

SAM-pose2seg: Pose-Guided Human Instance Segmentation in Crowds

Constantin Kolomiiets

Miroslav Purkrabek

Jiri Matas

Visual Recognition Group
Department of Cybernetics
Faculty of Electrical Engineering
Czech Technical University in Prague
kolomcon@fel.cvut.cz

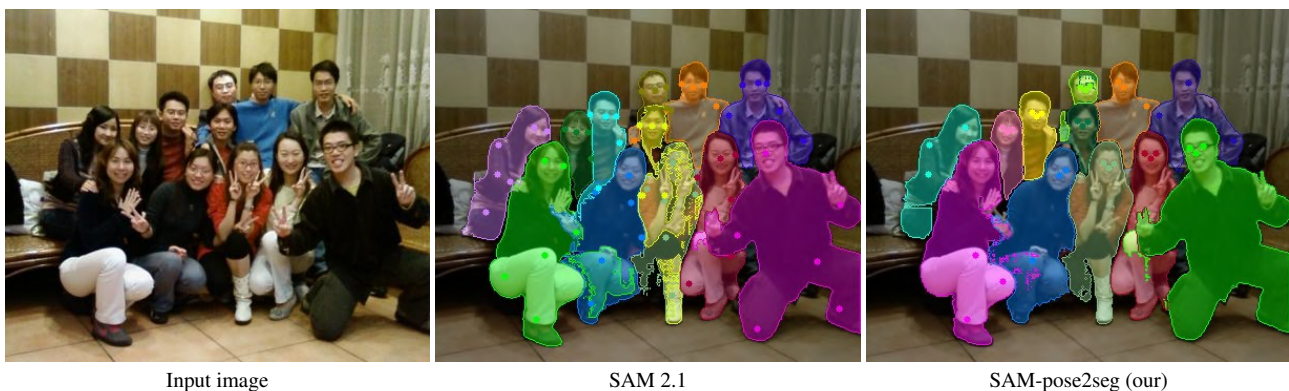


Figure 1. **SAM-pose2seg is superior to SAM 2** [21] for human instance segmentation, especially in crowded scenes. Both examples are generated from the same set of predicted keypoints from [19]. Notice the noise and incorrect masks in the middle of the crowd for SAM 2.1. Different prompts are used for SAM 2 and SAM-pose2seg since each model works best with different prompts.

Abstract

Segment Anything (SAM) provides an unprecedented foundation for human segmentation, but may struggle under occlusion, where keypoints may be partially or fully invisible. We adapt SAM 2.1 for pose-guided segmentation with minimal encoder modifications, retaining its strong generalization. Using a fine-tuning strategy called PoseMaskRefine, we incorporate pose keypoints with high visibility into the iterative correction process originally employed by SAM, yielding improved robustness and accuracy across multiple datasets. During inference, we simplify prompting by selecting only the three keypoints with the highest visibility. This strategy reduces sensitivity to common errors, such as missing body parts or misclassified clothing, and allows accurate mask prediction from as few as a single keypoint. Our results demonstrate that pose-guided fine-tuning of SAM enables effective, occlusion-aware human segmentation while preserving the generalization capabilities of the original model. The code and pretrained models will be available at the [project website](#)¹.

1. Introduction

A standard approach to human instance segmentation is to use a detector that predicts instance masks directly. However, in crowded scenes with heavy occlusion, these detectors often fail to separate overlapping instances. In contrast, human pose estimators are more robust in these conditions and produce structured keypoints that are easier to annotate and more stable under clutter.

Keypoint annotations are also significantly cheaper than per-pixel segmentations, making pose a practical intermediate representation for instance segmentation. In particular, human keypoints can serve as effective prompts for segmentation models. The task of *pose-guided human instance segmentation* takes either Ground-Truth or detected keypoints as input and outputs a segmentation mask for each instance.

This task was introduced with the OCHuman dataset [32] and the pose2seg model and was explored in many models [1, 2, 6, 27, 32, 33] since then. More recently, it has been used in the self-improving BMP loop (BBox-MaskPose, [19]), where pose-guided segmentation plays a key role in resolving multi-body ambiguity in heavily occluded scenes.

The Segment Anything Model (SAM, [11]) introduced a general segmentation framework trained on a large

¹MiraPurkrabek.github.io/BBox-Mask-Pose/

and diverse dataset of masks and images. SAM 2 [21] extended this with large-scale video training. While SAM shows strong generalization and has revolutionized prompted segmentation, it lacks semantic understanding and is not specialized for any class, such as humans.

We build on SAM’s generalization and introduce a method to adapt it for pose-guided human instance segmentation. Our model, SAM-pose2seg, incorporates semantic information into SAM through two main modifications. First, we fine-tune the decoder on human-only segmentation masks to specialize the predictions. Second, we replace SAM’s random point prompts with body keypoints during training. This aligns the prompt encoder and mask decoder more closely with the target task.

SAM-pose2seg is a pose-guided variant of SAM 2.1, fine-tuned for human instance segmentation. It introduces semantic specialization through decoder fine-tuning and aligns the prompt encoder with the task by training on human keypoints. This design enables robust segmentation from both the predicted pose and the Ground-Truth. In the remainder of the paper, we show that SAM-pose2seg achieves state-of-the-art performance on pose-guided human instance segmentation benchmarks and analyze design choices through ablation studies.

2. Related Work

Detection-Based Segmentation: The most direct approach to human instance segmentation uses a detector [9, 13, 16, 17, 22, 26, 26, 34]. A detector takes an image as input and outputs instance masks and class labels. Detectors work well in scenes similar to the training data, but struggle when instances overlap heavily. Under severe occlusion, they often merge multiple people into one mask, or assign different body parts to different instances. Tuning hyperparameters (e.g., non-maximum suppression) can reduce this, but increases false positives. **General Prompted Segmentation** requires an image and a human or automatically generated prompt. A prominent example is the SAM family [11, 21], trained on large image and video datasets with human-in-the-loop supervision. These models generalize well, but lack semantic understanding, making the task inherently ambiguous. As a result, they often segment skin, face, hair, clothing, or body parts instead of the full human instance.

Semantics-Aware Prompted Segmentation: Several works [8, 23, 28, 30] inject semantics into SAM. They take an image and a text prompt as input and detect and segment instances. This setting is often referred to as *open-vocabulary* or *zero-shot* segmentation. While these models capture semantics and segment entire instances, they do not provide localized cues and offer no control over which instance is chosen. Moreover, they are unsuitable for iterative loops such as BMP [19].

Pose-Aware Instance Segmentation is the most specific form of prompted segmentation. The model takes an image and a detected or Ground-Truth human pose and segments the corresponding person. These meth-

ods specialize in humans, sacrificing the generality of open-vocabulary or generic segmenters for robustness and precise, localized control. Test-time optimization methods [5, 6] refine predictions at inference without training. While training-free, they are computationally expensive and slow at test time. Standard pose-guided human segmentation methods [1, 2, 14, 27, 31–33] rely on small training datasets and thus generalize poorly, especially in crowded scenes. The closest related work is CrowdSAM [7], which builds a framework around SAM to automatically annotate bounding boxes in crowds. CrowdSAM uses SAM and DINO [18] as external tools and optimizes prompts, similar to [19]. In contrast, our SAM-pose2seg is end-to-end and, when used in an iterative loop, outperforms CrowdSAM, as shown in [19].

Datasets: There are not many datasets with annotated human instance segmentation masks and human poses. One of the first to offer this kind of data was COCO [15] and it became a standard for pose-guided segmentation training. COCO has two problems. First, it does not focus on overlapping instances, which is now the most challenging scenario (see [19]). Second, the human pose is annotated only for “large” instances, so the models are not trained or evaluated on small background persons. Overlapping instances were addressed in OCHuman dataset [32] but it is too small for training and contains only validation and testing sets. There are datasets such as [3, 10, 12, 25, 29] focusing on analysis of human body in crowded scenes, but none of these offer both segmentation and pose annotations. In this work, we worked with COCO [15] and CIHP [10] for training and OCHuman [32] for evaluation. For details about used datasets, see Sec. 4.

Iterative Methods: Pose-guided instance segmentation has been used in iterative pipelines [19, 24]. Our work builds on BMP [19], where instance segmentation is prompted by detected human keypoints. They use unmodified SAM 2 with a complex prompt selection strategy. We also build on SAM 2, but adapt it to this task. Our SAM-pose2seg is fine-tuned for human instances, removing the need for complex prompt selection. Compared to SAM 2 in BMP, SAM-pose2seg achieves better performance, with higher robustness and lower complexity.

3. Method

SAM-pose2seg builds on SAM2 [21], introducing task-specific modifications and a dedicated prompting strategy driven by pose keypoints. Our approach consists of three main components: adapting the base SAM architecture, fine-tuning with a pose-aware iterative sampling strategy, and designing an effective keypoint prompting mechanism for inference.

3.1. SAM

Segment Anything (SAM) is a powerful, robust tool with great generalization capacity. Its robustness to noise and point prompting capability offers a highly optimal foundation for our task. However, standard prompting tech-

niques often struggle in occluded scenarios, where detected keypoints lack unequivocal visibility information, making SAM unreliable.

Our goal is to adapt SAM (specifically, *SAM 2.1*) for pose-guided segmentation while retaining its ability to generalize. We introduce minimal modifications to the encoder structure to tailor the model specifically for this task.

3.2. Fine-tuning

We fine-tuned the *SAM 2.1 Hiera Base Plus* model using the official Meta training script, focusing specifically on the point selection mechanism. To preserve the generalization power of the backbone—trained on massive, diverse dataset—and to accelerate training, we froze the image encoder. Only the prompt encoder and mask decoder were optimized. No architectural changes were made to the mask decoder; however, it was also a part of the training process, as the prompt encoder’s only learned parameters are the weights for positive, negative, and bounding box embeddings. We adopted the sampling strategy *PoseMaskRefine*.

Method MaskRefine (Default Training Strategy).

Meta’s default training procedure serves as a strong baseline, which only builds on the Ground-Truth (GT) masks. It does not rely on pre-defined semantic or pose-based points, instead, they are sampled dynamically: the first is drawn uniformly from the Ground-Truth mask, while the remaining seven are sampled from the error region between the GT mask and the previous prediction (with a small probability, sampling is done from the GT mask instead). This totals eight points per iteration. By conditioning the model on both the previous logits and the newly sampled point, this method enables iterative refinement without dependence on explicit semantic cues.

Method PoseMaskRefine (Pose-Guided Refinement).

Building on the strong performance of MaskRefine, we introduce PoseMaskRefine to meet our specific task requirements. Here, the initial point is not sampled uniformly but is selected as an available pose keypoint with the highest visibility, unless no keypoint is available. In subsequent iterations, the remaining seven points are preferentially sampled from pose keypoints located within the current error region; with a small probability, sampling follows the original MaskRefine strategy and selects points uniformly from the Ground-Truth mask instead. If no such keypoints exist, sampling reverts to uniform selection from the error region. By incorporating pose cues while preserving the iterative, error-driven nature of MaskRefine, this strategy consistently outperformed the default approach (see Tab. 1).

The MaskRefine strategy effectively acts as hard negative mining; because samples are selected based on local error, boundary and occluded regions receive strong supervision. We attempted to simplify this by eliminating

the iterative correction process in favor of a pure keypoint-based prediction to reflect the inference prompting (see Sec. 4.3), but this yielded no significant improvement.

Fine-tuning with PoseMaskRefine makes the model significantly less sensitive to common errors, such as segmenting clothing only or omitting body parts (see Fig. 7).

3.3. Prompting

Effective prompting is crucial for interacting with SAM, yet pose keypoints do not perfectly align with the correction-based prompting concept inherent to SAM’s training. Therefore, our goal was to find the most effective way to pose prompting while sticking to the SAM’s structure.

Firstly, we incorporate the ProbPose [20] *visibility* metric as a reliable criterion for keypoint selection. This metric serves as a confidence measure for each keypoint, indicating whether the respective body part is observable in the picture and enabling the exclusion of keypoints that are likely occluded or noisy.

For the base model, selecting six keypoints based on both visibility and distance proved most stable (see Fig. 3), as this provided sufficient variability to recognize the whole body without the detriment caused by exceeding eight points.

However, across all fine-tuning experiments, PoseMaskRefine consistently yielded the best generalization and enabled a significant simplification of our inference strategy. Its flexibility allows us to replace complex heuristics with a simpler approach, reducing the risk of selecting occluded points (see Fig. 8). We now select only the **3 keypoints with the highest visibility scores**. This simplified strategy improved accuracy across all observed datasets using detected keypoints.

We employ different selection methods for the base SAM 2.1 and SAM-pose2seg in figures and tables, unless stated otherwise, to compare the optimal performance of each version.

4. Experiments

4.1. Implementation Details

Data. For training, we used the COCO and CIHP datasets together with ProbPose [20] keypoints to ensure a wide range of human poses and occlusion conditions, mirroring their main use at inference time. Training on COCO alone is not sufficient due to the pronounced domain shift relative to CIHP and OCHuman, both of which emphasize multi-person and heavily occluded scenarios. We generated detected pose keypoints for all instances (including small ones that are usually ignored for pose estimation) so that our model is robust to localization noise.

Architecture and Training Details. We fine-tuned the *SAM 2.1 Hiera Base Plus* checkpoint. We trained the prompt encoder and the mask decoder only, while the image encoder (backbone)

COCO	CIHP	train selection	test selection	# points	COCO AP	OCH AP	CIHP AP
✗	✗	mV	mV	6	37.7	29.4	66.4
✗	✗	mS	mS	6	41.2	29.5	71.6
✓	✗	mV	mV	6	41.2	27.0	62.2
✓	✗	mS	mS	6	42.6	26.1	64.3
✓	✓	mS	mS	6	43.3	29.3	67.7
✓	✗	mR	mV	3	43.7	29.8	68.9
✓	✓	mR	mV	3	43.7	34.1	71.7
✓	✓	P1mR	mV	3	44.5	34.5	72.7
✓	✓	PmR	mV	3	44.6	34.7	72.7

Table 1. **Ablation study: fine-tuning recipes.** Keypoint selection methods are MaxVis (mV), MaxSpread (mS), MaskRefine (mR), PoseMaskRefine (PmR) and Pose1MaskRefine (P1mR). First two rows are baselines SAM 2.1 without any fine-tuning. For mV and mS methods, the optimal # points is 6, for mR-based methods it is 3. The best point method is PmR with 3 points. Training on both COCO and CIHP is crucial for generalization to unseen OCHuman.

remained frozen. The configuration file, `SAM 2.1.hiera_b+MOSE_finetune.yaml`, was taken from the SAM 2 repository. All hyperparameters not mentioned here were left unchanged. The number of frames was set to 1. We trained for 15 epochs, though the model converged earlier; performance changes were marginal when varying epochs between 10 and 30. The probability of bounding box usage (`prob_to_use_box_input_for_train`) was set to 0.0, while the probability of point usage (`prob_to_use_pt_input_for_train`) was set to 1.0.

4.2. SOTA Comparison

The proposed *SAM-pose2seg* model represents a strong all-round solution for the pose-to-segmentation task. Across a wide range of datasets and evaluation settings, it consistently achieves competitive or superior performance. In particular, *SAM-pose2seg* demonstrates robust behavior under both detected and Ground-Truth pose inputs, making it well suited for practical deployment scenarios.

Using detected ProbPose [20] keypoints, we achieve **44.6 AP** on COCO val2017, **60.3 AP** on COCOPersons (COCO val2017 excluding small instances), **34.7 AP** on OCHuman test, and **72.7 AP** on CIHP val (for comparison with SOTA, see Tab. 2). These results demonstrate strong generalization across datasets with varying levels of occlusion and pose complexity, with particularly strong performance on the challenging OCHuman benchmark.

When evaluated with Ground-Truth poses, *SAM-pose2seg* further improves performance, reaching **61.6 AP** on COCOPersons, **70.0 AP** on OCHuman test, and **69.5 AP** on OCHuman val (for comparison with SOTA, see Tab. 3). For Ground-Truth keypoints, we apply spread-based keypoint selection, as their binary visibility annota-

Model (Prompting)	COCO val AP	COCOPersons val AP	OCHuman test AP
HQNet R-50	-	-	31.1
Pose2Seg	-	55.5	23.8
ExPoSeg ⁺	-	61.9	26.8
Occlusion C&P ⁺	-	-	28.3
Crowd-SAM ⁺	22.0	-	31.4
MultiPoseSeg	-	56.3	-
SAM 2.1	41.2	56.0	29.5
<i>SAM-pose2seg</i>	44.6	60.3	34.7

Table 2. **Evaluation of the base SAM model and *SAM-pose2seg* on detected ProbPose [20] keypoints.** The improvement is most visible on the OCHuman test set. Prompting methods are explained in Sec. 4.3.2. COCOPersons is a subset of COCO without Small category persons (no Ground-Truth poses are available there) for comparison with Pose2Seg. Models labeled with ⁺ estimate either masks or report detection AP. Every result except for SAM 2.1 was achieved on a different set of detected keypoints.

Model (Prompting)	COCOPersons val AP	OCHuman test AP	OCHuman val AP
Pose2Seg	58.2	55.2	54.4
base SAM 2.1	57.1	70.1	70.2
<i>SAM-pose2seg</i>	61.6	70.0	69.5

Table 3. **Evaluation of the base SAM model and *SAM-pose2seg* AP on the Ground-Truth pose keypoints.** MaxVis strategy does not make sense as the visibility classification is binary here, we therefore use selection by distance to optimize, while sorting out all keypoints with visibility set to zero. COCOPersons is a subset of COCO without Small category persons (no Ground-Truth poses are available there). MaxSpread₃ keypoint selection is used for *SAM-pose2seg*.

tions do not allow for visibility-based ranking.

4.3. Ablation Study

4.3.1. Correction Points In Training

We further review the correction process to prove why our model converges to the optimal usage of three keypoints during inference. Crucially, as the model adapts to human segmentation, it rarely requires all seven correction points. In practice, the fine-tuned *SAM-pose2seg* model often predicts a correct mask after the first sampled keypoint. When comparing Ground-Truth masks to masks in each of the correction iteration, it was discovered that the differences changes between the mask in each iteration and their IoU with the Ground-Truth mask are small – mean IoU on a small subset during before any correction iterations is already 81.0 % and only goes up by 6.4 percentage points after the last iteration (see Fig. 2). This is corroborated by our *Pose1MaskRefine* experiment, where only the first point is pose-derived, and subsequent correction points are sampled uniformly as in the default method. The marginal difference in AP between

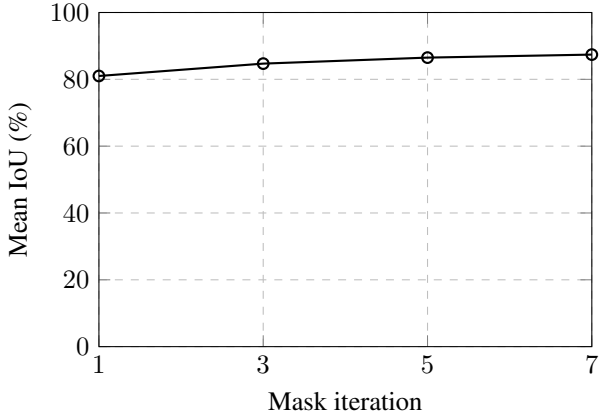


Figure 2. **Importance of correction points in the default SAM training method MaskRefine.** Comparison on a small COCO + CIHP training subset of mean IoU of the ground-truth mask. The first iteration contains only one point, any other i -th iteration is prompted by i keypoints and a correction mask. It seems that the refinement mostly causes minor changes in the overall mask shape.

Pose1MaskRefine and the full PoseMaskRefine confirms that the first keypoint is the primary driver of performance (see Tab. 1).

4.3.2. Prompting – Keypoint Selection

We evaluate two prompting strategies, MaxVis and MaxSpread, which differ in keypoint selection and interaction with Ground-Truth (GT) masks during inference and training (yet we stuck with PoseMaskRefine in the case of training).

Method MaxVis _{n} (Visibility-Based Keypoint Selection). This method selects the top n keypoints solely based on *visibility* scores, imposing no additional spatial, semantic, or structural constraints.

Method MaxSpread _{n} (Distance- and Visibility-Based Keypoint Selection). Originally introduced in Bbox-Mask-Pose [19] (specifically MaxSpread₆), this method extends visibility-based selection with a spatial diversification heuristic inspired by *k-means++* [4]. The first keypoint is chosen for maximum visibility; subsequent points are selected to maximize distance from previous ones. To reduce redundancy, number of used eye and nose keypoints is limited to one, and low-visibility points are excluded. This strategy generally outperforms MaxVis on the default model and is preferred for binary visibility classification.

MaxSpread₆ proved most effective for pose-to-segmentation, remaining consistent after incorporating ProbPose [20] visibility. For the base model, spatially distributed keypoints provide richer cues than simple visibility selection, which often overrepresents facial points and causes segmentation ambiguity (e.g., face vs. full body).

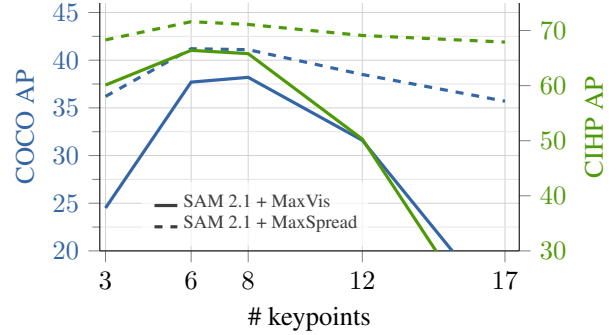


Figure 3. **Keypoint selection methods on SAM 2.1.** Prompting methods MaxVis (full) and MaxSpread (dashed) on COCO and CIHP datasets. 6 keypoints is the best for both methods. MaxSpread outperforms MaxVis as shown in [19].

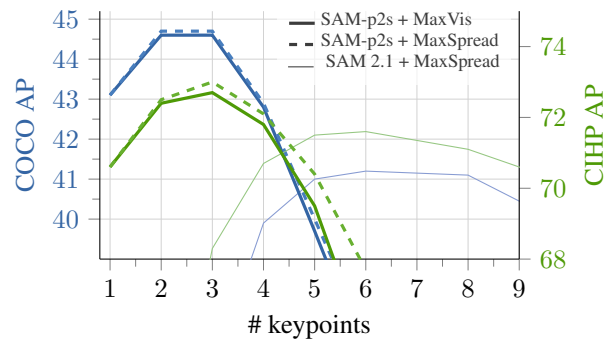


Figure 4. **Keypoint selection methods for SAM 2.1 and SAM-pose2seg on COCO and CIHP datasets.** SAM 2.1 with MaxSpread (thin line) peaks at 6 keypoints for both datasets. SAM-pose2seg with MaxSpread (dashed) and MaxVis (full) behaves the same on both datasets and peaks at 3 keypoints. We selected MaxVis method due to its simplicity. SAM-pose2seg outperforms SAM 2.1 with both selection methods on both datasets.

However, after fine-tuning with MaskRefine and PoseMaskRefine, the performance gap between MaxSpread _{n} and MaxVis _{n} becomes negligible (see Fig. 4).

We also briefly explored negative keypoint prompting; however, as it did not yield consistent gains and introduces strong context dependence, we defer a detailed analysis to the Supplementary.

4.3.3. Bounding Boxes

Bounding boxes are not part of our final pipeline, but we analyze their impact as an ablation to better understand SAM’s behavior under pose-guided prompting.

The SAM model does support the use of bounding boxes as prompts, and since ProbPose [20] provides bounding box predictions, it would be natural to incorporate them into our pipeline. Ground-Truth bounding boxes were found to be beneficial, as they define precise instance boundaries (see Tab. 4). Even in this case, difficult settings, where a correctly identified bounding box contains large parts of multiple people, might pose a challenge (see Fig. 5). Nevertheless, bounding boxes predicted by vari-



Figure 5. **Problematic bounding box usage** in SAM 2 when multiple people are included. Note: a part of the other person’s hand is recognized incorrectly in both cases due to an incorrect pose keypoint.

Model	Bbox type	COCO AP	CIHP AP
SAM 2.1	GT	50.5	76.4
SAM 2.1	inflated GT	17.5	52.4
SAM 2.1	none	41.2	71.6
SAM 1	GT	52.5	72.2
SAM 1	inflated GT	45.2	59.2
SAM 1	none	43.0	65.6

Table 4. **Usage of bounding boxes in SAM 1 and SAM 2.1.** on COCO val and CIHP val. Performance comparison of the base SAM 1 and SAM 2.1 models if ProbPose [20] keypoints are accompanied by bounding boxes. We use the keypoint selection method MaxSpread₆. While Ground-Truth bounding boxes seem to be helpful in both scenarios, if both dimensions are enlarged to simulate possible detector noise (by 50 % in each direction – 400 % area increase in total), SAM 2.1 is not reliable.

ous detectors (that would be a part of iterative pose refinement) may not be as precise [19]. As a result, segmentation precision degrades significantly when using SAM 2.1: the model may incorrectly include parts of other people or background regions in order to fill the provided bounding box (see Fig. 6).

4.3.4. Backbone Choice: SAM 1 vs. SAM 2

As we tried to analyze whether we could improve our prompting method, we also turned our attention to the older SAM 1 model. In both cases, the second largest models were used (*Hiera Base Plus* for SAM 2 and *Vit-L* for SAM 1). We conducted several tests showing its strengths and weaknesses, and it seems that for our task, SAM 2.1 is more suitable. SAM 1 performs slightly better

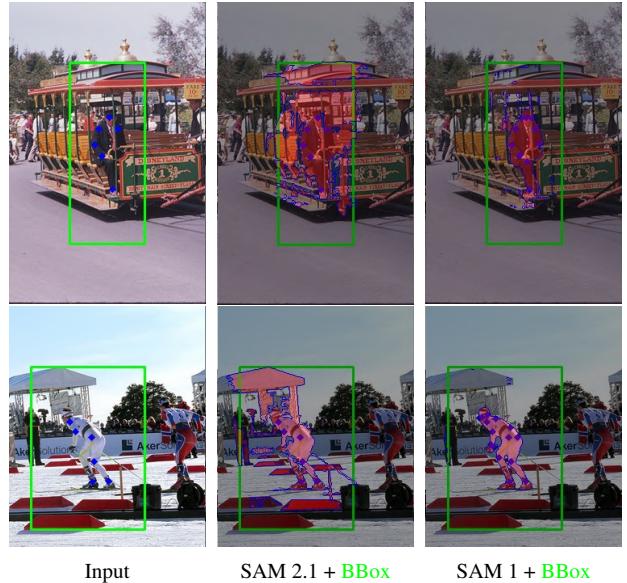


Figure 6. **Inflated bounding box usage** in base SAM 2.1 (middle) and SAM 1 (right). SAM 2.1 often incorporates unnecessary noise when the prompted bounding box does not fit the person exactly. In contrast, SAM 1 performance generally does not drop drastically.

on COCO dataset, where the poses are more visible, yet lags behind in OCHuman and CIHP datasets.

It also seems that SAM 1 is able to make better use of negative keypoints in the case of OCHuman and CIHP datasets. When sampling the closest pose keypoints, it serves better as an option to remove superficial parts. Imprecise bounding boxes also do not harm the prediction process significantly, which could not be said for SAM 2.1 that often includes other objects just to fill the boundaries.

However, in general, SAM 2.1 turned out to be a stronger base for our task, as the model itself is lighter (and, therefore, much faster during inference), the provided training code serves well for the fine-tuning and the experiments generally yielded worse results.

4.3.5. SAM-pose2seg Performance In Challenging Scenes

Incomplete segmentation. In contrast to the base SAM 2.1, incomplete segmentation occurs far less frequently. Since SAM-pose2seg is fine-tuned specifically for human segmentation, it learns to recover the full human body more consistently. This addresses a common ambiguity of the base SAM model, which in some cases produced masks corresponding only to clothing, skin regions, or isolated body parts rather than complete human instances.

Partially occluded pose. SAM-pose2seg demonstrates precise segmentation in scenarios where pose information is unreliable or heavily occluded. By prioritizing keypoints with the highest visibility scores, the model effectively leverages strong, reliable cues, reducing the impact of missing or ambiguous points. In cases where a



Figure 7. **Improved handling of incomplete segmentation**, with SAM-pose2seg correctly segmenting full human instances. More examples in Supplementary.

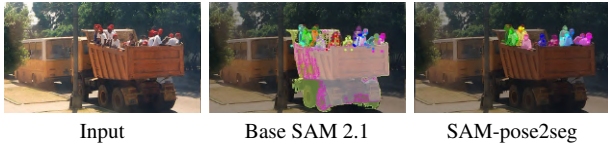


Figure 8. **Improved handling of partially occluded poses**. More examples in Supplementary.

single additional misdetected keypoint could compromise the segmentation, SAM-pose2seg maintains accuracy and produces consistent, complete masks.

5. Conclusions

To sum up, we achieved substantial improvements in the pose-to-segmentation task for humans by adapting and fine-tuning SAM. We identified and validated the most effective keypoint selection strategies, demonstrating that leveraging *visibility* scores allows the model to reliably choose strong, informative points. Simple and clear prompting methods can be particularly useful in occluded or crowded scenarios.

Our proposed SAM-pose2seg model shows strong generalization and robustness. It excels at recovering full human instances even when most keypoints are missing or unreliable, and is less sensitive to errors that would otherwise harm segmentation in the base SAM model. This makes it suitable both for iterative refinement in pose estimation pipelines and as a standalone tool for human instance segmentation. The deployment of SAM-pose2seg in the BMP loop [19] improves the segmentation accuracy from 31.8 AP to 33.7 AP on OCHuman.

Despite these improvements, some limitations remain. SAM-pose2seg can struggle when multiple people are closely overlapping and even high-*visibility* keypoints carry ambiguous semantic meaning, occasionally merging visually similar regions across instances. Future work could explore adaptive keypoint weighting, explicitly modeling the semantic meaning of each keypoint, or integrating the approach into SAM3 to leverage its language-based reasoning capabilities.

Overall, our results demonstrate that SAM-pose2seg provides a practical and reliable approach for human pose-guided segmentation, with potential applications as a pseudo-annotation tool for large-scale datasets or as a robust component in downstream vision pipelines.

Acknowledgements. This work was supported by the Ministry of the Interior of the Czech Republic project No. VJ02010041, the Technology Agency of the Czech Republic project CEDMO 2.0 No. FW10010387, the European Union’s Digital Europe Programme under Contract No. 101158609, and the Czech Technical University student grant SGS23/173/OHK3/3T/13.

References

- [1] Niaz Ahmad, Jawad Khan, Jeremy Yuhyun Kim, and Youngmoon Lee. Multiposeseg: Feedback knowledge transfer for multi-person pose estimation and instance segmentation. *2022 26th International Conference on Pattern Recognition (ICPR)*, pages 2086–2092, 2022. 1, 2
- [2] Niaz Ahmad, Jawad Khan, Jeremy Yuhyun Kim, and Youngmoon Lee. Joint human pose estimation and instance segmentation with poseplusseg. In *AAAI Conference on Artificial Intelligence*, 2022. 1, 2
- [3] 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)*, 2014. 2
- [4] David Arthur and Sergei Vassilvitskii. k-means++: The advantages of careful seeding. Technical report, Stanford, 2006. 5
- [5] Kambiz Azarian, Debasmit Das, Hyojin Park, and Fatih Porikli. Test-time adaptation vs. training-time generalization: A case study in human instance segmentation using keypoints estimation, 2022. 2
- [6] Kambiz Azarian, Debasmit Das, Hyojin Park, and Fatih Murat Porikli. Test-time adaptation vs. training-time generalization: A case study in human instance segmentation using keypoints estimation. *2023 IEEE/CVF Winter Conference on Applications of Computer Vision Workshops (WACVW)*, pages 411–420, 2022. 1, 2
- [7] Zhi Cai, Yingjie Gao, Yaoyan Zheng, Nan Zhou, and Di Huang. Crowd-sam: Sam as a smart annotator for object detection in crowded scenes. In *Proceedings of the European Conference on Computer Vision (ECCV)*, 2024. 2
- [8] Claudia Cattano, Gabriele Trivigno, Gabriele Rosi, Carlo Masone, and Giuseppe Averta. Samwise: Infusing wisdom in sam2 for text-driven video segmentation. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 3395–3405, 2025. 2
- [9] Z Ge. Yolox: Exceeding yolo series in 2021. *arXiv preprint arXiv:2107.08430*, 2021. 2
- [10] Ke Gong, Xiaodan Liang, Yicheng Li, Yimin Chen, Ming Yang, and Liang Lin. Instance-level human parsing via part grouping network. In *Proceedings of the European conference on computer vision (ECCV)*, pages 770–785, 2018. 2
- [11] Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete Xiao, Spencer Whitehead, Alexander C Berg, Wan-Yen Lo, et al. Segment anything. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 4015–4026, 2023. 1, 2
- [12] Jiefeng Li, Can Wang, Hao Zhu, Yihuan Mao, Hao-Shu Fang, and Cewu Lu. Crowdpose: Efficient crowded scenes pose estimation and a new benchmark. *arXiv preprint arXiv:1812.00324*, 2018. 2

- [13] Yanghao Li, Hanzi Mao, Ross Girshick, and Kaiming He. Exploring plain vision transformer backbones for object detection. In *European conference on computer vision*, pages 280–296. Springer, 2022. 2
- [14] Zhong Li, Xin Chen, Wangyiteng Zhou, Yingliang Zhang, and Jingyi Yu. Pose2body: Pose-guided human parts segmentation. In *2019 IEEE International Conference on Multimedia and Expo (ICME)*, pages 640–645, 2019. 2
- [15] Tsung-Yi Lin, Michael Maire, Serge J. 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*, 2014. 2
- [16] Zhuang Liu, Hanzi Mao, Chao-Yuan Wu, Christoph Feichtenhofer, Trevor Darrell, and Saining Xie. A convnet for the 2020s. *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022. 2
- [17] Chengqi Lyu, Wenwei Zhang, Haiyan Huang, Yue Zhou, Yudong Wang, Yanyi Liu, Shilong Zhang, and Kai Chen. RtmDET: An empirical study of designing real-time object detectors. *ArXiv*, abs/2212.07784, 2022. 2
- [18] Maxime Oquab, Timothée Darcet, Théo Moutakanni, Huy Vo, Marc Szafraniec, Vasil Khalidov, Pierre Fernandez, Daniel Haziza, Francisco Massa, Alaaeldin El-Nouby, et al. DINOv2: Learning robust visual features without supervision. *arXiv preprint arXiv:2304.07193*, 2023. 2
- [19] Miroslav Purkrabek and Jiri Matas. Detection, pose estimation and segmentation for multiple bodies: Closing the virtuous circle. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 9004–9013, 2025. 1, 2, 5, 6, 7
- [20] Miroslav Purkrabek and Jiri Matas. ProbPose: A probabilistic approach to 2d human pose estimation. In *Proceedings of the Computer Vision and Pattern Recognition Conference (CVPR)*, pages 27124–27133, 2025. 3, 4, 5, 6
- [21] Nikhila Ravi, Valentin Gabeur, Yuan-Ting Hu, Ronghang Hu, Chaitanya Ryali, Tengyu Ma, Haitham Khedr, Roman Rädle, Chloe Rolland, Laura Gustafson, et al. SAM 2: Segment anything in images and videos. *arXiv preprint arXiv:2408.00714*, 2024. 1, 2
- [22] Joseph Redmon and Ali Farhadi. YOLOv3: An incremental improvement. *arXiv preprint arXiv:1804.02767*, 2018. 2
- [23] Tianhe Ren, Shilong Liu, Ailing Zeng, Jing Lin, Kunchang Li, He Cao, Jiayu Chen, Xinyu Huang, Yukang Chen, Feng Yan, et al. Grounded sam: Assembling open-world models for diverse visual tasks. *arXiv preprint arXiv:2401.14159*, 2024. 2
- [24] Danila Rukhovich, Konstantin Sofiiuk, Danil Galeev, Olga Barinova, and Anton Konushin. IterDET: iterative scheme for object detection in crowded environments. In *Structural, syntactic, and statistical pattern recognition: Joint IAPR international workshops, s+ SSPR 2020, padua, Italy, January 21–22, 2021, proceedings*, pages 344–354. Springer, 2021. 2
- [25] Shuai Shao, Zijian Zhao, Boxun Li, Tete Xiao, Gang Yu, Xiangyu Zhang, and Jian Sun. CrowdHuman: A benchmark for detecting human in a crowd. *ArXiv*, abs/1805.00123, 2018. 2
- [26] Ke Sun, Bin Xiao, Dong Liu, and Jingdong Wang. Deep high-resolution representation learning for human pose estimation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 5693–5703, 2019. 2
- [27] Subarna Tripathi, Maxwell D. Collins, Matthew A. Brown, and Serge J. Belongie. Pose2Instance: Harnessing keypoints for person instance segmentation. *ArXiv*, abs/1704.01152, 2017. 1, 2
- [28] Zhaoyang Wei, Pengfei Chen, Xuehui Yu, Guorong Li, Jianbin Jiao, and Zhenjun Han. Semantic-aware sam for point-prompted instance segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 3585–3594, 2024. 2
- [29] Jiahong Wu, He Zheng, Bo Zhao, Yixin Li, Baoming Yan, Rui Liang, Wenjia Wang, Shipai Zhou, Guosen Lin, Yanwei Fu, Yizhou Wang, and Yonggang Wang. AI Challenger: A large-scale dataset for going deeper in image understanding. *ArXiv*, abs/1711.06475, 2017. 2
- [30] Shiting Xiao, Rishabh Kabra, Yuhang Li, Donghyun Lee, Joao Carreira, and Priyadarshini Panda. OpenWorldSAM: Extending sam2 for universal image segmentation with language prompts. *arXiv preprint arXiv:2507.05427*, 2025. 2
- [31] Lu Yang, Qing Song, Zhihui Wang, and Ming Jiang. Parsing r-cnn for instance-level human analysis. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019. 2
- [32] Song-Hai Zhang, Ruilong Li, Xin Dong, Paul Rosin, Zixi Cai, Xi Han, Dingcheng Yang, Haozhi Huang, and Shi-Min Hu. Pose2Seg: Detection free human instance segmentation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 889–898, 2019. 1, 2
- [33] Desen Zhou and Qian He. Poseg: Pose-aware refinement network for human instance segmentation. *IEEE Access*, 8:15007–15016, 2020. 1, 2
- [34] Zhuofan Zong, Guanglu Song, and Yu Liu. DETRS with collaborative hybrid assignments training. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 6748–6758, 2023. 2