

1 We thank the reviewers for their valuable feedback and note that all reviewers agree about the novelty of our work.

2 **AR1: Results for more than 10 cameras:** Triangulation converges when adding more cameras, thus narrowing ACTOR’s  
3 improvement over the baselines. Using many cameras however requires significant computation, hence we are interested  
4 in finding small yet informative subsets of cameras that yield good results. We chose to run the model for 1-10 cameras  
5 as evaluating the Oracle is slow for large camera sets, cf. Fig. 3 right. Results for 15 cameras on multi-people data:  
6 ACTOR = 69.99, Random = 74.71, Max-Azim = 73.79, Oracle = 57.24 (mm/joint hip-translated). Compared to  
7 running for 10 cameras, ACTOR keeps its gain over Max-Azim and significantly closes in on the Oracle, see Table 1.

8 **Minor issues:** (1) Fig. 5 is not missing; see page 8 and line 240. (2) We will clarify line ~ 223-224 in the paper.

9 **AR2: Runtimes and baseline comparisons:** Fig. 3-4 right show runtimes; OpenPose adds 0.134 s/image and our policy  
10 adds 0.005 s/image (25x faster), see line 237-240. ACTOR is almost as fast as the baselines while *significantly more*  
11 *accurate* when considering fewer views, e.g. 94 mm/joint (34 %) better than Max-Azim at 3 views (Table 1), yet only  
12 6% slower (Fig. 3 right). When using more cameras triangulation converges, decreasing the need for intelligent camera  
13 selection. The Oracle cheats by using 3d ground-truth when selecting views. It is greedy, as exhaustively evaluating the  
14 exponential set of views is infeasible. It is still very slow since it evaluates all cameras before selecting the next view.

15 **Dataset description:** We use the HD cameras, of which there are about 30 per scene. The HD cameras are used since  
16 they provide better image quality than VGA and are sufficiently dense, yet spread apart far enough to make each  
17 viewpoint unique. We will include this in an extended dataset description for the camera-ready.

18 **More realistic setup:** To mimic a single camera setup, we experimented on multi-people data assuming time elapses  
19 0.5 s between *every* selected view. Note that some joints which move drastically between views cannot be triangulated.  
20 Therefore we extend ACTOR with a monocular 3d pose estimator (DMHS, Popa et al., CVPR, 2017) as a fallback.  
21 *Without refining the policy* our results (in auto-mode) are: ACTOR = 123.3, Random = 138.7, Max-Azim = 150.1,  
22 Oracle = 106.1 (mm/joint hip-translated). ACTOR outperforms the baselines also in this setting (cf. Table 1). While  
23 these are only proof-of-concept results they clearly indicate the potential of our model for single camera setups.

24 **State-of-the-art:** Our results are significantly better than state-of-the-art 3d pose estimators such as MubyNet (Zanfir et  
25 al., NeurIPS, 2018). Their errors are about 150 mm/joint on multi-people data, while ACTOR obtains 96 mm/joint in  
26 auto-mode, see Table 1. There exists no prior work on active 3d human pose estimation, so results are hard to compare  
27 as they are framed in completely different settings (monocular 3d estimation vs active 2d-to-3d triangulation).

28 **AR3: No annotations needed for training:** We emphasize that our reward is *not* based on ground-truth visibility – it  
29 relies solely on the automatic 2d pose (body joint) detection. Thus our approach requires no annotations in training.

30 **Next-best-view (NBV) and object pose estimation:** NBV is a very general concept, and our formulation can be seen as  
31 one way of modeling and implementing it in the framework of deep reinforcement learning. We further address the  
32 novel problem of intelligent view selection for 3d human pose reconstruction (or generally articulated structures), which  
33 is more challenging than rigid 3d reconstruction and pose estimation or object detection<sup>1</sup>: i) Dimensionality is much  
34 higher than object detection and the problem is more complex than rigid 3d reconstruction - humans are articulated  
35 and deformable, scenes contain multiple people; ii) Realistic training data for supervised 3d human pose estimation is  
36 scarce, and our proposed methodology requires *no annotations* and instead uses self-supervision to learn the policy;  
37 iii) Reconstruction difficulty is impacted by occlusions and the pose of targets, hence our model adaptively selects a  
38 *variable number of views* depending on scene complexity; iv) NBV for e.g. object pose estimation is formulated such  
39 that, at each step, one selects views by assessing a utility function over a set of candidates. Most design choices are  
40 handcrafted, decisions are local. We frame the task as a deep RL problem where the policy is trained to maximize a  
41 *global objective*, searching over entire sequences of viewpoints, and by triangulating as many joints as possible.

42 **AR2, AR3: Practical applications:** We use our setup as a proxy for a moving observer and to develop viewpoint  
43 selection strategies for e.g. intelligent holoportation (AR1) while providing a framework for *reproducible* experiments.  
44 Practical developments of our methodology would include e.g. real-time intelligent processing of multi-camera  
45 (Panoptic/Lightstage) video feeds or control policies for a drone observer. In the latter case the model would further  
46 benefit from being extended to account for physical constraints, e.g. a single camera and limited speed. Our paper is a  
47 key first step since it presents fundamental methodology required for future applied research.

48 **AR1, AR3: Why not exhaustive triangulation:** Using all views requires significant computation, even with one of the  
49 fastest pose estimators such as OpenPose (0.134 s/image vs only 0.005 s/image for our policy, see line 238-239). Thus  
50 exhaustive computation does not suit real-time scenarios (processing 30 cameras takes over 4 s / frame). Also, assuming  
51 a physical observer (e.g. a drone), the need for a view selection strategy is crucial since covering all views is infeasible.

---

<sup>1</sup>Obviously there is no one size fits all – details on the output representation, level of detail, single versus multiple structures, occlusions, observation setup, or desired reward would affect formulation and modeling choices for each problem.