# Light3DPose: Real-time Multi-Person 3D Pose Estimation from Multiple Views

Alessio Elmi ⓡ Davide Mazzini ⓡ Pietro Tortella

Checkout Technologies s.r.l.

20100 Milan, Italy

Email: {alessio, davide, pietro}@checkoutfree.it

*Abstract*—We present an approach to perform 3D pose estimation of multiple people from a few calibrated camera views. Our architecture, leveraging the recently proposed unprojection layer, aggregates feature-maps from a 2D pose estimator backbone into a comprehensive representation of the 3D scene. Such intermediate representation is then elaborated by a fully-convolutional volumetric network and a decoding stage to extract 3D skeletons with sub-voxel accuracy. Our method achieves state of the art MPJPE on the CMU Panoptic dataset using a few *unseen* views and obtains competitive results even with a single input view. We also assess the transfer learning capabilities of the model by testing it against the publicly available Shelf dataset obtaining good performance metrics. The proposed method is inherently efficient: as a pure bottom-up approach, it is computationally independent of the number of people in the scene. Furthermore, even though the computational burden of the 2D part scales linearly with the number of input views, the overall architecture is able to exploit a very lightweight 2D backbone which is orders of magnitude faster than the volumetric counterpart, resulting in fast inference time. The system can run at 6 FPS, processing up to 10 camera views on a single 1080Ti GPU.

## I. INTRODUCTION

Multi-person 3D pose estimation is a complex problem, with many applications in different fields of computer vision, like people tracking or augmented reality. This problem is usually tackled with a two steps approach. At first, every view is processed independently in order to produce a set of 2D poses - or possibly, some intermediate feature representation. In this stage, all the achievements in 2D pose estimation field can be exploited (see [1] for a survey). Next, these poses have to be matched across views and eventually triangulated, in order to produce a final estimate of the 3D scene. Usually, occlusions between people - or even self-occlusions - are the main difficulties to deal with: crowded scenes and complex poses produce noisy 2D detections, which are hard to filter out or recover in the matching and triangulation phase.

Hence, the idea of creating a system which is able to handle occlusions in a global way, and that is not affected by the limitations brought by single-view inferences. Inspired by [2], [3] and [4], we developed a multi-person 3D reconstruction system, which takes a set of images capturing the scene from different views and outputs a set of 3D pose reconstructions in a global reference frame. Its main building block is a fully convolutional neural network, where low-level features of the input views are unprojected, fused and transformed

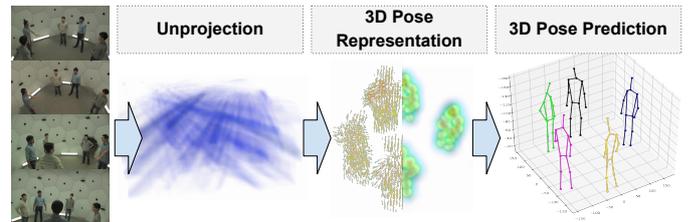ⓡ Equal contribution. Authors order determined by random function.



Fig. 1. This work proposes a fast and scalable approach for multi-person 3D pose estimation. Feature representations extracted from each views are aggregated and exploited to perform unified triangulation and pose estimation.

in order to produce a 3D representation of the probabilistic space. Following the general bottom-up approach of pose estimation, we extended the notion of *part affinity field* in three dimensions, making the pose reconstruction from density maps quick and agile. By doing this, we avoided all the limitations of the top-down strategies, where scalability is penalized as the number of people grows, and where inter-person occlusions and self-occlusions cannot be encoded - and recovered - by the network in a global way. On the contrary, thanks to the huge variety of pose configurations available in the CMU Panoptic dataset and with clever augmentation strategies about view-points, we could prove that our system is not affected by those limitations: our system can exploit *activations* and "shadows" in the feature space to estimate occlusions. Moreover, it does not depend on sophisticated algorithms of *detection-view assignment*, and it does not pay the computational burden of adding more views and subjects in the reconstruction process. Furthermore, we found that our system can produce good results even with just a single view suggesting that this approach can be further investigated also for monocular depth estimation tasks with multiple poses.

We conducted several experiments, which show the feasibility of our work and compare it to the other state-of-the-art approaches.

Our main contributions are the following:

- as far as we know this is the first complete bottom-up approach adopted in this context. In particular, it is capable of handling crowded scenes with good accuracy results and computational time.
- We show that even a very light backbone can produce good results. This implies that adding more views is almost computationally free.
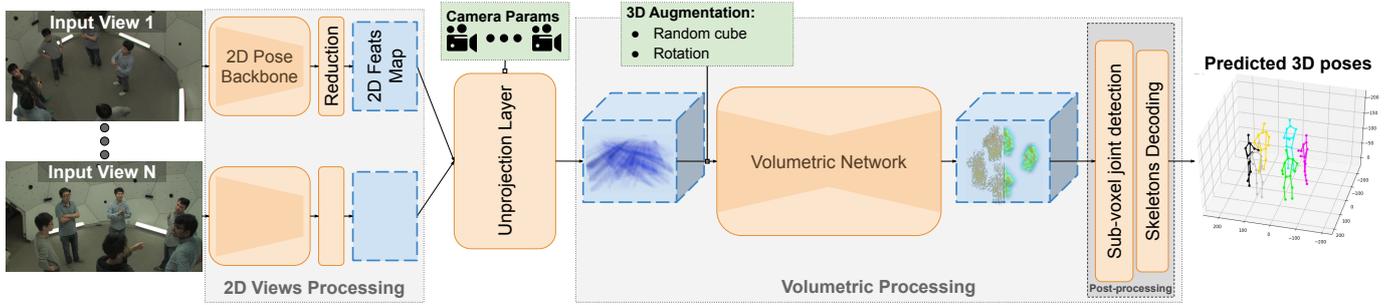
Fig. 2. An overall view of the complete processing pipeline. 2D pose backbone replicas process each view separately. Feature maps are then aggregated by the Unprojection layer into a 3D input representation of the scene. A volumetric network produces an output representation. A further decoding produces the final 3D pose estimations.

- We introduce 3D-data augmentation policies that greatly enhance the number of samples seen by the volumetric network.
- Our post-processing strategy leads to a sub-voxel localization, overcoming the issue of a quantized 3D space.

## II. RELATED WORKS

Multi-view, multi-person 3D pose estimation tries to fuse the achievements coming from 2D pose estimation, structure from motion and monocular depth estimation research fields. All of them are very well studied and pretty active topics nowadays.

Pose estimation from a single image is usually tackled following one of these two main strategies: *bottom-up* or *top-down* approaches. The former try to infer all key-points (i.e. *parts*) and/or limbs simultaneously and aggregate them eventually, using specific post-process logic. These methods can claim higher speed over their competitors since neural inference is done only once. At the same time, they usually have to deal with down-scaled feature maps, which limit the accuracy in terms of localization. In this group, we cannot omit [4]. Inspired by the work of [5], they introduced the notion of part-affinity fields. Their work has been extended by [6], [7] leading to a better part association, and by [8], where stronger descriptors led to a finer sub-pixel resolution. Other insights on the resolution issue were provided by [9], with heatmap encoding/decoding refinements. On the contrary, top-down approaches [10], [11], [12], possibly combined with multi-scale strategies [13], [14], rely on object detectors to identify humans in the scene, then a single-person neural inference is performed for each of them. These techniques generally outperform their bottom-up competitors on public challenges, while suffering in scalability with increasing number of subjects. Some hybrid approaches have emerged as well [15]. Finally, we want to mention some attempts [16], [17] to reduce the computational burden of pose estimation networks.

Three-dimensional pose estimation has emerged following two different tracks. The first one aims to recover the third dimension from a monocular view [18], [19], [20], [21], [22], [23], [24]. These methods usually start from 2D pose estimations, and lift them in a second stage in order to obtain their depth. In particular, they all deal with single-pose scenarios. We mention two attempts to extend this task to a multiple poses: Moon et al. [25] adopted a top-down strategy; Rogez et al. [26] introduced pose proposals (from anchor-poses) in the spirit of the Faster R-CNN approach. The second research track takes advantage of multiple views and claims to reconstruct 3D poses in a global reference frame. Sometimes this is the initial step of detection-to-track pipelines, like in [27], [28], where temporal evolution can be exploited in order to refine predictions. Multiple-view pose reconstruction may focus on single [29], [30], [31], [32], [2] or multiple poses [33], and they can exploit geometrical constraints [34], [35], in pair with visual features [36]. In particular, we highlight two works where multi-view projections have been combined with deep learning. [32] exploited the epipolar geometry in order to refine 2D pose estimation model, and consequently improve the final single pose 3D reconstruction. [2] showed that 2D features of each view can be fused and processed into a volumetric representation, which is analyzed to achieve a neat 3D reconstruction of the pose - again - for a single subject. However, to the best of our knowledge there is not any attempt to extend this approach to a multi-person scenario.

## III. METHOD

We call Light3DPose our system. In this section, we outline the architecture of Light3DPose, followed by a detailed explanation of all its components.

We are given a *detection space* $S$ with fixed boundaries and a set of fixed *setup cameras* $\{C_i\}_{i=1,\dots,N_c}$ whose intrinsic and extrinsic parameters are known. In particular, the projections $P_i : S \to F_i$ are known, where $F_i$ denotes the frame of the camera $C_i$. The cameras are synchronized, so for any time $t$ we have a set of images $I_1^t, \dots, I_{N_c}^t$, one for each camera. We will assume the time fixed throughout the paper, and omit the superscript $t$.

The input of Light3DPose is a set of pairs $\{I_i, C_{\nu_i}\}_{i=1,\dots,m}$ where $I_i$ is an image and $C_{\nu_i}$ is one of the setup cameras. The number of input pairs $m$ is variable and can range from 1 to $N_c$. In Section V we study both from the performance and computational sides the impact of the number of input views.

The output of Light3DPose is a set of 3D human poses $\{A_1, \dots, A_k\}$, with $k$ an arbitrary number. A 3D human pose

$A_i$ is a list $(a_i^l)_{l \in pose\_layout}$ of *joints*. Each joint is a pair composed of a point in the space $S$ and a label identifying the joint type, but when no confusion arises we identify the joint with the underlying point in $S$. The joint type ranges in a *pose layout* named CMU14, described in Section IV-C.

The internal pipeline of Light3DPose is composed of three main stages (see Figure 2):

- a *2D Views Processing* stage which returns a 2D feature map for each camera;
- an *Unprojection layer* [2] which aggregates the information coming from all the 2D views into a 3D features space representation;
- a *Volumetric Processing* that process the aggregated 3D representation and produces the output;

and each of these stages are composed of a different number of modules.

### A. 2D Views Processing

This processing stage takes as input one image $I$ and produces a 2D activation $\mathcal{R}(\mathcal{B}(I))$. When Light3DPose processes a set of pairs $\{(I_i, C_i)\}$, each $I_i$ is fed independently to the 2D Views processing stage. The different 2D View Processing stages share the same weight.

The stage is composed of two modules: a *2D Pose Backbone* followed by a *Reduction module*.

*1) 2D Backbone:* The input to the 2D Backbone module is an image $I$, and the output is a 2D feature map $\mathcal{B}(I)$. The 2D backbone is a MobileNet V1 [37] with some modifications from [16] on the latest layers. The stride of *conv4_2dw* has been removed and all succeeding convolutions have been set to dilation 2. This operation makes the network global stride to be 16 instead of 32 which is common for classification networks. We used weights pretrained on COCO dataset from [16].

*2) Reduction Module:* Input to the reduction module is the 2D feature map $\mathcal{B}(I)$, and the output is a 2D feature map $\mathcal{R}(\mathcal{B}(I))$. The purpose of this module is to project the feature space produced by the 2D Backbone to a lower-dimensional feature space. This module is crucial in order to encode the information of the backbone into a lighter feature map, to maintain the computations performed by the Volumetric Network feasible. Our Reduction Module is essentially a residual module composed of three depth-wise convolutions + ReLUs. We borrow this architecture from [16].

### B. Unprojection

This processing stage represents the contact point between the 2D feature maps and the 3D model of the scene, collecting the result of the 2D Processing stage into a 3D feature map representation. This is the only stage of Light3DPose that uses the calibration parameters of the cameras, and it has no trainable parameters.

Fix integer numbers $Q_x, Q_y, Q_z$, and a positive float value $Q_{size}$. Construct a cube $C \subseteq S$ composed of $Q_x \times Q_y \times Q_z$ voxels with edge of length $Q_{size}$. In C one has the integer coordinate system $(i_x, i_y, i_z)$ corresponding to the index of the voxels of C, and we denote by $\iota : C \to S$ the embedding.

The input to this stage is a set of pairs $\{(\mathcal{R}_i, C_{\nu_i})\}_{i=1,\ldots,m}$, where each $\mathcal{R}_i$ is the output of one of the 2D View Processing modules, and $C_{\nu_i}$ is one of the setup cameras.

The output of the unprojection stage is a 3D feature map $\mathcal{U}^\iota$ with shape $Q_x \times Q_y \times Q_z \times N_{feats}$, where $N_{feats}$ is the number of channels of the 2D feature map $\mathcal{R}$.

To compute the value of the $j$-th feature of the voxel $\mathcal{U}^\iota(i_x, i_y, i_z)$ we use the formula:

$$\mathcal{U}^\iota(i_x, i_y, i_z)^j = \frac{1}{m} \cdot \sum_{i=1}^m \mathcal{R}_i(P_{\nu_i}(\iota(i_x, i_y, i_z)))^j \quad (1)$$

where recall that $P_i$ denotes the projection associated with the camera $C_i$, and by $\mathcal{R}_i(u, v)^j$ we denote the $j$-th channel of the 2D feature map $\mathcal{R}_i$ at the point with frame coordinates $(u, v)$.

This layer is a generalization of the Unprojection introduced in [2] where a cube is built around each person. It can be efficiently implemented using vectorized operations and a differentiable sampling operator [38].

### C. Volumetric Processing

Input to this stage is the 3D feature map $\mathcal{U}$ output of the Unprojection layer.

The output of this stage is a set of 3D human poses $\{A_1, \ldots, A_k\}$.

The stage is composed of three modules:

- the *Volumetric Network*,
- the *Sub-voxel Joint Detection*
- the *Skeleton Decoder*.

The approach is similar to OpenPose [4]: the neural part of the network is trained to predict a Gaussian centered on each joint; the network should also predict a set of Part Affinity Fields (PAFs) that are used by the decoder to efficiently build the skeletons. Our method directly predicts 3D poses, thus the main differences between our volumetric processing part and OpenPose are in the use of a different neural architecture to handle 3D volumes data, and an adaptation to the 3D setting of the decoding of the output of the Volumetric network. Moreover, we introduce a Sub-voxel Peak Detector module to increase the accuracy of the joints predictions.

*1) Volumetric Network:* This is the trainable neural part of the volumetric processing. The purpose is to predict a set of 3D Gaussians centered on every joint and a set of 3D PAFs for the skeleton reconstruction.

The input to this module is the 3D activation $\mathcal{U}$ output of the unprojection layer with shape $Q_x \times Q_y \times Q_z \times N_{feats}$.

Output of this module is a 3D activation $\mathcal{V}$ with shape $Q_x \times Q_y \times Q_z \times N_{gt}$, where $N_{gt} = N_{joints} + 3 \cdot N_{PAF}$, where $N_{joints}$ is the number of joints of the pose layout, and $N_{PAF}$ is the number of PAF. This output can also be seen as a pair of collections $\mathcal{V} = ((H^l)_{l \in pose\_layout}, (\mathbf{V}^s)_{s \in PAFs})$ where each $H^l$ is a 3D feature map corresponding to a heatmap and each $\mathbf{V}^s$ is a collection of 3 (one per each of to the 3-dimensional

directions of the vector) 3D features map corresponding to a vectormap.

We adopt a V2V network from [39], but we set the minimum number of channels of the earliest and latest layers to $64$ in wherever layer the original network has 32 channels. We name this modified V2V network: *V2V64*. We also experimented with 32 and 96 channels architectures. Results are reported in Section V.

The output $\mathcal{V}$ of the module is then confronted with the ground-truth with an appropriate loss function, which is used to perform the training of Light3DPose. The dataset labels are lists of poses of persons in the 3D space. The procedure to create ground-truth heatmaps and vectormaps is a generalization to 3D space of the one in [4], so we omit the details. We opted to use a SmoothL1 loss function and to weight equally the loss coming from the heatmap and the vectormap. We experimented different loss functions and weights between heatmap and vectormap, the results are reported in V.

*2) Sub-voxel Joint Detector:* Several state-of-the-art works on single-person pose estimation are based on a variation of the Integral Regression Framework [21], [2], [14] which represents the unifying approach between heatmap and regression-based methods. The Integral Pose Regression framework assumes that the point to be localized follows a unimodal distribution. This is not the case of multiple poses scenarios, where more than one peak need to be estimated. We present an alternative formulation of such framework which, under the correct assumptions, can be used in a multi-person setup.

The sub-voxel joint detector module takes as input one heatmap $H$ ouput of the Volumetric network, and outputs a list of peaks $\mathcal{S}(H) = \{p_i\}$. The module is applied to each joint heatmap $\{H^l\}$, obtaining a set of peak for each joint type $\{\{p_i^l\}_{i=1,\dots,n_l}\}_l$.

In order to simplify the notation, we discuss the 1D case, but operators can be intuitively extended to 2D or 3D. Given a learned heatmap $H$, for each spatial location $x$ the values $H(x)$ represent the probability of such location of being a joint. We fix a neighbour function $N : C \to subsets(C)$ that associates to each point a neighbor of it (typically, an interval of a given radius centered at $x$). Define the non-local maxima suppression $\mathbf{P} : C \to \{0, 1\}$ via the formula :

$$\mathbf{P}(x) = \delta\left(\left(\max_{\bar{x} \in N(x)} \mathbf{H}(\bar{x})\right) = \mathbf{H}(x)\right) \quad (2)$$

where $\delta$ is a Dirac function. $\mathbf{P}(x) = 1$ if and only if $x$ is a maximum of $H_{|N(x)}$. Define the *pixel-peaks* as

$$\tilde{R} = \{x \in C \mid P(x) = 1\}$$

For each $x \in \hat{R}$, define the localized heatmap

$$L_x H = \frac{1}{\sum_{\bar{x} \in N(x)} H(\bar{x})} \cdot H_{|N(x)}$$

Finally, define the *sub-pixel peaks* as

$$\mathcal{S}(H) = \left\{ \sum_{\bar{x} \in N(x)} \bar{x} \cdot (L_x H)(\bar{x}) \mid x \in \tilde{R} \right\}$$

| | | MPJPE (cm) | Head | Torso | Up Arm | Lo Arm | Up Leg | Lo Leg | Avg |
|---|---|---|---|---|---|---|---|---|---|
| Cube | Rotation | | | | 3D Augmentations | | | | |
| | | 8.236 | 99.1 | 99.3 | 87.8 | 65.4 | 96.9 | 88.3 | 89.2 |
| ✓ | | 4.598 | 99.6 | **99.7** | 98.5 | 90.1 | **99.3** | 98.5 | 97.7 |
| | ✓ | 5.350 | 99.6 | **99.7** | 98.6 | 91.1 | 99.0 | 94.9 | 97.3 |
| ✓ | ✓ | **3.859** | **99.7** | **99.7** | **99.5** | 95.6 | 99.3 | **98.8** | **98.8** |
| | | | | | Number of Volumetric Features | | | | |
| 32 | | 4.760 | 99.6 | 99.7 | 97.1 | 78.9 | 99.5 | 98.6 | 95.9 |
| 64 | | **3.859** | **99.7** | **99.7** | **99.5** | 95.6 | **99.3** | **98.8** | 98.8 |
| 96 | | 3.975 | **99.7** | **99.7** | **99.5** | 96.2 | **99.3** | 98.7 | **98.9** |
| | | | | | Loss Type | | | | |
| L1 | | 4.106 | 99.6 | **99.7** | 99.2 | 96.2 | 99.0 | 98.0 | 98.7 |
| L2 | | 4.125 | 99.6 | **99.7** | **99.5** | 96.6 | **99.4** | **98.9** | **99.0** |
| SmoothL1 | | **3.859** | **99.7** | **99.7** | **99.5** | 95.6 | 99.3 | 98.8 | 98.8 |
| | | | | | Heatmap / Vectormap Loss Ratio | | | | |
| 1 | | **3.859** | **99.7** | **99.7** | **99.5** | 95.6 | 99.3 | **98.8** | 98.8 |
| 3 | | 4.074 | **99.7** | **99.7** | 99.1 | **96.6** | **99.5** | 98.6 | **98.9** |
| 10 | | 3.935 | **99.7** | **99.7** | 98.0 | 90.9 | **99.5** | **98.8** | 97.9 |
| | | | | | Sub-voxel refinement | | | | |
| | | 4.899 | **99.7** | **99.7** | 99.4 | 94.9 | **99.3** | **98.8** | 98.6 |
| | ✓ | **3.859** | **99.7** | **99.7** | **99.5** | 95.6 | **99.3** | **98.8** | 98.8 |

The assumption we rely on is that for every $x$, in the neighbour $N(x)$ there should be at most one local maximum. In general, this assumption holds if the radius is small enough w.r.t. the quantization constant $Q_{size}$. In practice, we obtain good results by choosing $N(x)$ to be a 1 or 2 voxels radius interval centered at $x$, see Section V.

*3) Skeletons decoder:* This module takes as input the peaks $\{\mathcal{S}(H^l)\}_{l \in pose\_layout}$ of the sub-pixel joint detection and the vectormaps $\{\mathbf{V}^s\}_{s \in PAFs}$ output of the Volumetric Network and outputs a list of 3D poses. Our algorithm is a direct extension of the one proposed by OpenPose [4], with the only difference that line integrals are computed over three-dimensional vector fields.

## IV. EXPERIMENTAL SETUP

### A. Datasets

*1) CMU Panoptic dataset [3]:* it consists of 31 Full-HD and 480 VGA video streams from synchronized cameras at 29.97 FPS; various scenes (65 sequences with multiple people, social interactions, and a wide range of actions) for a total duration of 5.5 hours. The dataset includes robustly labeled 3D poses, computed using all the camera views. This dataset is perhaps the most complete, open and free to use dataset available for the task of 3D pose estimation. However,
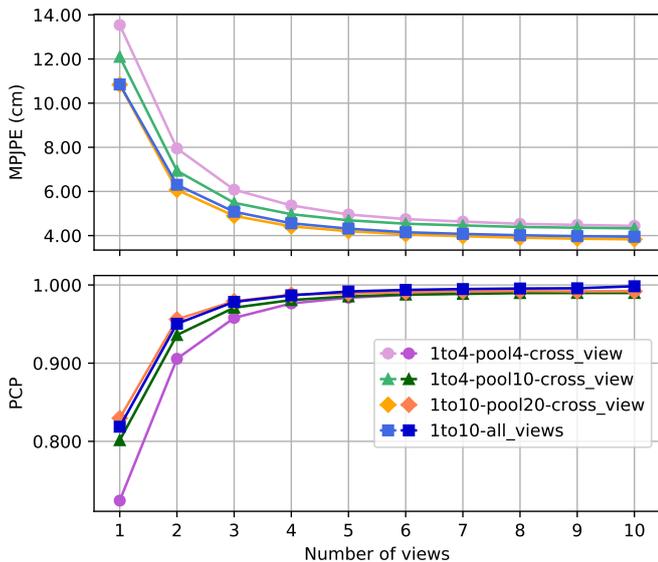
Fig. 3. Accuracy vs number of views. *1to4-pool4-cross_view*: training with 1 to 4 simultaneous views from a pool of 4; *1to4-pool10-cross_view*: 1 to 4 from a pool of 10; *1to10-pool20-cross_view*: 1 to 10 from a pool of 20; *1to10-all_views*: 1 to 10 from all available views. Configurations 1-2-3 are *cross-view*.

considering that they released annotations quite recently, most works in literature use it only for qualitative evaluations [36] or for single-person pose detection [2] discarding multi-person scenes. To the best of our knowledge only [33] makes use of CMU Panoptic dataset to train and evaluate multi-person 3D pose estimation. We adopt the same subset of scenes and the same train/val/test split of CMU Panoptic used in [33]: 20 scenes (343k images) of which 10, 4 and 6 scenes for training, validation and test respectively. Only HD cameras are used with data frame rates downsampled to 2 FPS. Since one of our concerns is to assess the cross-view generalization of our model, we split the dataset by scene and by view. Val and test splits use cameras 2, 13, 16, 18 while the train split uses all (or a subset of) the remaining 27 cameras. This is the same camera split used by [2]. We name this dataset: *PanopticD2D*.

*2) Shelf [34]:* we adopt this dataset to evaluate the ability of our model to transfer to a completely unseen setup. It consists of a single scene of four people disassembling a shelf at a close range. Video streams are from five calibrated cameras. The dataset includes 3D annotated groundtruth skeletons.

### B. Evaluation metrics

We employed two commonly used metrics that capture different types of errors in models prediction:

- *MPJPE: Mean Per Joint Precision Error*. Given a pair of skeletons, MPJPE is defined as the average of the square distance of the predicted joints from the corresponding ground-truth joints.
- *PCP: Percentage of Correct estimated Parts*. We implemented this metric according to [36]. A body part is correct if the average distance of the two joints is less than

a threshold from the corresponding groundtruth joints locations. The threshold is computed as the 50% of the length of the groundtruth body part.

Before computing these metrics we associate for each scene the predicted skeletons to the groundtruth skeletons using linear assignment.

### C. Implementation details

*1) Pose Layout:* We used a simplified pose layout of 14 keypoints. Apart from the canonical 12 parts of arms and legs, we only added *neck* and *nose*. Sometimes, a layout conversion was needed across different datasets and labeling standards. Moreover, we defined 13 PAFs; starting from the *neck*, a tree-structure along arms, legs and nose has been defined. In our setup, increasing excessively the number of joints or PAFs would not make sense due to the limitations of our quantized space.

*2) Skeleton Decoding:* Parameters have been found by performing a grid search on Panoptic D2D validation set. Eventually, we opted for an interpolation over a region of size $5 \times 5 \times 5$ voxels. Then, all local maxima with a score lower than 0.3 are discarded; every PAF where the linear integral is on average lower than 0.2 is also removed. Finally, only candidate poses with more than 7 keypoints are retained.

*3) 3D space quantization:* we set the size of the quantization voxel to 7.5 cm. This allows us to maintain a quantization of $64 \times 64 \times 32$ voxels on Panoptic dataset to efficiently cover the whole scene of approximately $5 \times 5 \times 2.5$ meters, the last dimension being the vertical axis.

*4) Training recipe:* Models have been trained with Adam optimizer. We set the initial learning rate to 0.002 and used a step decay policy of 0.3 every 50 epochs. All models have been trained for 200 epochs with a batch size of 8. We implemented the architecture in PyTorch.

## V. ABLATION ANALYSIS

### A. 3D Augmentation

We applied 3D data augmentation techniques to the 3D feature space between the Unprojection layer and the Volumetric network. In particular, we implemented the followings:

*1) Random cube embedding:* During the training, we consider $C \subseteq S$ to be strictly smaller, and to be randomly embedded. This corresponds to take a random crop of the 3D crop of the scene to be considered for the parameters update.

From the volumetric network point of view, this reflects into a data augmentation strategy, since moving the cube inside $S$ corresponds to a change of the observed scene and a change in the extrinsic parameters of the cameras.

We set C to have $32 \times 32 \times 32$ voxels, and we change the embedding at the start of each epoch.

*2) Random rotation:* we implement rotations along the vertical axis of $90°, 180°, 270°$, to allow a fast implementation. One should take care of the fact that rotation of the 3D space is not reflected into images transformation, so when a rotation is applied we cut the back-propagation graph just before the unprojection layer. In our specific architecture, this sparse

| Model | MPJPE (cm) | | | PCP |
|---|---|---|---|---|
| | single | multi | avg | avg |
| ACTOR [33] (2 views)* | 17.21 | 50.24 | 33.72 | - |
| ACTOR (4 views)* | 8.19 | 20.10 | 14.14 | - |
| ACTOR (10 views)* | 6.13 | 12.21 | 9.17 | - |
| Oracle [33] (using GT to select cameras)* | **4.24** | 9.19 | 6.71 | - |
| Ours (1 unseen view) | 10.34 | 9.32 | 9.43 | 80.8 |
| Ours (2 to 4 unseen views depending on scene) | 5.30 | **4.09** | **4.22** | 98.2 |
| Ours (10 views, from training view pool) | **3.50** | **3.56** | **3.55** | **98.6** |

*ACTOR: number in brackets refers to maximum number of views to choose from. Oracle means: best views to triangulate are selected using groundtruth.

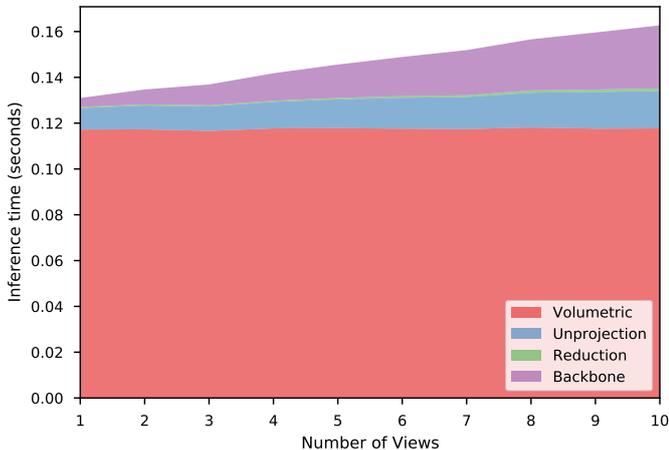| Model | Actor 1 | Actor 2 | Actor 3 | Avg | Speed(s) |
|---|---|---|---|---|---|
| Belagiannis et al. [34] | 66.1 | 65.0 | 83.2 | 71.4 | - |
| Belagiannis et al. [40] | 75.0 | 67.0 | 86.0 | 76.0 | - |
| Belagiannis et al. [41] | 75.3 | 69.7 | 87.6 | 77.5 | - |
| Ershadi et al. [42] | 93.3 | 75.9 | 94.8 | 88.0 | - |
| Dong et al. [36] | **98.8** | **94.1** | **97.8** | **96.9** | .465 |
| Ours | 94.3 | 78.4 | 96.8 | 89.8 | **.146** |



Fig. 4. Adding more views increases computational time by a linear factor. However, only few modules are affected by this growth. The main CNN block (in red) has a $O(1)$ complexity, both in the number of views and people. Inference time is measured on a single NVIDIA GeForce GTX 1080Ti.

back-prop signal does not drastically affect the training since the only trainable part before the volumetric network is the Reduction layer which has a limited number of parameters.

### B. Architecture

In Table I we reported the results of different experiments to evaluate the contribution of our architectural choices.

*1) Number of volumetric features:* it refers to the channels of the volumetric input: it involves the 2D feature maps, the input/output of the unprojection and the volumetric network. For 32 features we used the original V2V network whereas for 64 and 96 we modified it as described in Section III-C1. Models with 64 and 96 channels achieve similar MPJPE and PCP

values but 64 is an obvious choice for being computationally lighter.

*2) Loss:* we run experiments with different loss types and weighted differently the heatmap and vectormap losses. By evaluating separately PAFs and Peaks quality we noticed that good peaks have a stronger impact in the final metrics than good PAFs, thus we weighted more the Peak part of the loss. Results seem to suggest that the task of predicting good peaks should be tackled with a more elaborate approach than simply differentiate loss weights.

*3) Sub-voxel refinement:* by activating it we achieve a lower MPJPE. It has almost no effect on PCP since it improves the sub-voxel localization but does not reduce false positives.

### C. Study on the number of input views

These experiments have a two-fold goal. On one side, we wanted to understand better the impact on the accuracy of a short/large number of views in the training pool; on the other hand, we wanted to check how well our augmentation strategies could compensate/emulate unseen angles. In Figure 3 we reported four experiments where we varied the number of views and the number of simultaneous angles used on each training inference. In particular, they show that even a few cameras can produce mildly good results; also, after a certain number, adding more views gives unnoticeable improvements.

## VI. COMPARISON WITH STATE-OF-THE-ART

In Table II we report a comparison between our method and the results in [33] (ACTOR and ORACLE). We remark that the task that authors of [33] are trying to solve is different from ours. They train an agent to find what are the best views to use to triangulate *that* particular scene. We consider it to be a good baseline even if the core task of [33] is not the triangulation algorithm itself. We select 4 fixed validation views and we never train on those. Since some recordings have fewer views available, it turned out that only 36.2% of the test set has 4 views, 31.3% has 3 and 32.5% has just two angles available. The evaluation metric is the MPJPE expressed in $cm$. The MPJPE of our method is more than 3 times lower compared to ACTOR with 4 views and on average *lower* than the Oracle.

We also run our model on the Shelf dataset in order to test it in a completely new environment with unseen views, camera parameters, sensors and every other detail that can bias the evaluation. Results are reported in Table III. Our method obtains good results even if not on-par with the work by Dong et al. [36]. However, their approach is much slower being based on top-down 2D backbones. We detail a speed comparison between the two methods in Section VI-A.

### A. Inference speed

Being a pure bottom-up approach, our method can scale well when increasing the number of views and subjects. Even though our complexity is $O(n)$ in the number of views and $O(n^2)$ in the number of people, adding more cameras affects only the *Backbone*, *Reduction* and *Unprojection* modules, which are a small fraction of the cumulative computation
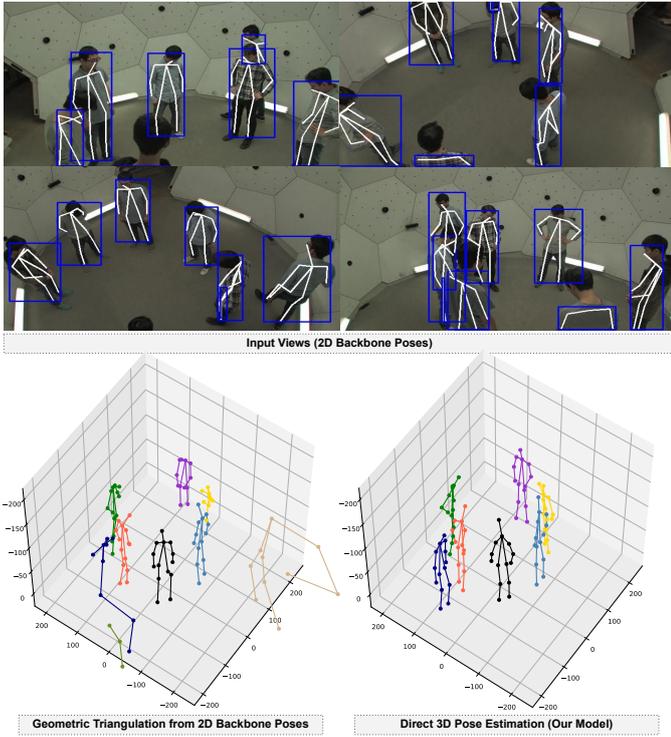
Fig. 5. Left: geometric triangulation using 2D poses from Lightweight OpenPose [16] (same 2D backbone weights as ours) and iterative greedy matching [35]. Right: direct 3D pose estimation with our model.
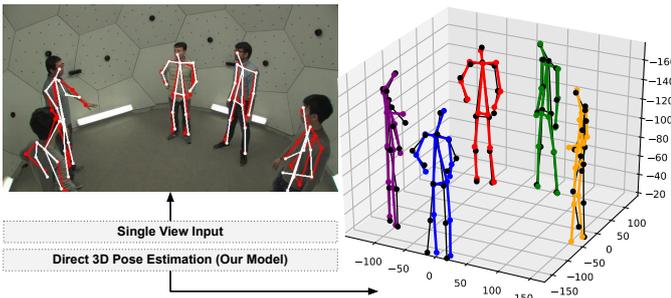


Fig. 6. Direct 3D pose predictions by our model from a single camera view. On the frame we projected in red the groundtruth, in white our predictions. In the 3D plot: predictions in color, groundtruth in dashed-black. The network "hallucinates" straight legs of non visible body parts relying on a strong learned prior.

burden (e.g. for 10 views they take all together only 45 ms, see Figure 4). On the other hand, post-processing the CNN output costs even less; Cao et al. [4], implemented an optimized version which takes 0.58 ms for a 9 people image. For reference we can compare our method with the one presented in [36], see Table III. Their approach starts with a person detector [43], which takes around 10 ms per view. Then, each detection is forwarded to two branches, of which the 2D pose estimation [11] is most expensive (we measured 67 ms). From here, the final 3D pose inference takes around 80 ms. We can estimate that a 5 views scenario with 5 people will take $(10 * 5) + (67 * 5) + 80 \approx 465ms$, which is about

3.2 times our implementation.

*B. Qualitative results*

In figure 5 we show a comparison between our model which performs direct 3D estimations and the result of the geometric triangulation using the 2D skeletons predicted by Lightweight OpenPose[16] and the iterative greedy matching by [35]. Notice that our 2D backbone has exactly the same weights as the backbone of [16] since we do not train nor finetune such part of the network. This highlights the power of estimating directly 3D poses: our volumetric architecture can learn strong pose priors and implicitly discards false detections. By exploiting the 3D representation of the space, it is less prone to occlusion-related errors and it can better deal with crowded scenes. This behavior is even more evident in Figure 6 where our method correctly predicts all 3D poses from a monocular view. In particular, notice that even the legs of the blue skeleton are predicted even if they are not visible from that particular view. (View and scene from the validation set). We suppose that the model hallucinates straight up legs since most of the people in Panoptic D2D training set are standing.

## VII. CONCLUSION

We present a method for multi-person human pose estimation from calibrated views. Our neural architecture is able to predict 3D pose representations *directly* from raw camera views. To the best of our knowledge, this is the first attempt to tackle such a task in a completely bottom-up fashion. The proposed method exhibits good computational scalability properties: in particular, it is essentially independent of the number of people in the scene. Moreover, it scales linearly with the number of input views.

Conducted experiments show state-of-the-art performance on the Panoptic D2D dataset as well as a good generalization on the unseen Shelf dataset. We hope that our work can open new research lines and new scenarios. The method visibly benefits from a wide variety of configurations of people, cameras, and environments during training. Simple 3D data augmentation techniques have been explored and proven effective in enhancing the performance; however, larger datasets, both real and synthetic, could significantly increase the model capabilities.

## REFERENCES

[1] Q. Dang, J. Yin, B. Wang, and W. Zheng, "Deep learning based 2d human pose estimation: A survey," *Tsinghua Science and Technology*, vol. 24, no. 6, pp. 663–676, 2019.

[2] K. Iskakov, E. Burkov, V. Lempitsky, and Y. Malkov, "Learnable triangulation of human pose," in *Proceedings of the IEEE International Conference on Computer Vision*, 2019, pp. 7718–7727.

[3] H. Joo, H. Liu, L. Tan, L. Gui, B. Nabbe, I. Matthews, T. Kanade, S. Nobuhara, and Y. Sheikh, "Panoptic studio: A massively multiview system for social motion capture," in *Proceedings of the IEEE International Conference on Computer Vision*, 2015, pp. 3334–3342.

[4] Z. Cao, T. Simon, S.-E. Wei, and Y. Sheikh, "Realtime multi-person 2d pose estimation using part affinity fields," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 7291–7299.

[5] S.-E. Wei, V. Ramakrishna, T. Kanade, and Y. Sheikh, "Convolutional pose machines," in *Proceedings of the IEEE conference on Computer Vision and Pattern Recognition*, 2016, pp. 4724–4732.

[6] A. Newell, Z. Huang, and J. Deng, "Associative embedding: End-to-end learning for joint detection and grouping," in *Advances in neural information processing systems*, 2017, pp. 2277–2287.

[7] G. Papandreou, T. Zhu, L.-C. Chen, S. Gidaris, J. Tompson, and K. Murphy, "Personlab: Person pose estimation and instance segmentation with a bottom-up, part-based, geometric embedding model," in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018, pp. 269–286.

[8] S. Kreiss, L. Bertoni, and A. Alahi, "Pifpaf: Composite fields for human pose estimation," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 11 977–11 986.

[9] F. Zhang, X. Zhu, H. Dai, M. Ye, and C. Zhu, "Distribution-aware coordinate representation for human pose estimation," 2019.

[10] A. Newell, K. Yang, and J. Deng, "Stacked hourglass networks for human pose estimation," in *European conference on computer vision*. Springer, 2016, pp. 483–499.

[11] Y. Chen, Z. Wang, Y. Peng, Z. Zhang, G. Yu, and J. Sun, "Cascaded pyramid network for multi-person pose estimation," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 7103–7112.

[12] B. Xiao, H. Wu, and Y. Wei, "Simple baselines for human pose estimation and tracking," in *Proceedings of the European conference on computer vision (ECCV)*, 2018, pp. 466–481.

[13] L. Ke, M.-C. Chang, H. Qi, and S. Lyu, "Multi-scale structure-aware network for human pose estimation," in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018, pp. 713–728.

[14] K. Sun, B. Xiao, D. Liu, and J. Wang, "Deep high-resolution representation learning for human pose estimation," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 5693–5703.

[15] M. Kocabas, S. Karagoz, and E. Akbas, "Multiposenet: Fast multi-person pose estimation using pose residual network," in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018, pp. 417–433.

[16] D. Osokin, "Real-time 2d multi-person pose estimation on cpu: Lightweight openpose," in *arXiv preprint arXiv:1811.12004*, 2018.

[17] F. Zhang, X. Zhu, and M. Ye, "Fast human pose estimation," in *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2019.

[18] D. Tome, C. Russell, and L. Agapito, "Lifting from the deep: Convolutional 3d pose estimation from a single image," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 2500–2509.

[19] D. Pavllo, C. Feichtenhofer, D. Grangier, and M. Auli, "3d human pose estimation in video with temporal convolutions and semi-supervised training," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 7753–7762.

[20] W. Yang, W. Ouyang, X. Wang, J. Ren, H. Li, and X. Wang, "3d human pose estimation in the wild by adversarial learning," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 5255–5264.

[21] X. Sun, B. Xiao, F. Wei, S. Liang, and Y. Wei, "Integral human pose regression," in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018, pp. 529–545.

[22] A. Nibali, Z. He, S. Morgan, and L. Prendergast, "3d human pose estimation with 2d marginal heatmaps," in *2019 IEEE Winter Conference on Applications of Computer Vision (WACV)*. IEEE, 2019, pp. 1477–1485.

[23] A. Kanazawa, J. Y. Zhang, P. Felsen, and J. Malik, "Learning 3d human dynamics from video," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 5614–5623.

[24] I. Habibie, W. Xu, D. Mehta, G. Pons-Moll, and C. Theobalt, "In the wild human pose estimation using explicit 2d features and intermediate 3d representations," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 10 905–10 914.

[25] G. Moon, J. Y. Chang, and K. M. Lee, "Camera distance-aware top-down approach for 3d multi-person pose estimation from a single rgb image," in *Proceedings of the IEEE International Conference on Computer Vision*, 2019, pp. 10 133–10 142.

[26] G. Rogez, P. Weinzaepfel, and C. Schmid, "Lcr-net++: Multi-person 2d and 3d pose detection in natural images," *IEEE transactions on pattern analysis and machine intelligence*, 2019.

[27] L. Bridgeman, M. Volino, J.-Y. Guillemaut, and A. Hilton, "Multi-person 3d pose estimation and tracking in sports," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*, 2019, pp. 0–0.

[28] T. Ohashi, Y. Ikegami, and Y. Nakamura, "Synergetic reconstruction from 2d pose and 3d motion for wide-space multi-person video motion capture in the wild," *arXiv preprint arXiv:2001.05613*, 2020.

[29] J. Martinez, R. Hossain, J. Romero, and J. J. Little, "A simple yet effective baseline for 3d human pose estimation," in *Proceedings of the IEEE International Conference on Computer Vision*, 2017, pp. 2640–2649.

[30] G. Pavlakos, X. Zhou, K. G. Derpanis, and K. Daniilidis, "Harvesting multiple views for marker-less 3d human pose annotations," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 6988–6997.

[31] D. Tome, M. Toso, L. Agapito, and C. Russell, "Rethinking pose in 3d: Multi-stage refinement and recovery for markerless motion capture," in *2018 international conference on 3D vision (3DV)*. IEEE, 2018, pp. 474–483.

[32] H. Qiu, C. Wang, J. Wang, N. Wang, and W. Zeng, "Cross view fusion for 3d human pose estimation," in *Proceedings of the IEEE International Conference on Computer Vision*, 2019, pp. 4342–4351.

[33] A. Pirinen, E. Gärtner, and C. Sminchisescu, "Domes to drones: Self-supervised active triangulation for 3d human pose reconstruction," in *Advances in Neural Information Processing Systems*, 2019, pp. 3907–3917.

[34] V. Belagiannis, S. Amin, M. Andriluka, B. Schiele, N. Navab, and S. Ilic, "3d pictorial structures for multiple human pose estimation," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2014, pp. 1669–1676.

[35] J. Tanke and J. Gall, "Iterative greedy matching for 3d human pose tracking from multiple views," in *German Conference on Pattern Recognition*. Springer, 2019, pp. 537–550.

[36] J. Dong, W. Jiang, Q. Huang, H. Bao, and X. Zhou, "Fast and robust multi-person 3d pose estimation from multiple views," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 7792–7801.

[37] A. G. Howard, M. Zhu, B. Chen, D. Kalenichenko, W. Wang, T. Weyand, M. Andreetto, and H. Adam, "Mobilenets: Efficient convolutional neural networks for mobile vision applications," *arXiv preprint arXiv:1704.04861*, 2017.

[38] M. Jaderberg, K. Simonyan, A. Zisserman *et al.*, "Spatial transformer networks," in *Advances in neural information processing systems*, 2015, pp. 2017–2025.

[39] G. Moon, J. Yong Chang, and K. Mu Lee, "V2v-posenet: Voxel-to-voxel prediction network for accurate 3d hand and human pose estimation from a single depth map," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 5079–5088.

[40] V. Belagiannis, X. Wang, B. Schiele, P. Fua, S. Ilic, and N. Navab, "Multiple human pose estimation with temporally consistent 3d pictorial structures," in *European Conference on Computer Vision*. Springer, 2014, pp. 742–754.

[41] V. Belagiannis, S. Amin, M. Andriluka, B. Schiele, N. Navab, and S. Ilic, "3d pictorial structures revisited: Multiple human pose estimation," *IEEE transactions on pattern analysis and machine intelligence*, vol. 38, no. 10, pp. 1929–1942, 2015.

[42] S. Ershadi-Nasab, E. Noury, S. Kasaei, and E. Sanaei, "Multiple human 3d pose estimation from multiview images," *Multimedia Tools and Applications*, vol. 77, no. 12, pp. 15 573–15 601, 2018.

[43] Z. Li, C. Peng, G. Yu, X. Zhang, Y. Deng, and J. Sun, "Light-head r-cnn: In defense of two-stage object detector," *arXiv preprint arXiv:1711.07264*, 2017.