathcal{R}(\omega, \alpha, \beta)]$
 $step = step + 1$
end

loss function \mathcal{L} is the mean squared error calculated between the regressed heatmap and the ground truth. The search algorithm proceeds the update of supernet weights ω and the cell-level architecture α in an alternative way, shown in function *cell_update* in Algo. 1.

3.3. Architecture Derivation

Network-level Architecture We adopt the optimized policy controller to sample the scale vector of each stage for N times, then evaluate each status of parallel branches on the validation set by using the same shared weights in supernet. The one with the best performance will be selected as the discovered network-level architecture. We choose $N = 10$ by default, and a larger N can be adopted to sample more status of parallel branches.

Cell-level Architecture The optimal operators in cell-level architecture are derived by taking the argmax over $\alpha^{(i,j)}$. Different from previous works [21, 22], we preserve the connection from all predecessors for every node.

4. Experiments

In all experiments, we use Nvidia Titan RTX for benchmarking the computing power. All experiments are performed under CUDA 9.0 and CUDNN V7, using Pytorch.

4.1. Datasets

COCO Dataset We adopt the COCO dataset [20] as our main testbed, which contains more than 200k images and 250k person instances with 17 keypoints per instance. Our models are trained on the training dataset only, including 56k images and 150k person instances. During training, each image is resized to 256×192 . We apply random flipping and rotation ($[-45^\circ, 45^\circ]$) as data augmentation. The validation set contains about 5k images. Results are reported using the official evaluation metric of the COCO Keypoint Challenge: average precision (AP) and average recall (AR), based on object keypoints similarity (OKS):

$$\text{OKS} = \frac{\sum_i [\exp(-d_i^2/2s^2k_i^2)\delta(v_i > 0)]}{\sum_i [\delta(v_i > 0)]} \quad (5)$$

Here, d_i is the Euclidean distance between ground-truth and detected keypoint, v_i is the visibility flag of the ground-truth keypoint. Each kind of keypoint also has a scale s which is defined as the square root of the object segment area. k_i is a predefined constant provided by the COCO Keypoint Challenge 2017 which is used to control falloff.

MPII Dataset The MPII dataset [1] contains around 25k images with over 40k people with annotated body joints. Note that we use MPII to study the transferability of COCO-search architectures in this paper, rather than conducting any search on MPII directly. We trained our discovered model on training split only, and evaluate the single person pose estimation result on the validation dataset. We adopt the same data augmentation as on the COCO dataset. The input image size is set to 256×256 . The evaluation metric is the Percentage of Correct Keypoints (PCK) [46].

4.2. Architecture Search and Implementations

Stem and Head Structure In both our supernet and derived architecture, from the input image we leverage a pre-defined stem module to extract feature maps of different scales serving as inputs into the multi-scale modules. We design the “stem” structure with two 3×3 convolutions (with stride 2) followed by a residual bottleneck [16], which produces a feature map of $4\times$ scale. We sequentially stack more 3×3 convolutions with stride 2 to provide feature maps of higher scales. At the end of the network, the outputs of different branches are upsampled via nearest-neighbor interpolation to the highest resolution (“ $4\times$ ” branch in our case), and are aggregated with a concatenation. The concatenated feature map is then processed by a 3×3 convolution of stride 1 towards the final output.

Searching Details The number of branches B is set to 4, where the scales of branches are $\{4\times, 8\times, 16\times, 32\times\}$ and the number of channels are 16, 32, 64, 128, respectively. The number of nodes in each individual cell is set as $K_I = 4$, while for the aggregation cell $K_M = 2$. Each stage has 2 multi-scale modules, where all stages will share the same cell-level architecture parameters $\alpha^{(i,j)}$, following previous works [21, 22, 50]. We adopt Adam optimizer to train the cell-level architecture parameters $\alpha^{(i,j)}$ with a learning rate of $3e-3$. During our search we set the total number of epochs $N_{epoch} = 50$, and $N_{iter} = 1100$ iterations per epoch. At the end of each epoch for the cell-level search, the scale controller will be updated for $N_{ctr} = 10$ times, with a batch size of 2. The controller is also optimized by an Adam optimizer, with a learning rate of $3.5e-3$. For the reward to our controller, we use the OKS-based AP score on the validation dataset. Note that we start training the controller at the 5-th epochs, to avoid the bias caused by the under-trained supernet.

Training Details The number of channels in the derived architecture is doubled from the supernet used in the search. The network is trained for 210 epochs. The initial learning rate is set to 0.001, which will be multiplied by a decay factor of 0.1 at epoch 170 and epoch 200. Adam optimizer is adopted for the network, with $\beta_1 = 0.9$, $\beta_2 = 0.99$. The batch size is set to 128.

Testing Details We borrow the person detector from HTC [4] to get the person bounding box for COCO dataset testing. We follow the test procedure in [27], where the output heatmap is computed by averaging the output heatmap of the original and its flipped one.

4.3. Results

Fig. 4 visualizes the best network-level architecture discovered by our AutoPose framework. Diversified scales are preferred at the beginning. Branches with different scales receive balanced activations, with feature maps of lower resolutions ($8\times, 16\times, 32\times$) are processed earlier and features of the highest resolution ($4\times$) are later processed. Please see our supplementary materials for the details of the searched cell structures.

Results on COCO Tab. 2 shows the performance of the discovered architecture on COCO validation dataset, compared with other state-of-the-art methods. Without any pre-training, our AP score significantly outperforms 8-stage Hourglass [27] by 5.6 and CPN [7] by 5, with the same input size. We surpass the SimpleBaseline [42] by 1.6 with 32.1% less GFLOPs. Compared with the state-of-the-art HRNets, our model achieves competitive AP scores. AutoPose also shows clearly superior performance over the previous NAS for pose estimation work, PoseNFS [45], with lower input image resolution and FLOPs.

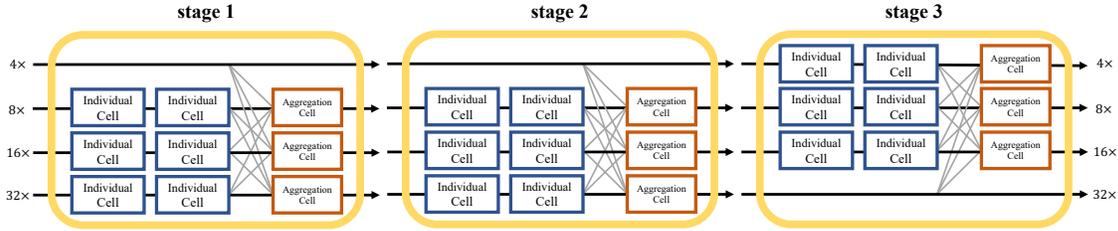


Figure 4. The discovered network-level macro architecture of our proposed AutoPose framework.

Table 2. Comparison on the COCO validation set. The symbol † indicates the result obtained by reproducing with the author’s official released code, as the original paper did not report the same setting. Models in the first block at table are hand-crafted models, while models in the second block are obtained by neural architecture search.

Method	Pretrain	Image size	GFLOPs	AP	AP ⁵⁰	AP ⁷⁵	AP ^M	AP ^L	AR
8-stage Hourglass[27]	N	256 × 192	14.30	66.9	-	-	-	-	-
8-stage Hourglass[27]	N	256 × 256	-	67.1	-	-	-	-	-
CPN[7]	Y	256 × 192	6.20	68.6	-	-	-	-	-
Posefix[24]	Y	256 × 192	-	71.5	88.0	77.6	68.0	78.1	-
SimpleBaseline-50[42]	Y	256 × 192	8.90	70.4	88.6	78.3	67.1	77.2	76.3
SimpleBaseline-101[42]	Y	256 × 192	12.40	71.4	89.3	79.3	68.1	78.1	77.1
SimpleBaseline-152[42]	Y	256 × 192	15.70	72.0	89.3	79.8	68.7	78.9	77.8
HRNet-W32[38]	N	256 × 192	7.10	73.4	89.5	80.7	70.2	80.1	78.9
HRNet-W48[38] †	N	256 × 192	14.60	73.5	89.5	80.6	70.4	79.8	79.2
PoseNFS-3 [45]	Y	384 × 288	14.8	73.0	-	-	-	-	-
AutoPose (Ours)	N	256 × 192	10.65	73.6	90.6	80.1	69.8	79.7	78.1

Results on MPII We directly transfer the discovered architecture on the COCO dataset to train on the MPII dataset. We note that images from the MPII dataset are larger (256 × 256) than those in COCO (256 × 192), which is to our discovered architecture’s disadvantage. Without any pre-training, the architecture designed by AutoPose surpasses “ASR+AHO” [29]. We found that recent works (Simple-Baseline [42] and HRNet [38]) adopt ImageNet pretraining and enjoy improved performance. Further discussion with the authors of HRNet [38] helped us confirm that adopting a pretraining backbone is vital to the performance boost and convergence for complicated pose estimation network, especially on relatively small datasets like MPII.

This reveals a challenging issue when extending NAS to advanced computer vision tasks: as automatically designed architectures will require extra efforts to be pre-trained on ImageNet, training from scratch is an unfair setting to NAS when compared with human-designed networks which adopt mature backbones pre-trained on extra data (e.g. ResNet101 [16]). One remedy might be still leveraging a pretrained hand-crafted backbone followed by searched aggregation [5, 13]; but that will limit the flexibility of NAS and somehow contradict to the philosophy of AutoML. We leave this open challenge for future work.

4.4. Ablation Study

Multi-Scale Search Space Recently, Liu *et al.* proposed Auto-Deeplab [21] to search for both the cell-level and network-level architecture. However, their search architecture only allowed one scale to be activated per layer.

Table 3. Transferability results on MPII validation dataset. No multi-scale testing technique is applied on our result.

Method	Pretrain	Head	Shoulder	Elbow	Wrist	Hip	Knee	Ankle	Total
“ASR+AHO” [29]	Y	97.3	95.1	88.7	84.7	88.4	82.5	78.1	87.8
Hourglass[27]	N	96.5	96.0	90.3	85.4	88.8	85.0	81.9	89.2
SimpleBaseline[42]	Y	97.0	95.9	90.3	85.0	89.2	85.3	81.3	89.6
HRNet-W32[38]	Y	97.1	95.9	90.3	86.4	89.1	87.1	83.3	90.3
AutoPose (Ours)	N	96.6	95.0	88.3	83.2	87.2	82.8	78.9	88.0

Table 4. Evaluating our multi-scale parallel branches and aggregation cells on the COCO validation set.

Aggregation Cell	Multi-branch	AP	AP ⁵⁰	AP ⁷⁵	AP ^M	AP ^L	AR
searched	searched (single)	63.9	86.5	70.2	61.2	68.9	69.0
N/A	searched	70.8	90.4	77.6	67.2	77.0	75.3
searched	fixed to four	71.8	89.5	79.4	69.6	75.5	74.8
searched	searched	73.6	90.6	80.1	69.8	79.7	78.1

To demonstrate the effectiveness of the proposed multi-scale search space, we construct two variants of our network. First, we limit the maximum number of activated branch in each stage to one during search, i.e., degrading our parallel branches back to a single chain-like backbone. Second, we allow our multi-scale branches, but remove the context aggregations in our aggregation cells and instead adopt plain concatenation operations. Third, we manually keep all branches in all stages to be activated, and only search for the cell-level architecture.

The ablation results are shown in Tab. 4. When we only search one single scale per stage, the AP score dramatically drops to 63.9%, demonstrating the key contribution of maintaining multiple scales throughout the network. When we activate multiple scales but bypass our aggregation cells, although the average precision under OKS = 0.5 achieves the similar bar (90.4% v.s. 90.6%), large OKS (0.75) will drop, implying the cross-scale context aggrega-

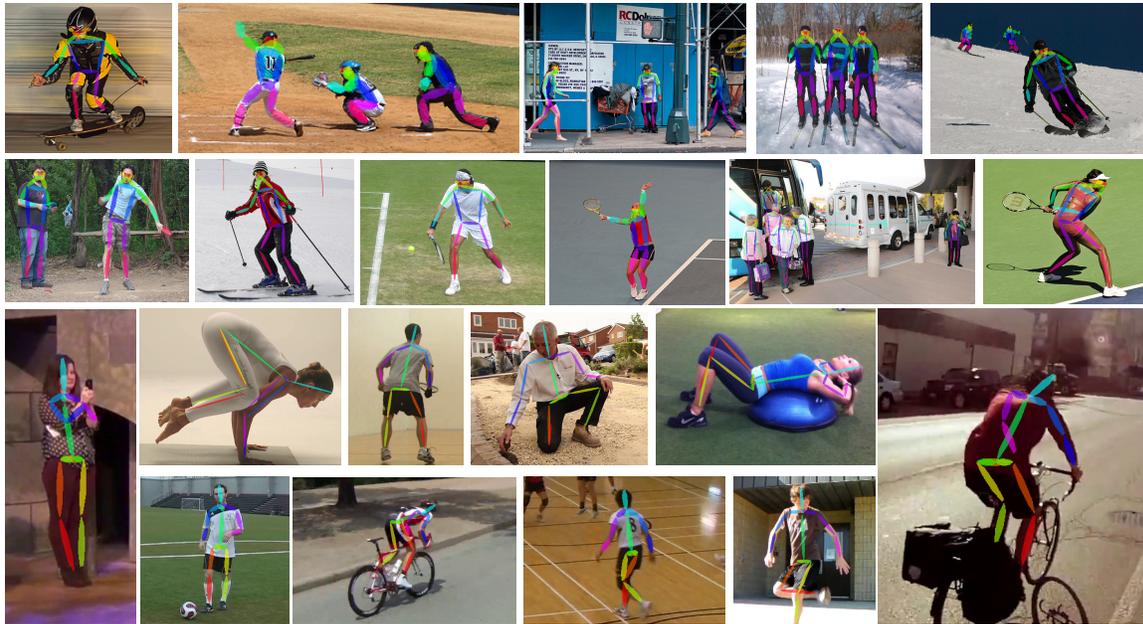


Figure 5. The visualizations of predicted keypoints. The first and last two rows are the results on COCO and MPII, respectively.

tion to be vital here. When we manually activate all four branches across all stages during search and evaluation, and only search for the micro cell-level structure, the performance remains to be lower than our proposed method.

Bi-level Optimization As existing NAS works only adopt a single search algorithm throughout their framework, we are the first to leverage a bi-level optimization for different aspects of our search space. Fig. 6 visualizes the convergences and performance during our search process. We can see that the validation accuracy steadily improves throughout the search process, while both the cell-level and network-level optimizations are converged to a small subset of the search space, indicated by the decreasing entropies.

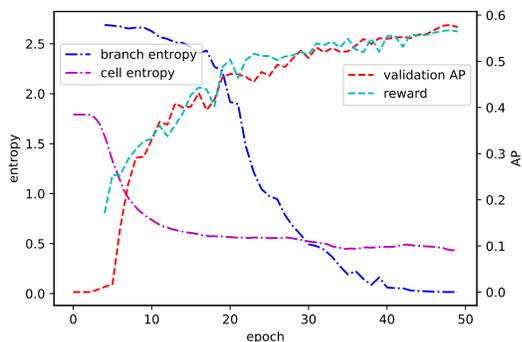


Figure 6. During the search process, the entropies of both cell-level and network-level architecture are decreasing (i.e. the supernet is becoming confident in its derived architecture), while the validation AP and controller’s reward keep increasing (i.e. the performance of our supernet is being improved). Note that we begin to train the controller at 5-th epoch.

Image Size during searching Previous work [21] reduced the input image size during the search phase to accelerate NAS. However, we argue that for dense image prediction tasks like pose estimation which requires fine details in context, adopting reduced input size during search amplifies the gap between the search and the train from scratch processes. We conduct an experiment where we use the size of 128×96 during searching, while still keep the size of 256×192 during training from scratch. In Tab. 5, we can see an obvious performance drop when the image size is reduced during searching.

Table 5. Impact of image sizes during search (evaluated on the COCO validation set).

Image size		AP	AP ⁵⁰	AP ⁷⁵	AP ^M	AP ^L	AR
Search	Train						
128×96	256×192	70.1	89.6	75.9	66.1	75.8	74.5
256×192	256×192	73.6	90.6	80.1	69.8	79.7	78.1

5. Conclusion

We propose AutoPose for the 2D human pose estimation task. Our high-capacity multi-scale search space supports the discovery of multiple parallel branches of diverse scales. By leveraging both gradient-based cell-level search and reinforcement learning based network-level search, our AutoPose successfully discovers the architecture that achieves very competitive performance on COCO and MPII dataset. We believe that pursuing an optimal combination of both cell-level and network-level structure is vital to accurate pose estimation, and wish our work to further advocate more NAS practice in advanced computer vision applications, with more reflections in the pros and cons.

References

- [1] Mykhaylo Andriluka, Leonid Pishchulin, Peter Gehler, and Bernt Schiele. 2d human pose estimation: New benchmark and state of the art analysis. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2014.
- [2] Randy Ardywibowo, Shahin Boluki, Xinyu Gong, Zhangyang Wang, and Xiaoning Qian. Nads: Neural architecture distribution search for uncertainty awareness. *arXiv preprint arXiv:2006.06646*, 2020.
- [3] Zhe Cao, Tomas Simon, Shih-En Wei, and Yaser Sheikh. Realtime multi-person 2d pose estimation using part affinity fields. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 7291–7299, 2017.
- [4] Kai Chen, Jiangmiao Pang, Jiaqi Wang, Yu Xiong, Xiaoxiao Li, Shuyang Sun, Wansen Feng, Ziwei Liu, Jianping Shi, Wanli Ouyang, Chen Change Loy, and Dahua Lin. Hybrid task cascade for instance segmentation. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2019.
- [5] Liang-Chieh Chen, Maxwell Collins, Yukun Zhu, George Papandreou, Barret Zoph, Florian Schroff, Hartwig Adam, and Jon Shlens. Searching for efficient multi-scale architectures for dense image prediction. In *Advances in Neural Information Processing Systems*, pages 8699–8710, 2018.
- [6] Wuyang Chen, Xinyu Gong, Xianming Liu, Qian Zhang, Yuan Li, and Zhangyang Wang. Fasterseg: Searching for faster real-time semantic segmentation. *arXiv preprint arXiv:1912.10917*, 2019.
- [7] Yilun Chen, Zhicheng Wang, Yuxiang Peng, Zhiqiang Zhang, Gang Yu, and Jian Sun. Cascaded pyramid network for multi-person pose estimation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 7103–7112, 2018.
- [8] Yukang Chen, Tong Yang, Xiangyu Zhang, Gaofeng Meng, Chunhong Pan, and Jian Sun. Detnas: Neural architecture search on object detection. *arXiv preprint arXiv:1903.10979*, 2019.
- [9] Shaojin Ding, Tianlong Chen, Xinyu Gong, Weiwei Zha, and Zhangyang Wang. Autospeech: Neural architecture search for speaker recognition. *arXiv preprint arXiv:2005.03215*, 2020.
- [10] Xuanyi Dong and Yi Yang. Searching for a robust neural architecture in four gpu hours. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1761–1770, 2019.
- [11] Hao-Shu Fang, Shuqin Xie, Yu-Wing Tai, and Cewu Lu. Rmpe: Regional multi-person pose estimation. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 2334–2343, 2017.
- [12] Yonggan Fu, Wuyang Chen, Haotao Wang, Haoran Li, Yingyan Lin, and Zhangyang Wang. Autogan-distiller: Searching to compress generative adversarial networks. *arXiv preprint arXiv:2006.08198*, 2020.
- [13] Golnaz Ghiasi, Tsung-Yi Lin, and Quoc V Le. Nas-fpn: Learning scalable feature pyramid architecture for object detection. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 7036–7045, 2019.
- [14] Xinyu Gong, Shiyu Chang, Yifan Jiang, and Zhangyang Wang. Autogan: Neural architecture search for generative adversarial networks. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 3224–3234, 2019.
- [15] Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross Girshick. Mask r-cnn. In *Proceedings of the IEEE international conference on computer vision*, pages 2961–2969, 2017.
- [16] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016.
- [17] Shaoli Huang, Mingming Gong, and Dacheng Tao. A coarse-fine network for keypoint localization. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 3028–3037, 2017.
- [18] Eldar Insafutdinov, Leonid Pishchulin, Bjoern Andres, Mykhaylo Andriluka, and Bernt Schiele. Deepercut: A deeper, stronger, and faster multi-person pose estimation model. In *European Conference on Computer Vision*, pages 34–50. Springer, 2016.
- [19] Shian-Ru Ke, LiangJia Zhu, Jenq-Neng Hwang, Hung-I Pai, Kung-Ming Lan, and Chih-Pin Liao. Real-time 3d human pose estimation from monocular view with applications to event detection and video gaming. In *2010 7th IEEE International Conference on Advanced Video and Signal Based Surveillance*, pages 489–496. IEEE, 2010.
- [20] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *European conference on computer vision*, pages 740–755. Springer, 2014.
- [21] Chenxi Liu, Liang-Chieh Chen, Florian Schroff, Hartwig Adam, Wei Hua, Alan L Yuille, and Li Fei-Fei. Auto-deeplab: Hierarchical neural architecture search for semantic image segmentation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 82–92, 2019.
- [22] Hanxiao Liu, Karen Simonyan, and Yiming Yang. Darts: Differentiable architecture search. *arXiv preprint arXiv:1806.09055*, 2018.
- [23] Renqian Luo, Fei Tian, Tao Qin, Enhong Chen, and Tie-Yan Liu. Neural architecture optimization. In *Advances in neural information processing systems*, pages 7816–7827, 2018.
- [24] Gyeongsik Moon, Ju Yong Chang, and Kyoung Mu Lee. Posefix: Model-agnostic general human pose refinement network. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 7773–7781, 2019.
- [25] Renato Negrinho and Geoff Gordon. Deeparchitect: Automatically designing and training deep architectures. *arXiv preprint arXiv:1704.08792*, 2017.
- [26] Vladimir Nekrasov, Hao Chen, Chunhua Shen, and Ian Reid. Fast neural architecture search of compact semantic segmentation models via auxiliary cells. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 9126–9135, 2019.

- [27] Alejandro Newell, Kaiyu Yang, and Jia Deng. Stacked hour-glass networks for human pose estimation. In *European conference on computer vision*, pages 483–499. Springer, 2016.
- [28] George Papandreou, Tyler Zhu, Nori Kanazawa, Alexander Toshev, Jonathan Tompson, Chris Bregler, and Kevin Murphy. Towards accurate multi-person pose estimation in the wild. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 4903–4911, 2017.
- [29] Xi Peng, Zhiqiang Tang, Fei Yang, Rogerio S Feris, and Dimitris Metaxas. Jointly optimize data augmentation and network training: Adversarial data augmentation in human pose estimation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2226–2234, 2018.
- [30] Hieu Pham, Melody Y Guan, Barret Zoph, Quoc V Le, and Jeff Dean. Efficient neural architecture search via parameter sharing. *arXiv preprint arXiv:1802.03268*, 2018.
- [31] Ruijie Quan, Xuanyi Dong, Yu Wu, Linchao Zhu, and Yi Yang. Auto-reid: Searching for a part-aware convnet for person re-identification. *arXiv preprint arXiv:1903.09776*, 2019.
- [32] Yaadhav Raaj, Haroon Idrees, Gines Hidalgo, and Yaser Sheikh. Efficient online multi-person 2d pose tracking with recurrent spatio-temporal affinity fields. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 4620–4628, 2019.
- [33] Esteban Real, Alok Aggarwal, Yanping Huang, and Quoc V Le. Regularized evolution for image classifier architecture search. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, pages 4780–4789, 2019.
- [34] Mikel D Rodriguez, Javed Ahmed, and Mubarak Shah. Action mach a spatio-temporal maximum average correlation height filter for action recognition. In *CVPR*, volume 1, page 6, 2008.
- [35] Michael S Ryoo and Jake K Aggarwal. Spatio-temporal relationship match: video structure comparison for recognition of complex human activities. In *ICCV*, volume 1, page 2. Citeseer, 2009.
- [36] Christian Schuldt, Ivan Laptev, and Barbara Caputo. Recognizing human actions: a local svm approach. In *Proceedings of the 17th International Conference on Pattern Recognition, 2004. ICPR 2004.*, volume 3, pages 32–36. IEEE, 2004.
- [37] Dehua Song, Chang Xu, Xu Jia, Yiyi Chen, Chunjing Xu, and Yunhe Wang. Efficient residual dense block search for image super-resolution.
- [38] Ke Sun, Bin Xiao, Dong Liu, and Jingdong Wang. Deep high-resolution representation learning for human pose estimation. *arXiv preprint arXiv:1902.09212*, 2019.
- [39] Mingxing Tan, Bo Chen, Ruoming Pang, Vijay Vasudevan, Mark Sandler, Andrew Howard, and Quoc V Le. Mnasnet: Platform-aware neural architecture search for mobile. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2820–2828, 2019.
- [40] Yunhe Wang, Yixing Xu, and Dacheng Tao. Dc-nas: Divide-and-conquer neural architecture search. *arXiv preprint arXiv:2005.14456*, 2020.
- [41] Ronald J Williams. Simple statistical gradient-following algorithms for connectionist reinforcement learning. *Machine learning*, 8(3-4):229–256, 1992.
- [42] Bin Xiao, Haiping Wu, and Yichen Wei. Simple baselines for human pose estimation and tracking. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 466–481, 2018.
- [43] Sirui Xie, Hehui Zheng, Chunxiao Liu, and Liang Lin. Snas: stochastic neural architecture search. *arXiv preprint arXiv:1812.09926*, 2018.
- [44] Yuliang Xiu, Jiefeng Li, Haoyu Wang, Yinghong Fang, and Cewu Lu. Pose flow: Efficient online pose tracking. *arXiv preprint arXiv:1802.00977*, 2018.
- [45] Sen Yang, Wankou Yang, and Zhen Cui. Pose neural fabrics search. *arXiv preprint arXiv:1909.07068*, 2019.
- [46] Yi Yang and Deva Ramanan. Articulated human detection with flexible mixtures of parts. *IEEE transactions on pattern analysis and machine intelligence*, 35(12):2878–2890, 2012.
- [47] Zhaohui Yang, Yunhe Wang, Xinghao Chen, Boxin Shi, Chao Xu, Chunjing Xu, Qi Tian, and Chang Xu. Cars: Continuous evolution for efficient neural architecture search. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1829–1838, 2020.
- [48] Zhuotun Zhu, Chenxi Liu, Dong Yang, Alan Yuille, and Daguang Xu. V-nas: Neural architecture search for volumetric medical image segmentation. In *2019 International Conference on 3D Vision (3DV)*, pages 240–248. IEEE, 2019.
- [49] Barret Zoph and Quoc V Le. Neural architecture search with reinforcement learning. *arXiv preprint arXiv:1611.01578*, 2016.
- [50] Barret Zoph, Vijay Vasudevan, Jonathon Shlens, and Quoc V Le. Learning transferable architectures for scalable image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 8697–8710, 2018.