

# SDPose: EXPLOITING DIFFUSION PRIORS FOR OUT-OF-DOMAIN AND ROBUST POSE ESTIMATION

000  
001  
002  
003  
004  
005 **Anonymous authors**  
006 Paper under double-blind review  
007  
008  
009  
010  
011  
012  
013  
014  
015  
016  
017  
018  
019  
020  
021  
022  
023  
024  
025  
026  
027  
028  
029  
030  
031  
032  
033  
034  
035  
036  
037  
038  
039  
040  
041  
042  
043  
044  
045  
046  
047  
048  
049  
050  
051  
052  
053  
SDPose: EXPLOITING DIFFUSION PRIORS FOR  
OUT-OF-DOMAIN AND ROBUST POSE ESTIMATION



Figure 1: **SDPose: OOD-robust pose via diffusion priors.** On stylized paintings, SDOurs surpasses Sapiens and ViTPose++-H, matching SoTA on COCO and setting new records on HumanArt and COCO-OOD; yellow boxes show baseline failures.

## ABSTRACT

Pre-trained diffusion models provide rich multi-scale latent features and are emerging as powerful vision backbones. While recent works such as Marigold (Ke et al., 2024) and Lotus (He et al., 2024) adapt diffusion priors for dense prediction with strong cross-domain generalization, their poten-

054        tial for structured outputs (e.g., human pose estimation) remains underex-  
 055        plored. In this paper, we propose **SDPose**, a fine-tuning framework built  
 056        upon Stable Diffusion to fully exploit pre-trained diffusion priors for human  
 057        pose estimation. First, rather than modifying cross-attention modules or  
 058        introducing learnable embeddings, we directly predict keypoint heatmaps  
 059        in the SD U-Net’s image latent space to preserve the original generative pri-  
 060        ors. Second, we map these latent features into keypoint heatmaps through a  
 061        lightweight convolutional pose head, which avoids disrupting the pre-trained  
 062        backbone. Finally, to prevent overfitting and enhance out-of-distribution  
 063        robustness, we incorporate an auxiliary RGB reconstruction branch that  
 064        preserves domain-transferable generative semantics. To evaluate robustness  
 065        under domain shift, we further construct **COCO-OOD**, a style-transferred  
 066        variant of COCO with preserved annotations. With just one-fifth of the  
 067        training schedule used by Sapiens on COCO, SDPose attains parity with  
 068        Sapiens-1B/2B on the COCO validation set and establishes a new state  
 069        of the art on the cross-domain benchmarks HumanArt and COCO-OOD.  
 070        **Extensive ablations highlight the importance of diffusion priors, RGB re-**  
 071        **construction, and multi-scale SD U-Net features for cross-domain general-**  
 072        **ization, and t-SNE analyses further explain SD’s domain-invariant latent**  
 073        **structure. We also show that SDPose serves as an effective zero-shot pose**  
 074        **annotator for controllable image and video generation.**

## 1 INTRODUCTION

078        With the recent rise of embodied AI, video generation, and 3D asset rendering, the need for  
 079        cross-domain-robust human pose estimation has become critical in robotics as well as in film,  
 080        animation, and game production. Although recent advances on academic benchmarks such  
 081        as MS COCO (Lin et al., 2014) using models such as DWPose (Yang et al., 2023), RTM-  
 082        Pose (Jiang et al., 2023) and OpenPose (Martinez, 2019), as well as approaches leveraging  
 083        large pretrained backbones such as ViTPose (Xu et al., 2022; 2023) and Sapiens (Khirodkar  
 084        et al., 2024), have achieved strong in-domain accuracy, they often exhibit severe performance  
 085        degradation under domain shifts and require substantial fine-tuning efforts.

086        Recently, pre-trained diffusion models such as Stable Diffusion (Rombach et al., 2022) have  
 087        emerged as robust vision backbones. A growing body of work has shown that with fine-  
 088        tuning and adaptation, diffusion priors can be repurposed for 3D generation (Cheng et al.,  
 089        2023; Lin et al., 2025; Long et al., 2024), segmentation (Karmann & Urfalioglu, 2025), and  
 090        dense prediction tasks (Ke et al., 2024; He et al., 2024), while consistently demonstrating  
 091        strong cross-domain robustness and highlighting their potential for leveraging intra-visual  
 092        multimodality in generative priors. However, their potential for structured and semantically  
 093        aware outputs, particularly in human pose estimation, remains largely unexplored. Concur-  
 094        rent efforts like GenLoc (Wang et al., 2025a) and Diff-Tracker (Zhang et al., 2024) indicate  
 095        that generative priors can benefit keypoint localization and tracking by steering learnable  
 096        condition embeddings and adapting the diffusion model’s cross-attention. We instead ex-  
 097        amine a complementary axis: can one rely purely on SD U-Net latent features, without  
 098        attention read-outs or condition tokens, to produce reliable pose heatmaps?

099        **To bridge this gap and investigate how SD’s rich latent representations can be effectively**  
 100        **leveraged for robust cross-domain pose estimation, our contributions are as follows:**

101        (1) We propose **SDPose**, a fine-tuning framework with three key components: **(i) Latent-**  
 102        **space preservation.** We operate entirely in the SD U-Net’s image latent space without  
 103        modifying cross-attention modules or adding learnable embeddings, thus preserving pre-  
 104        trained visual semantics and feature geometry. **(ii) Lightweight pose decoder head.**  
 105        We introduce a minimal decoder that maps SD U-Net features to keypoint heatmaps with  
 106        only a shallow convolutional head, ensuring low overhead and minimal disturbance to the  
 107        pretrained latent representations. **(iii) RGB reconstruction regularization.** We add  
 108        an auxiliary reconstruction branch that regularizes fine-tuning, helping to maintain domain-  
 109        transferable generative semantics and improve out-of-distribution robustness.

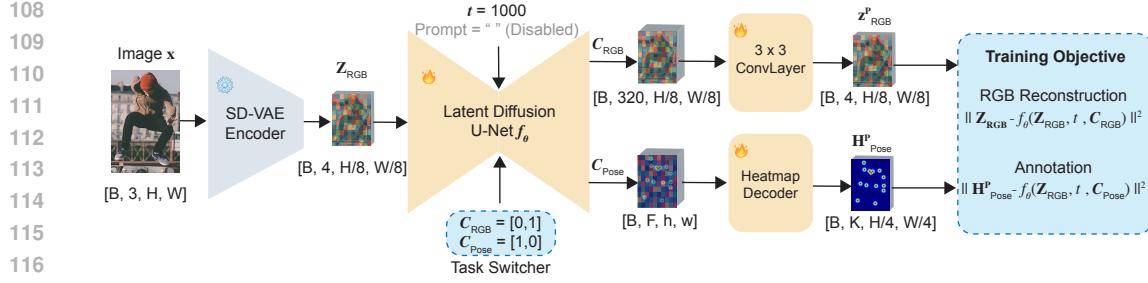


Figure 2: **Training Pipeline of SDPose.** The input RGB image is first encoded into the latent space by a pre-trained VAE. The U-Net is conditioned for multi-task learning via a class embedding. When the class label is set to  $[0,1]$ , the U-Net predicts the reconstructed RGB latent; when set to  $[1,0]$ , it produces features for heatmap prediction. The output layer of the U-Net is task-specific: the original convolutional output layer is retained for RGB latent reconstruction, while a lightweight heatmap decoder is used to process the U-Net’s intermediate features for keypoint heatmap prediction.

(2) To systematically evaluate robustness under domain shift, we introduce **COCO-OOD**, a style-transferred extension of COCO that includes oil-painting, ukiyo-e, and color sketch domains. This dataset fills an important gap in benchmarking generalization robustness.

(3) We conduct extensive ablation studies to understand how diffusion priors and our RGB reconstruction branch contribute to cross-domain generalization in pose estimation. We further compare multi-scale features from different upsampling blocks of the SD-UUnet, identifying the feature level that yields the strongest robustness under artistic domain shifts. Finally, through a latent-space comparison with Sapiens using t-SNE visualizations, we observe that SD’s pretrained latent features naturally capture domain-invariant structures, which is highly beneficial for cross-domain perception tasks.

On COCO (Lin et al., 2014) and COCO-WholeBody (Jin et al., 2020), SDPose delivers in-domain performance on par with the current SoTA, Sapiens (Khirodkar et al., 2024). Under domain shift (HumanArt (Ju et al., 2023), COCO-OOD), it sets a new state of the art while using only one-fifth of Sapiens’s fine-tuning epochs, highlighting the efficiency and cross-domain robustness of generative priors. Beyond quantitative benchmarks, we further demonstrate SDPose as a zero-shot pose annotator for downstream controllable generation tasks, including ControlNet-based image synthesis and video generation, where it provides reliable and qualitatively superior pose guidance.

## 2 RELATED WORKS

### 2.1 LATENT DIFFUSION MODELS

Latent diffusion models (LDMs), built on DDPM and further advanced by ODE and SDE samplers (Ho et al., 2020; Song et al., 2020b;a; Lu et al., 2022; Rombach et al., 2022), have gained traction over the past few years. Classic architectures, such as the UNet-based Stable Diffusion and Diffusion Transformers (DiT) (Peebles & Xie, 2023; Esser et al., 2024), have demonstrated strong performance across diverse conditional generation tasks (Zhang et al., 2023). Pretrained on large-scale datasets such as LAION-5B (Schuhmann et al., 2022), generative models like Stable Diffusion provide rich visual priors that can be effectively leveraged for a wide range of tasks. Recent advances in flow-matching (Lipman et al., 2022; Esser et al., 2024; Xie et al., 2024) further show that latent diffusion models can achieve high-quality synthesis with only a few sampling steps. These developments highlight the power of latent generative priors as a strong visual foundation.

### 2.2 LEVERAGING DIFFUSION PRIORS FOR PREDICTION TASKS

162 A growing body of work has explored repurposing  
 163 pretrained latent diffusion priors for dense  
 164 prediction tasks. Marigold (Ke et al., 2024)  
 165 adapts Stable Diffusion by fine-tuning only the  
 166 denoising U-Net using synthetic data, delivering  
 167 high-quality depth results. Later, subsequent  
 168 methods such as Lotus (He et al., 2024) and  
 169 GenPercept (Xu et al., 2024) both adopt a deter-  
 170 ministic one-step fine-tuning strategy, removing

171 the multi-step stochastic diffusion process and directly predicting task annotations, which  
 172 significantly improves both accuracy and inference speed. In contrast, leveraging latent  
 173 diffusion priors for structured outputs (e.g., human pose) remains underexplored. Prior  
 174 works such as GenLoc (Wang et al., 2025a) and Diff-Tracker (Zhang et al., 2024) freeze  
 175 Stable Diffusion backbone and use learnable condition or prompt embeddings to read cross-  
 176 attention maps, rather than decoding from U-Net latent features, aiming at zero-/few-shot  
 177 and schema-flexible generalization. We instead remain in the image latent space, treat the  
 178 SD U-Net as a multi-scale backbone, and attach a minimal convolutional head to produce  
 179 keypoint heatmaps, retaining SD-native visual semantics and improves robustness under  
 domain shift.

### 180 2.3 HUMAN POSE ESTIMATION

181 Human pose estimation is a classic and fundamental task in computer vision. Early ap-  
 182 proaches predominantly relied on CNN backbones such as HRNet (Sun et al., 2020) and  
 183 CSPNeXt (Chen et al., 2024), coupled with heuristically designed decoding heads. Mod-  
 184 els like RTMPose (Jiang et al., 2023) and DWPose (Yang et al., 2023) achieved strong  
 185 performance on academic benchmarks such as COCO and COCO WholeBody. However,  
 186 these models exhibit limited generalization when transferring from real human figures to  
 187 out-of-domain cases, such as anime characters. More recently, fine-tuned methods built  
 188 on extensive pre-trained vision backbones, such as ViTPose (Xu et al., 2022; 2023) and  
 189 Sapiens (Khirodkar et al., 2024), have achieved SoTA results on standard benchmarks,  
 190 demonstrating the benefit of leveraging pre-trained foundation models for pose estimation.  
 191 Nevertheless, these methods incur high fine-tuning costs, as they require large task-specific  
 192 datasets and lengthy training schedules to achieve competitive performance. In this paper,  
 193 we demonstrate that fine-tuning the Stable Diffusion pipeline with minimal architectural  
 194 modifications can address both the generalization gap and the high fine-tuning cost.

## 196 3 PRELIMINARIES

### 197 3.1 HEATMAP REPRESENTATION AND UNBIASED DATA PROCESSING (UDP)

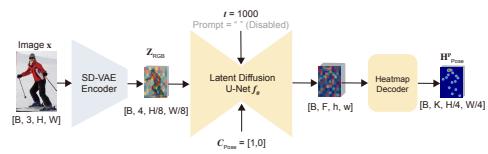
198 Let  $(x_i, y_i)$  denote the  $i$ -th ground-truth keypoint in an  $H \times W$  image. The standard heatmap  
 199 representation encodes each keypoint as

$$200 H_i(u, v) = \exp\left(-\frac{(u - x_i)^2 + (v - y_i)^2}{2\sigma^2}\right), \quad (\hat{x}_i, \hat{y}_i) = \arg \max_{u, v} H_i(u, v).$$

201 While widely adopted, this discrete pixel-space formulation suffers from quantization bias:  
 202 predicted coordinates become misaligned under flips, scales, or rotations since the argmax  
 203 operation only yields integer positions. To address this issue, we adopt the Unbiased Data  
 204 Processing (UDP) method (Huang et al., 2020), which removes quantization bias by es-  
 205 timating keypoints in a continuous domain. Following common practice, the heatmap is  
 206 generated at one-quarter resolution of the input image, which balances localization accu-  
 207 racy with computational efficiency.

### 208 3.2 PARAMETERIZATION FOR LATENT DIFFUSION MODEL

209 Traditional latent diffusion models (LDMs) (Ho et al., 2020) adopt the  $\epsilon$ -prediction pa-  
 210 rameterization, where the denoiser  $f_\theta$  is trained to predict the Gaussian noise  $\epsilon_t$  added at



211 Figure 3: **SDPose Inference Pipeline**.

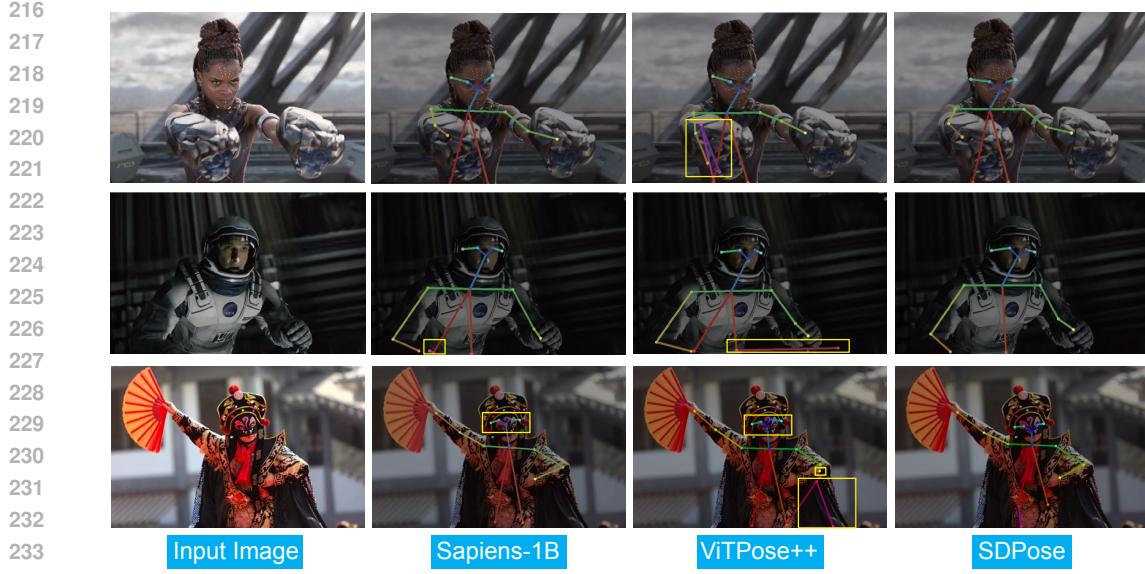


Figure 4: **Qualitative results on real-world photographs.** The yellow boxes highlight regions where baselines fail to predict accurate poses.

timestep  $t$ :

$$\hat{\epsilon}_t = f_\theta(z_t, t),$$

with  $z_t$  denoting the noisy latent at step  $t$ . The clean latent  $x_0$  can then be recovered by

$$\hat{x}_0 = \frac{z_t - \sqrt{1 - \alpha_t} \hat{\epsilon}_t}{\sqrt{\alpha_t}},$$

where  $\alpha_t = \prod_{s=1}^t (1 - \beta_s)$  is the cumulative product of the noise schedule.

However, Lotus (He et al., 2024) shows that  $\epsilon$ -prediction injects unnecessary stochastic variation, which accumulates across multiple denoising steps, degrading dense prediction quality. Lotus therefore advocates a **deterministic** adaptation, directly predicting the clean annotation latent  $x_0$  in a single step:

$$\hat{x}_0 = g_\theta(z_T),$$

where  $T$  is a fixed timestep and  $g_\theta$  is fine-tuned U-Net applied once. This formulation eliminates the prediction variance introduced by multiple-step denoising, simplifies optimization, and significantly accelerates inference. In our approach, we similarly avoid the diffusion chain and adopt this  $x_0$ -prediction design.

## 4 METHODOLOGY

### 4.1 LEVERAGING THE MULTI-SCALE LATENT FEATURES FOR POSE ESTIMATION

We directly leverage the multi-scale latent features of SD U-Net for the pose estimation task. The input image is encoded by the frozen SD-VAE encoder and then fed into the SD U-Net, from which we extract multi-scale features at the upsampling stage. **These multi-scale features provide rich and robust visual representations for downstream keypoint prediction.** The specific feature level used for each task configuration is discussed in Sec. 5.5.2.

### 4.2 THE U-NET CONVOLUTIONAL OUTPUT LAYER FORMS AN INFORMATION BOTTLENECK

Stable Diffusion’s U-Net outputs a 4-channel latent  $z \in \mathbb{R}^{4 \times h \times w}$  for the VAE through a single convolutional layer. In contrast, pose estimation requires  $K$ -channel heatmaps

Table 1: **Quantitative comparison across COCO, HumanArt, COCO-OOD Monet, and COCO-OOD Ukiyo-e.** All models are trained on COCO. Full quantitative comparison across various models are in the supplementary.

Model Variant	Model	Pre-trained Backbone	Params	COCO		HumanArt		COCO-OOD		COCO-OOD	
				AP	AR	AP	AR	Monet	AR	Ukiyo-e	AR
Body	Sapiens-1B (Khirodkar et al., 2024)	Sapiens ViT	1.169B	82.1	85.9	64.3	67.4	58.8	63.3	61.5	66.2
	Sapiens-2B (Khirodkar et al., 2024)	Sapiens ViT	2.163B	<b>82.2</b>	<b>86.0</b>	69.6	72.2	59.6	64.0	62.3	66.8
	GenLoc (Wang et al., 2025a)	Stable Diffusion-v1.5	0.95B	77.6	80.7	67.0	70.8	N/A	N/A	N/A	N/A
	SDPose (Ours)	Stable Diffusion-v2	0.95B	81.3	85.2	<b>71.2</b>	<b>73.9</b>	<b>64.3</b>	<b>68.9</b>	<b>66.0</b>	<b>70.7</b>
Wholebody	Sapiens-1B (Khirodkar et al., 2024)	Sapiens ViT	1.169B	72.7	79.2	N/A	N/A	38.7	46.8	40.5	49.4
	Sapiens-2B (Khirodkar et al., 2024)	Sapiens ViT	2.163B	<b>74.4</b>	<b>81.0</b>	N/A	N/A	44.4	55.5	46.6	55.5
	SDPose (Ours)	Stable Diffusion-v2	0.95B	71.5	78.4	N/A	N/A	<b>46.6</b>	<b>54.8</b>	<b>47.7</b>	<b>56.4</b>

$H \in \mathbb{R}^{K \times H' \times W'}$  with  $K \gg 4$ , making the 4-channel latent a severe information bottleneck. To address this, we replace the original 4-channel head with a lightweight heatmap decoder (Xiao et al., 2018). The decoder consists of deconvolution layers for upsampling, followed by convolutions that output  $K$ -channel heatmaps (Fig. 2). This modification removes the bottleneck and shortens the supervision path to keypoints.

### 4.3 AUXILIARY RGB RECONSTRUCTION

To preserve the fine-detail representation capability of diffusion priors and to avoid overfitting to the pose estimation domain, we adopt the *Detail Preserver* strategy from Lotus (He et al., 2024). Concretely, we introduce a class embedding  $C \in \{C_{\text{RGB}}, C_{\text{Pose}}\}$  that controls the behavior of the denoising U-Net  $f_\theta$ . When  $C_{\text{RGB}}$  is provided, the network is trained to reconstruct the RGB latent  $z_{\text{RGB}}$ ; when  $C_{\text{Pose}}$  is provided, it learns to reconstruct the ground-truth heatmap  $H_{\text{Pose}}$ . The overall objective is

$$L = \|z_{\text{RGB}} - f_{\theta}(z_{\text{input}}, t, C_{\text{RGB}})\|^2 + \|H_{\text{Pose}} - f_{\theta}(z_{\text{input}}, t, C_{\text{Pose}})\|^2,$$

where  $z_{\text{input}}$  is the latent encoded from the input image by the SD-VAE, and  $t$  is fixed to  $t = 1000$  in our experiments.

## 4.4 INFERENCE

As illustrated in Fig. 3, the input RGB image  $x$  is encoded by the SD-VAE into the latent representation  $z_{\text{RGB}}$ . The latent diffusion U-Net then performs a single-step regression with the timestep fixed at  $t = 1000$ , using the class label  $C_{\text{Pose}}$  to execute the pose estimation task. The text condition is disabled by feeding an empty text embedding to the U-Net.

## 5 EXPERIMENTS

## 5.1 EXPERIMENT SETTINGS

**Implementation Details.** We train SDPose based on Stable Diffusion V2 (Rombach et al., 2022), with text conditioning disabled. During training, we fix the timestep  $t = 1000$ . For more details, please see the supplementary materials.

**Training Datasets.** We train two variants, SDPose Body (17-keypoints) and SDPose Wholebody (133-keypoints), on MS COCO (Lin et al., 2014) and COCO-WholeBody (Jin et al., 2020), respectively. All images are processed using standard top-down augmentations, with the input resolution set to  $1024 \times 768$ . Further details are provided in the supplementary materials.

**Validation Datasets and Metrics.** (1) For the Body variant, we evaluate SDPose on MS COCO (Lin et al., 2014) for real-world images, and on HumanArt (Ju et al., 2023) and COCO-OOD for cross-domain benchmarks. (2) For the Wholebody variant, we evaluate SDPose on COCO-WholeBody (Jin et al., 2020) and the extended COCO-OOD. Further details of the evaluation datasets and metrics are provided in the supplementary materials.

**COCO-OOD.** To complement HumanArt and enable OOD evaluation with matched content and labels, we translate all COCO val images into three artistic domains: Monet-style

paintings using the official StyTR2 framework (Deng et al., 2022), ukiyo-e style using the official CycleGAN implementation (Zhu et al., 2017), and color-sketch style using Nano Banana (Nano Banana). For all stylized images, we reuse the original COCO val annotations (bounding boxes and keypoints). Please refer to the supplementary materials for details.

## 5.2 QUANTITATIVE AND QUALITATIVE COMPARISON ON REAL-WORLD SCENES

For the Body variant, SDPose achieves 81.3 AP / 85.2 AR on the COCO validation set (Table 1) with only 40 training epochs using a 0.95B SD-v2 backbone. It matches the accuracy of Sapiens (82.1–82.2 AP) despite requiring 5× fewer epochs and a smaller backbone, and surpasses GenLoc (+3.7 AP, +4.5 AR). Figure 4 further illustrates robustness on real-world photos, where SDPose rivals Sapiens and corrects its failure cases (e.g., Sichuan opera eye keypoints). For the Wholebody variant, SDPose achieves competitive performance with Sapiens-1B on the COCO-WholeBody validation set. Further details are provided in the supplementary materials.



Figure 5: **Illustration of the COCO-OOD dataset.**

## 5.3 SDPOSE’S STRONG OOD ROBUSTNESS

In this section, we demonstrate the superior OOD robustness of SDPose using quantitative evaluation on HumanArt and our COCO-OOD benchmark. As shown in Table 1, SDPose achieves state-of-the-art results on HumanArt and COCO-OOD with fewer training epochs and a smaller parameter budget. On COCO-OOD WholeBody, SDPose continues to demonstrate strong out-of-domain robustness. As shown in Fig. 1, SDPose achieves more accurate body pose estimation across diverse animation styles and humanoid robots compared with baseline models. Additional qualitative results on whole-body pose estimation in stylized paintings are provided in the supplementary materials.

## 5.4 WHY DIFFUSION PRIORS EXHIBIT STRONG OOD ROBUSTNESS?

In this subsection, we investigate why the diffusion-based prior (Stable Diffusion v2) (Rombach et al., 2022) exhibits stronger cross-domain generalization than the human-centric Sapiens ViT backbone (Khirodkar et al., 2024). We compare latent feature distributions from the pretrained Sapiens-1B ViT and Stable Diffusion v2 U-Net, as well as their pose-finetuned counterparts, Sapiens-1B Pose and our SDPose model (finetuned on the last four upsampling blocks). We sample 300 person instances from the COCO and COCO-OOD validation sets. For each instance, we crop the person region using the ground-truth bounding box and collect four stylistic variants: the original COCO image plus its ukiyo-e, Monet-oil, and color-sketch versions from COCO-OOD, yielding 1200 crops in total. These images are passed through the corresponding backbones to extract latent features, on which we run t-SNE (Maaten & Hinton, 2008) and compute silhouette scores (Rousseeuw, 1987).

Fig. 6(a–e) show the t-SNE visualizations of the *pretrained* priors, and Fig. 6(f–j) show the *pose-finetuned* models. In Fig. 6(a), the pretrained Sapiens features form clear style-driven clusters (silhouette by style = 0.3469) with a negative silhouette by person-instance ( $-0.1608$ ), indicating that its representation is dominated by artistic appearance rather than instance structure. In contrast, Stable Diffusion U-Net features in Fig. 6(b–e) exhibit much weaker style separation and gradually stronger person-instance coherence; deeper upsampling blocks yield increasingly person-consistent and style-invariant distributions, reflected by rising person-instance silhouettes (up to  $\approx 0.22$ ) and near-zero style silhouettes. After pose finetuning, both Sapiens-Pose (Fig. 6(f)) and SDPose (Fig. 6(g–j)) produce visible person-instance clusters, but SDPose features are noticeably tighter and more focused, with the mid-level SDPose blocks achieving the highest person silhouettes ( $\approx 0.45$ –0.48) while

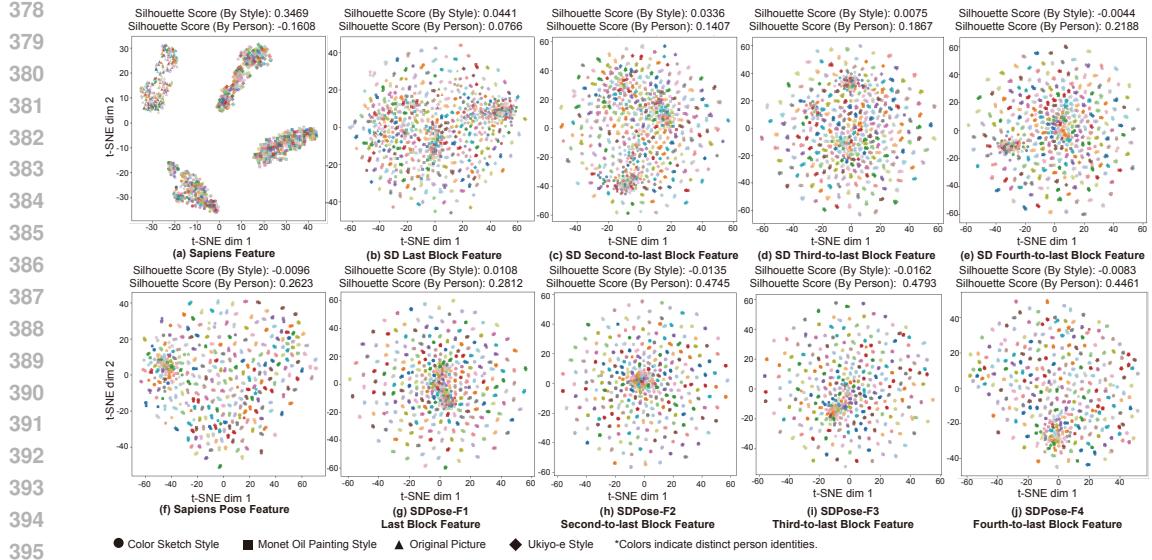


Figure 6: **t-SNE visualization of features from Sapiens ViT, Stable Diffusion U-Net blocks, Sapiens Pose, and SDPose across four visual domains.** Each point corresponds to an image sample; colors represent person instances and marker shapes denote artistic styles.

keeping style silhouettes close to zero. Overall, the SD U-Net (especially after SDPose fine-tuning) provides more instance-coherent and style-disentangled representations than Sapiens, which explains why SD-based features offer stronger cross-domain robustness under artistic style shifts.

## 5.5 ABLATION STUDY

### 5.5.1 ABLATION ON AUXILIARY RECONSTRUCTION TASK AND DIFFUSION PRIORS

We conduct ablations to validate our designs. For the “w/o diffusion priors” variant, we train the U-Net from scratch (no pretrained priors). For the “w/o RGB recon.” variant, we disable only the auxiliary RGB reconstruction branch; all other settings remain identical. From Table 2, two trends emerge. First, removing the RGB branch yields a consistent but modest AP/AR drop on COCO that becomes more pronounced on HumanArt and COCO-OOD, indicating that the auxiliary reconstruction acts as a useful regularizer and improves robustness under domain shift. Second, removing diffusion priors causes a much larger degradation, especially on the OOD benchmarks, highlighting that the pretrained generative priors are the primary source of SDPose’s generalization.

### 5.5.2 ABLATION ON MULTI-SCALE FEATURES FROM THE SD U-NET

Prior work (Liu et al., 2023b; 2024; 2023a; Wang et al., 2025b) suggests that penultimate features often transfer better than final ones. As shown in Table 2, we evaluate SD U-Net upsampling features from the last four blocks. On COCO (17-keypoint), the last-block (F1) and penultimate (F2) features perform very similarly, and F1 is slightly better on both COCO and HumanArt, which contain a mixture of natural and artistic content. In contrast, on COCO-OOD Monet, where the domain shift is purely stylistic, F2 achieves the best performance, suggesting that it captures style-invariant cues more effectively than F1. For the 133-keypoint WholeBody setting, F2 consistently outperforms all other feature levels on both COCO-WholeBody and COCO-OOD, indicating that it offers a better balance between semantic robustness and spatial detail in this more fine-grained regime. Together with the latent-space analysis in Section 5.4, these results show that deeper SD features (F3 and F4) encode even stronger instance-consistent, style-invariant semantics but operate at much

432

433  
434  
435  
Table 2: **Ablation studies on diffusion priors, RGB reconstruction, and U-Net  
feature selection.** All experiments are trained on COCO with 40 epochs (body-17 key-  
points) or 42 epochs (wholebody-133 keypoints).

Model Variant	Ablation Setting	COCO		HumanArt		COCO-OOD Monet	
		AP	AR	AP	AR	AP	AR
SDPose-Body	Full Model (Last Block)	<b>81.3</b>	<b>85.2</b>	<b>71.2</b>	<b>73.9</b>	<b>63.5</b>	<b>68.2</b>
	w/o RGB Reconstruction	80.8 (-0.5)	84.9 (-0.3)	69.8 (-1.4)	72.6 (-1.3)	62.5 (-1.0)	67.3 (-0.9)
	w/o Diffusion Priors	74.9 (-6.4)	79.4 (-5.8)	53.8 (-17.4)	58.0 (-15.9)	52.7 (-10.8)	57.9 (-10.3)
	Last Block (F1)	<b>81.3</b>	<b>85.2</b>	<b>71.2</b>	<b>73.9</b>	63.5	68.2
SDPose-Wholebody	Second-to-last Block (F2)	81.1	85.0	70.4	73.3	<b>64.3</b>	<b>68.9</b>
	Third-to-last Block (F3)	81.0	85.1	70.6	73.3	63.5	68.2
	Fourth-to-last Block (F4)	79.2	83.4	65.0	68.1	58.1	62.8
	Last Block (F1)	70.5	77.5	N/A	N/A	44.7	53.0
	Second-to-last Block (F2)	<b>71.5</b>	<b>78.4</b>	N/A	N/A	<b>46.6</b>	<b>54.8</b>
	Third-to-last Block (F3)	70.4	77.6	N/A	N/A	45.5	53.8
	Fourth-to-last Block (F4)	64.6	72.1	N/A	N/A	37.1	45.4

446

447

448  
449  
450  
451  
452  
coarser spatial resolutions (H/16 and H/32), which removes the geometric details needed for  
453  
454  
455  
456  
457  
458  
459  
460  
461  
462  
463  
464  
465  
466  
467  
468  
469  
470  
471  
472  
473  
474  
475  
476  
477  
478  
479  
480  
481  
482  
483  
484  
485  
486  
487  
488  
489  
490  
491  
492  
493  
494  
495  
496  
497  
498  
499  
500  
501  
502  
503  
504  
505  
506  
507  
508  
509  
510  
511  
512  
513  
514  
515  
516  
517  
518  
519  
520  
521  
522  
523  
524  
525  
526  
527  
528  
529  
530  
531  
532  
533  
534  
535  
536  
537  
538  
539  
540  
541  
542  
543  
544  
545  
546  
547  
548  
549  
550  
551  
552  
553  
554  
555  
556  
557  
558  
559  
560  
561  
562  
563  
564  
565  
566  
567  
568  
569  
570  
571  
572  
573  
574  
575  
576  
577  
578  
579  
580  
581  
582  
583  
584  
585  
586  
587  
588  
589  
590  
591  
592  
593  
594  
595  
596  
597  
598  
599  
600  
601  
602  
603  
604  
605  
606  
607  
608  
609  
610  
611  
612  
613  
614  
615  
616  
617  
618  
619  
620  
621  
622  
623  
624  
625  
626  
627  
628  
629  
630  
631  
632  
633  
634  
635  
636  
637  
638  
639  
640  
641  
642  
643  
644  
645  
646  
647  
648  
649  
650  
651  
652  
653  
654  
655  
656  
657  
658  
659  
660  
661  
662  
663  
664  
665  
666  
667  
668  
669  
670  
671  
672  
673  
674  
675  
676  
677  
678  
679  
680  
681  
682  
683  
684  
685  
686  
687  
688  
689  
690  
691  
692  
693  
694  
695  
696  
697  
698  
699  
700  
701  
702  
703  
704  
705  
706  
707  
708  
709  
710  
711  
712  
713  
714  
715  
716  
717  
718  
719  
720  
721  
722  
723  
724  
725  
726  
727  
728  
729  
730  
731  
732  
733  
734  
735  
736  
737  
738  
739  
740  
741  
742  
743  
744  
745  
746  
747  
748  
749  
750  
751  
752  
753  
754  
755  
756  
757  
758  
759  
759  
760  
761  
762  
763  
764  
765  
766  
767  
768  
769  
770  
771  
772  
773  
774  
775  
776  
777  
778  
779  
779  
780  
781  
782  
783  
784  
785  
786  
787  
788  
789  
789  
790  
791  
792  
793  
794  
795  
796  
797  
798  
799  
800  
801  
802  
803  
804  
805  
806  
807  
808  
809  
809  
810  
811  
812  
813  
814  
815  
816  
817  
818  
819  
819  
820  
821  
822  
823  
824  
825  
826  
827  
828  
829  
829  
830  
831  
832  
833  
834  
835  
836  
837  
838  
839  
839  
840  
841  
842  
843  
844  
845  
846  
847  
848  
849  
850  
851  
852  
853  
854  
855  
856  
857  
858  
859  
859  
860  
861  
862  
863  
864  
865  
866  
867  
868  
869  
869  
870  
871  
872  
873  
874  
875  
876  
877  
878  
879  
879  
880  
881  
882  
883  
884  
885  
886  
887  
888  
889  
889  
890  
891  
892  
893  
894  
895  
896  
897  
898  
899  
900  
901  
902  
903  
904  
905  
906  
907  
908  
909  
909  
910  
911  
912  
913  
914  
915  
916  
917  
918  
919  
919  
920  
921  
922  
923  
924  
925  
926  
927  
928  
929  
929  
930  
931  
932  
933  
934  
935  
936  
937  
938  
939  
939  
940  
941  
942  
943  
944  
945  
946  
947  
948  
949  
950  
951  
952  
953  
954  
955  
956  
957  
958  
959  
960  
961  
962  
963  
964  
965  
966  
967  
968  
969  
970  
971  
972  
973  
974  
975  
976  
977  
978  
979  
980  
981  
982  
983  
984  
985  
986  
987  
988  
989  
989  
990  
991  
992  
993  
994  
995  
996  
997  
998  
999  
1000  
1001  
1002  
1003  
1004  
1005  
1006  
1007  
1008  
1009  
1009  
1010  
1011  
1012  
1013  
1014  
1015  
1016  
1017  
1018  
1019  
1019  
1020  
1021  
1022  
1023  
1024  
1025  
1026  
1027  
1028  
1029  
1029  
1030  
1031  
1032  
1033  
1034  
1035  
1036  
1037  
1038  
1039  
1039  
1040  
1041  
1042  
1043  
1044  
1045  
1046  
1047  
1048  
1049  
1049  
1050  
1051  
1052  
1053  
1054  
1055  
1056  
1057  
1058  
1059  
1059  
1060  
1061  
1062  
1063  
1064  
1065  
1066  
1067  
1068  
1069  
1069  
1070  
1071  
1072  
1073  
1074  
1075  
1076  
1077  
1078  
1079  
1079  
1080  
1081  
1082  
1083  
1084  
1085  
1086  
1087  
1088  
1089  
1089  
1090  
1091  
1092  
1093  
1094  
1095  
1096  
1097  
1098  
1099  
1099  
1100  
1101  
1102  
1103  
1104  
1105  
1106  
1107  
1108  
1109  
1109  
1110  
1111  
1112  
1113  
1114  
1115  
1116  
1117  
1118  
1119  
1119  
1120  
1121  
1122  
1123  
1124  
1125  
1126  
1127  
1128  
1129  
1129  
1130  
1131  
1132  
1133  
1134  
1135  
1136  
1137  
1138  
1139  
1139  
1140  
1141  
1142  
1143  
1144  
1145  
1146  
1147  
1148  
1149  
1149  
1150  
1151  
1152  
1153  
1154  
1155  
1156  
1157  
1158  
1159  
1159  
1160  
1161  
1162  
1163  
1164  
1165  
1166  
1167  
1168  
1169  
1169  
1170  
1171  
1172  
1173  
1174  
1175  
1176  
1177  
1178  
1179  
1179  
1180  
1181  
1182  
1183  
1184  
1185  
1186  
1187  
1188  
1189  
1189  
1190  
1191  
1192  
1193  
1194  
1195  
1196  
1197  
1198  
1198  
1199  
1199  
1200  
1201  
1202  
1203  
1204  
1205  
1206  
1207  
1208  
1209  
1209  
1210  
1211  
1212  
1213  
1214  
1215  
1216  
1217  
1218  
1219  
1219  
1220  
1221  
1222  
1223  
1224  
1225  
1226  
1227  
1228  
1229  
1229  
1230  
1231  
1232  
1233  
1234  
1235  
1236  
1237  
1238  
1239  
1239  
1240  
1241  
1242  
1243  
1244  
1245  
1246  
1247  
1248  
1249  
1249  
1250  
1251  
1252  
1253  
1254  
1255  
1256  
1257  
1258  
1259  
1259  
1260  
1261  
1262  
1263  
1264  
1265  
1266  
1267  
1268  
1269  
1269  
1270  
1271  
1272  
1273  
1274  
1275  
1276  
1277  
1278  
1279  
1279  
1280  
1281  
1282  
1283  
1284  
1285  
1286  
1287  
1288  
1289  
1289  
1290  
1291  
1292  
1293  
1294  
1295  
1296  
1297  
1298  
1298  
1299  
1299  
1300  
1301  
1302  
1303  
1304  
1305  
1306  
1307  
1308  
1309  
1309  
1310  
1311  
1312  
1313  
1314  
1315  
1316  
1317  
1318  
1319  
1319  
1320  
1321  
1322  
1323  
1324  
1325  
1326  
1327  
1328  
1329  
1329  
1330  
1331  
1332  
1333  
1334  
1335  
1336  
1337  
1338  
1339  
1339  
1340  
1341  
1342  
1343  
1344  
1345  
1346  
1347  
1348  
1349  
1349  
1350  
1351  
1352  
1353  
1354  
1355  
1356  
1357  
1358  
1359  
1359  
1360  
1361  
1362  
1363  
1364  
1365  
1366  
1367  
1368  
1369  
1369  
1370  
1371  
1372  
1373  
1374  
1375  
1376  
1377  
1378  
1379  
1379  
1380  
1381  
1382  
1383  
1384  
1385  
1386  
1387  
1388  
1389  
1389  
1390  
1391  
1392  
1393  
1394  
1395  
1396  
1397  
1398  
1398  
1399  
1399  
1400  
1401  
1402  
1403  
1404  
1405  
1406  
1407  
1408  
1409  
1409  
1410  
1411  
1412  
1413  
1414  
1415  
1416  
1417  
1418  
1419  
1419  
1420  
1421  
1422  
1423  
1424  
1425  
1426  
1427  
1428  
1429  
1429  
1430  
1431  
1432  
1433  
1434  
1435  
1436  
1437  
1438  
1439  
1440  
1441  
1442  
1443  
1444  
1445  
1446  
1447  
1448  
1449  
1450  
1451  
1452  
1453  
1454  
1455  
1456  
1457  
1458  
1459  
1460  
1461  
1462  
1463  
1464  
1465  
1466  
1467  
1468  
1469  
1470  
1471  
1472  
1473  
1474  
1475  
1476  
1477  
1478  
1479  
1480  
1481  
1482  
1483  
1484  
1485  
1486  
1487  
1488  
1489  
1490  
1491  
1492  
1493  
1494  
1495  
1496  
1497  
1498  
1499  
1500  
1501  
1502  
1503  
1504  
1505  
1506  
1507  
1508  
1509  
1509  
1510  
1511  
1512  
1513  
1514  
1515  
1516  
1517  
1518  
1519  
1519  
1520  
1521  
1522  
1523  
1524  
1525  
1526  
1527  
1528  
1529  
1529  
1530  
1531  
1532  
1533  
1534  
1535  
1536  
1537  
1538  
1539  
1539  
1540  
1541  
1542  
1543  
1544  
1545  
1546  
1547  
1548  
1549  
1549  
1550  
1551  
1552  
1553  
1554  
1555  
1556  
1557  
1558  
1559  
1559  
1560  
1561  
1562  
1563  
1564  
1565  
1566  
1567  
1568  
1569  
1569  
1570  
1571  
1572  
1573  
1574  
1575  
1576  
1577  
1578  
1579  
1579  
1580  
1581  
1582  
1583  
1584  
1585  
1586  
1587  
1588  
1589  
1589  
1590  
1591  
1592  
1593  
1594  
1595  
1596  
1597  
1598  
1598  
1599  
1599  
1600  
1601  
1602  
1603  
1604  
1605  
1606  
1607  
1608  
1609  
1609  
1610  
1611  
1612  
1613  
1614  
1615  
1616  
1617  
1618  
1619  
1619  
1620  
1621  
1622  
1623  
1624  
1625  
1626  
1627  
1628  
1629  
1629  
1630  
1631  
1632  
1633  
1634  
1635  
1636  
1637  
1638  
1639  
1639  
1640  
1641  
1642  
1643  
1644  
1645  
1646  
1647  
1648  
1649  
1649  
1650  
1651  
1652  
1653  
1654  
1655  
1656  
1657  
1658  
1659  
1659  
1660  
1661  
1662  
1663  
1664  
1665  
1666  
1667  
1668  
1669  
1669  
1670  
1671  
1672  
1673  
1674  
1675  
1676  
1677  
1678  
1679  
1679  
1680  
1681  
1682  
1683  
1684  
1685  
1686  
1687  
1688  
1689  
1689  
1690  
1691  
1692  
1693  
1694  
1695  
1696  
1697  
1698  
1698  
1699  
1699  
1700  
1701  
1702  
1703  
1704  
1705  
1706  
1707  
1708  
1709  
1709  
1710  
1711  
1712  
1713  
1714  
1715  
1716  
1717  
1718  
1719  
1719  
1720  
1721  
1722  
1723  
1724  
1725  
1726  
1727  
1728  
1729  
1729  
1730  
1731  
1732  
1733  
1734  
1735  
1736  
1737  
1738  
1739  
1739  
1740  
1741  
1742  
1743  
1744  
1745  
1746  
1747  
1748  
1749  
1749  
1750  
1751  
1752  
1753  
1754  
1755  
1756  
1757  
1758  
1759  
1759  
1760  
1761  
1762  
1763  
1764  
1765  
1766  
1767  
1768  
1769  
1769  
1770  
1771  
1772  
1773  
1774  
1775  
1776  
1777  
1778  
1779  
1779  
1780  
1781  
1782  
1783  
1784  
1785  
1786  
1787  
1788  
1789  
1789  
1790  
1791  
1792  
1793  
1794  
1795  
1796  
1797  
1798  
1798  
1799  
1799  
1800  
1801  
1802  
1803  
1804  
1805  
1806  
1807  
1808  
1809  
1809  
1810  
1811  
1812  
1813  
1814  
1815  
1816  
1817  
1818  
1819  
1819  
1820  
1821  
1822  
1823  
1824  
1825  
1826  
1827  
1828  
1829  
1829  
1830  
1831  
1832  
1833  
1834  
1835  
1836  
1837  
1838  
1839  
1839  
1840  
1841  
1842  
1843  
1844  
1845  
1846  
1847  
1848  
1849  
1849  
1850  
1851  
1852  
1853  
1854  
1855  
1856  
1857  
1858  
1859  
1859  
1860  
1861  
1862  
1863  
1864  
1865  
1866  
1867  
1868  
1869  
1869  
1870  
1871  
1872  
1873  
1874  
1875  
1876  
1877  
1878  
1879  
1879  
1880  
1881  
1882  
1883  
1884  
1885  
1886  
1887  
1888  
1889  
1889  
1890  
1891  
1892  
1893  
1894  
1895  
1896  
1897  
1898  
1898  
1899  
1899  
1900  
1901  
1902  
1903  
1904  
1905  
1906  
1907  
1908  
1909  
1909  
1910  
1911  
1912  
1913  
1914  
1915  
1916  
1917  
1918  
1919  
1919  
1920  
1921  
1922  
1923  
1924  
1925  
1926  
1927  
1928  
1929  
1929  
1930  
1931  
1932  
1933  
1934  
1935  
1936  
1937  
1938  
1939  
1939  
1940  
1941  
1942  
1943  
1944  
1945  
1946  
1947  
1948  
1949  
1949  
1950  
1951  
1952  
1953  
1954  
1955  
1956  
1957  
1958  
1959  
1959  
1960  
1961  
1962  
1963  
1964  
1965  
1966  
1967  
1968  
1969  
1969  
1970  
1971  
1972  
1973  
1974  
1975  
1976  
1977  
1978  
1979  
1979  
1980  
1981  
1982  
1983  
1984  
1985  
1986  
1987  
1988  
1989  
1989  
1990  
1991  
1992  
1993  
1994  
1995  
1996  
1997  
1998  
1998  
1999  
1999  
2000  
2001  
2002  
2003  
2004  
2005  
2006  
2007  
2008  
2009  
2009  
2010  
2011  
2012  
2013  
2014  
2015  
2016  
2017  
2018  
2019  
2019  
2020  
2021  
2022  
2023  
2024  
2025  
2026<br

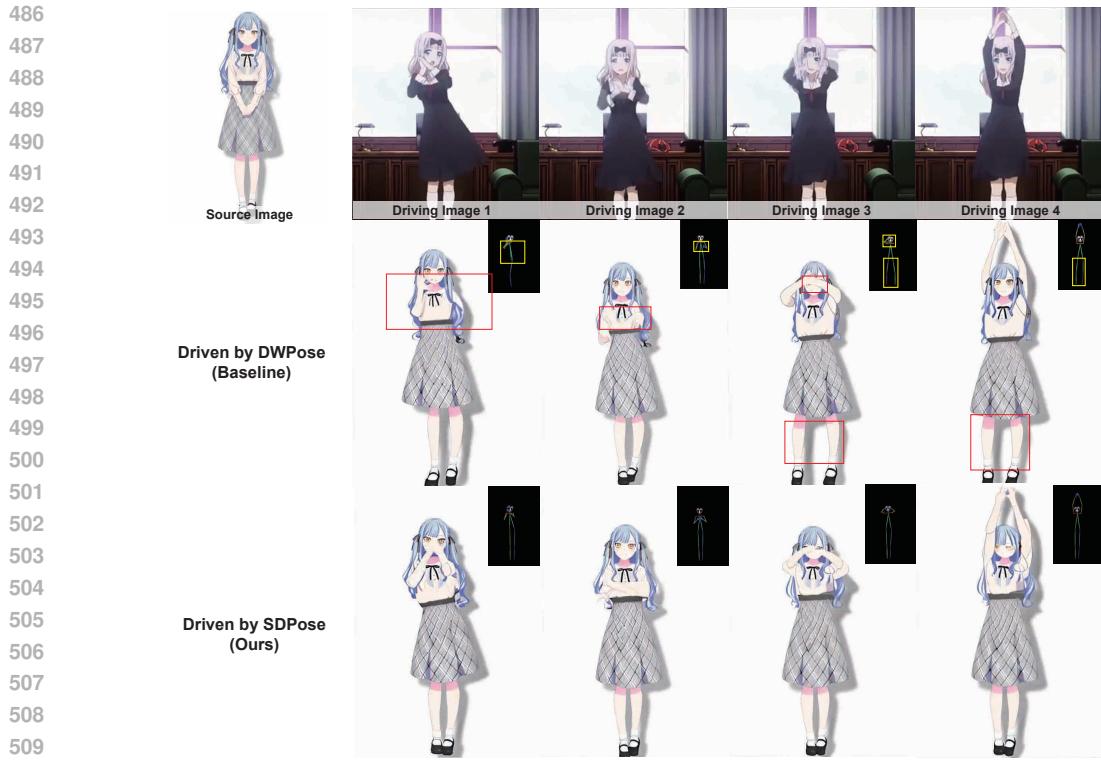


Figure 8: **Qualitative comparison for pose-controlled video generation in the wild.** The first row shows the source image and frames from the driving video. The second row shows output video frames generated from the pose sequence estimated by the baseline model DWPose, while the third row shows the results guided by our SDPose. Red boxes highlight failures in the generated video, and yellow boxes highlight errors in pose estimation.

## 6.2 POSE-GUIDED VIDEO GENERATION

Recent advances in controlled video generation have gained significant traction (Hu, 2024; Guo et al., 2023; Kim et al., 2024). Despite the progress of video generation models in producing higher-quality outputs, extracting reliable control conditions remains critical for achieving high-quality results. As shown in Fig. 8, our SDPose provides more accurate poses for the driving frames, enabling more reliable pose-sequence transfer from animations to animations. Video frames are generated by Moore-Animated Anyone<sup>1</sup>.

## 7 CONCLUSION

In this paper, we present **SDPose**, an SD-native fine-tuning framework for human pose estimation. SDPose preserves the original U-Net with only lightweight task-specific components and adapts diffusion latent priors for keypoint prediction through an RGB reconstruction branch and a heatmap decoder. We further introduce **COCO-OOD**, a style-transferred extension of COCO for evaluating robustness under domain shifts. With only one-fifth of the fine-tuning cost of Sapiens and a smaller backbone, SDPose matches its in-domain accuracy on COCO and achieves state-of-the-art results on COCO-OOD and HumanArt. Our ablations and latent-space analyses show that diffusion priors and multi-scale SD features naturally encode domain-invariant structure, enabling strong generalization in both pose estimation and zero-shot pose annotation for controllable generation.

<sup>1</sup><https://github.com/MooreThreads/Moore-AnimateAnyone>

540  
541 ETHICS STATEMENT

542 This work builds upon publicly available datasets (COCO, COCO Wholebody), all of which  
 543 have established licenses and annotation protocols. No private or personally identifiable  
 544 information is used. Our method focuses on improving the robustness of pose estimation  
 545 under domain shifts, which can benefit applications such as animation and embodied AI.

546  
547 REPRODUCIBILITY STATEMENT

548 We have made every effort to ensure reproducibility. All datasets used are publicly available,  
 549 and we detail dataset splits, preprocessing steps, and evaluation protocols in Sec. 4. Our  
 550 training settings, hyperparameters, and model architectures are fully described in Sec. 4,  
 551 Sec. 5 and Appendix A. Code and scripts to reproduce our experiments will be released  
 552 upon publication.

553  
554 REFERENCES

555 Xiangqi Chen, Chengzhan Yang, Jiashuaizi Mo, Yaxin Sun, Hicham Karmouni, Yunliang  
 556 Jiang, and Zhonglong Zheng. Cspnext: A new efficient token hybrid backbone. *Engineering  
 557 Applications of Artificial Intelligence*, 132:107886, 2024.

558 Yen-Chi Cheng, Hsin-Ying Lee, Sergey Tulyakov, Alexander G Schwing, and Liang-Yan  
 559 Gui. Sdfusion: Multimodal 3d shape completion, reconstruction, and generation. In  
 560 *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp.  
 561 4456–4465, 2023.

562 Yingying Deng, Fan Tang, Weiming Dong, Chongyang Ma, Xingjia Pan, Lei Wang, and  
 563 Changsheng Xu. Stytr2: Image style transfer with transformers. In *IEEE Conference on  
 564 Computer Vision and Pattern Recognition (CVPR)*, 2022.

565 Patrick Esser, Sumith Kulal, Andreas Blattmann, Rahim Entezari, Jonas Müller, Harry  
 566 Saini, Yam Levi, Dominik Lorenz, Axel Sauer, Frederic Boesel, et al. Scaling rectified flow  
 567 transformers for high-resolution image synthesis. In *Forty-first international conference  
 568 on machine learning*, 2024.

569 Yuwei Guo, Ceyuan Yang, Anyi Rao, Zhengyang Liang, Yaohui Wang, Yu Qiao, Maneesh  
 570 Agrawala, Dahua Lin, and Bo Dai. Animatediff: Animate your personalized text-to-image  
 571 diffusion models without specific tuning. *arXiv preprint arXiv:2307.04725*, 2023.

572 Jing He, Haodong Li, Wei Yin, Yixun Liang, Leheng Li, Kaiqiang Zhou, Hongbo Liu,  
 573 Bingbing Liu, and Ying-Cong Chen. Lotus: Diffusion-based visual foundation model for  
 574 high-quality dense prediction. *arXiv preprint arXiv:2409.18124*, 2024.

575 Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Ad-  
 576 vances in neural information processing systems*, 33:6840–6851, 2020.

577 Li Hu. Animate anyone: Consistent and controllable image-to-video synthesis for character  
 578 animation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern  
 579 Recognition*, pp. 8153–8163, 2024.

580 Junjie Huang, Zheng Zhu, Feng Guo, and Guan Huang. The devil is in the details: Delving  
 581 into unbiased data processing for human pose estimation. In *Proceedings of the IEEE/CVF  
 582 conference on computer vision and pattern recognition*, pp. 5700–5709, 2020.

583 Tao Jiang, Peng Lu, Li Zhang, Ningsheng Ma, Rui Han, Chengqi Lyu, Yining Li, and  
 584 Kai Chen. Rtmpose: Real-time multi-person pose estimation based on mmpose. *arXiv  
 585 preprint arXiv:2303.07399*, 2023.

586 Sheng Jin, Lumin Xu, Jin Xu, Can Wang, Wentao Liu, Chen Qian, Wanli Ouyang, and Ping  
 587 Luo. Whole-body human pose estimation in the wild. In *Proceedings of the European  
 588 Conference on Computer Vision (ECCV)*, 2020.

594 Xuan Ju, Ailing Zeng, Jianan Wang, Qiang Xu, and Lei Zhang. Human-art: A versa-  
 595 tile human-centric dataset bridging natural and artificial scenes. In *Proceedings of the*  
 596 *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023.

597

598 Markus Karmann and Onay Urfalioglu. Repurposing stable diffusion attention for training-  
 599 free unsupervised interactive segmentation. In *Proceedings of the Computer Vision and*  
 600 *Pattern Recognition Conference*, pp. 24518–24528, 2025.

601 Bingxin Ke, Anton Obukhov, Shengyu Huang, Nando Metzger, Rodrigo Caye Daudt, and  
 602 Konrad Schindler. Repurposing diffusion-based image generators for monocular depth  
 603 estimation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern*  
 604 *recognition*, pp. 9492–9502, 2024.

605

606 Rawal Khirodkar, Timur Bagautdinov, Julieta Martinez, Su Zhaoen, Austin James, Peter  
 607 Selednik, Stuart Anderson, and Shunsuke Saito. Sapiens: Foundation for human vision  
 608 models. In *European Conference on Computer Vision*, pp. 206–228. Springer, 2024.

609 Jeongho Kim, Min-Jung Kim, Junsoo Lee, and Jaegul Choo. Tcan: Animating human  
 610 images with temporally consistent pose guidance using diffusion models. In *European*  
 611 *Conference on Computer Vision*, pp. 326–342. Springer, 2024.

612 Jiantao Lin, Xin Yang, Meixi Chen, Yingjie Xu, Dongyu Yan, Leyi Wu, Xinli Xu, Lie Xu,  
 613 Shunsi Zhang, and Ying-Cong Chen. Kiss3dgen: Repurposing image diffusion models  
 614 for 3d asset generation. In *Proceedings of the Computer Vision and Pattern Recognition*  
 615 *Conference*, pp. 5870–5880, 2025.

616

617 Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan,  
 618 Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In  
 619 *European conference on computer vision*, pp. 740–755. Springer, 2014.

620 Yaron Lipman, Ricky TQ Chen, Heli Ben-Hamu, Maximilian Nickel, and Matt Le. Flow  
 621 matching for generative modeling. *arXiv preprint arXiv:2210.02747*, 2022.

622

623 Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. Improved baselines with visual  
 624 instruction tuning, 2023a.

625 Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning,  
 626 2023b.

627

628 Haotian Liu, Chunyuan Li, Yuheng Li, Bo Li, Yuanhan Zhang, Sheng Shen, and Yong Jae  
 629 Lee. Llava-next: Improved reasoning, ocr, and world knowledge, January 2024. URL  
 630 <https://llava-vl.github.io/blog/2024-01-30-llava-next/>.

631 Xiaoxiao Long, Yuan-Chen Guo, Cheng Lin, Yuan Liu, Zhiyang Dou, Lingjie Liu, Yuexin  
 632 Ma, Song-Hai Zhang, Marc Habermann, Christian Theobalt, et al. Wonder3d: Single  
 633 image to 3d using cross-domain diffusion. In *Proceedings of the IEEE/CVF conference on*  
 634 *computer vision and pattern recognition*, pp. 9970–9980, 2024.

635

636 Cheng Lu, Yuhao Zhou, Fan Bao, Jianfei Chen, Chongxuan Li, and Jun Zhu. Dpm-solver:  
 637 A fast ode solver for diffusion probabilistic model sampling in around 10 steps. *Advances*  
 638 *in neural information processing systems*, 35:5775–5787, 2022.

639 Laurens van der Maaten and Geoffrey Hinton. Visualizing data using t-sne. *Journal of*  
 640 *machine learning research*, 9(Nov):2579–2605, 2008.

641

642 Ginés Hidalgo Martinez. *Openpose: Whole-body pose estimation*. PhD thesis, Carnegie  
 643 Mellon University Pittsburgh, PA, USA, 2019.

644 Nano Banana. Nano banana: Ai image editor and style-transfer tool. <https://nanobanana.ai>. Accessed: 2025-11-12.

645

646 William Peebles and Saining Xie. Scalable diffusion models with transformers. In *Proceedings*  
 647 *of the IEEE/CVF international conference on computer vision*, pp. 4195–4205, 2023.

648 Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Om-  
 649 mer. High-resolution image synthesis with latent diffusion models. In *Proceedings of*  
 650 *the IEEE/CVF conference on computer vision and pattern recognition*, pp. 10684–10695,  
 651 2022.

652 Peter J Rousseeuw. Silhouettes: a graphical aid to the interpretation and validation of  
 653 cluster analysis. *Journal of computational and applied mathematics*, 20:53–65, 1987.

654 Christoph Schuhmann, Romain Beaumont, Richard Vencu, Cade Gordon, Ross Wightman,  
 655 Mehdi Cherti, Theo Coombes, Aarush Katta, Clayton Mullis, Mitchell Wortsman, et al.  
 656 Laion-5b: An open large-scale dataset for training next generation image-text models.  
 657 *Advances in neural information processing systems*, 35:25278–25294, 2022.

658 Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models.  
 659 *arXiv preprint arXiv:2010.02502*, 2020a.

660 Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon,  
 661 and Ben Poole. Score-based generative modeling through stochastic differential equations.  
 662 *arXiv preprint arXiv:2011.13456*, 2020b.

663 Ke Sun, Bin Xiao, Dong Liu, and Jingdong Wang. Deep high-resolution representation  
 664 learning for human pose estimation. In *Proceedings of the IEEE/CVF conference on*  
 665 *computer vision and pattern recognition*, pp. 5693–5703, 2019.

666 Ke Sun, Zigang Geng, Depu Meng, Bin Xiao, Dong Liu, Zhaoxiang Zhang, and Jingdong  
 667 Wang. Bottom-up human pose estimation by ranking heatmap-guided adaptive keypoint  
 668 estimates. *arXiv preprint arXiv:2006.15480*, 2020.

669 Dongkai Wang, Jiang Duan, Liangjian Wen, Shiyu Xuan, Hao Chen, and Shiliang Zhang.  
 670 Generalizable object keypoint localization from generative priors. In *Proceedings of the*  
 671 *Computer Vision and Pattern Recognition Conference*, pp. 20265–20274, 2025a.

672 Zeqing Wang, Qingyang Ma, Wentao Wan, Haojie Li, Keze Wang, and Yonghong Tian.  
 673 Is this generated person existed in real-world? fine-grained detecting and calibrating  
 674 abnormal human-body. In *Proceedings of the Computer Vision and Pattern Recognition*  
 675 *Conference*, pp. 21226–21237, 2025b.

676 Bin Xiao, Haiping Wu, and Yichen Wei. Simple baselines for human pose estimation and  
 677 tracking. In *Proceedings of the European conference on computer vision (ECCV)*, pp.  
 678 466–481, 2018.

679 Enze Xie, Junsong Chen, Junyu Chen, Han Cai, Haotian Tang, Yujun Lin, Zhekai Zhang,  
 680 Muyang Li, Ligeng Zhu, Yao Lu, et al. Sana: Efficient high-resolution image synthesis  
 681 with linear diffusion transformers. *arXiv preprint arXiv:2410.10629*, 2024.

682 Guangkai Xu, Yongtao Ge, Mingyu Liu, Chengxiang Fan, Kangyang Xie, Zhiyue Zhao, Hao  
 683 Chen, and Chunhua Shen. What matters when repurposing diffusion models for general  
 684 dense perception tasks? *arXiv preprint arXiv:2403.06090*, 2024.

685 Yufei Xu, Jing Zhang, Qiming Zhang, and Dacheng Tao. Vitpose: Simple vision transformer  
 686 baselines for human pose estimation. *Advances in neural information processing systems*,  
 687 35:38571–38584, 2022.

688 Yufei Xu, Jing Zhang, Qiming Zhang, and Dacheng Tao. Vitpose++: Vision transformer  
 689 for generic body pose estimation. *IEEE Transactions on Pattern Analysis and Machine*  
 690 *Intelligence*, 46(2):1212–1230, 2023.

691 Zhendong Yang, Ailing Zeng, Chun Yuan, and Yu Li. Effective whole-body pose estimation  
 692 with two-stages distillation. In *Proceedings of the IEEE/CVF International Conference*  
 693 *on Computer Vision*, pp. 4210–4220, 2023.

694 Yuhui Yuan, Rao Fu, Lang Huang, Weihong Lin, Chao Zhang, Xilin Chen, and Jingdong  
 695 Wang. Hrformer: High-resolution transformer for dense prediction. *NeurIPS*, 2021.

702 Lvmi Zhang, Anyi Rao, and Maneesh Agrawala. Adding conditional control to text-to-  
703 image diffusion models. In *Proceedings of the IEEE/CVF International Conference on*  
704 *Computer Vision*, pp. 3836–3847, 2023.

705

706 Zhengbo Zhang, Li Xu, Duo Peng, Hossein Rahmani, and Jun Liu. Diff-tracker: text-to-  
707 image diffusion models are unsupervised trackers. In *European Conference on Computer*  
708 *Vision*, pp. 319–337. Springer, 2024.

709

710 Jun-Yan Zhu, Taesung Park, Phillip Isola, and Alexei A Efros. Unpaired image-to-image  
711 translation using cycle-consistent adversarial networkss. In *Computer Vision (ICCV),*  
712 *2017 IEEE International Conference on*, 2017.

713

714

715

716

717

718

719

720

721

722

723

724

725

726

727

728

729

730

731

732

733

734

735

736

737

738

739

740

741

742

743

744

745

746

747

748

749

750

751

752

753

754

755

756 SUPPLEMENTARY MATERIALS FOR **SDPose**: EXPLOITING DIFFUSION  
 757 PRIORS FOR OUT-OF-DOMAIN AND ROBUST POSE ESTIMATION  
 758

759 **A EXPERIMENT SETTINGS**

760 **A.1 IMPLEMENTATION DETAILS**

763 We train SDPose on the COCO-2017 person keypoints *train2017* split only (no extra data),  
 764 with text prompts disabled. The diffusion timestep is fixed at  $t = 1000$ . We use AdamW  
 765 with a learning rate of  $3 \times 10^{-5}$ . All experiments are run on 8 NVIDIA A100-NVLink  
 766 GPUs with a total batch size of 128, without gradient accumulation. Inputs are resized to  
 767  $1024 \times 768$  with standard top-down augmentations. The 17-keypoint model is trained for 40  
 768 epochs (approximately 3 days), and the 133-keypoint model for 42 epochs (approximately  
 769 3 days and a half).

770 **A.2 TRAINING DATASETS**

773 **COCO 2017 Keypoint Detection** We train the 17-keypoint variant on the COCO-2017  
 774 person keypoint detection dataset (Lin et al., 2014). The full COCO release contains more  
 775 than 200,000 images and about 250,000 person instances. Person keypoint annotations  
 776 follow the 17-point format (nose, eyes, ears, shoulders, elbows, wrists, hips, knees, and  
 777 ankles).

778 **COCO Wholebody** To further evaluate large-scale whole-body keypoint estimation, we  
 779 adopt COCO-WholeBody (Jin et al., 2020), an extended benchmark built on top of COCO  
 780 images. COCO-WholeBody augments the original 17 body joints with fine-grained annota-  
 781 tions of foot (6 keypoints), face (68 keypoints), and hands (42 keypoints for hands), resulting  
 782 in a total of 133 keypoints per person. The dataset provides consistent whole-body annota-  
 783 tions across the same training and validation splits as COCO-2017, enabling both fair  
 784 comparison with standard pose estimation methods and comprehensive evaluation under  
 785 the whole-body setting.

786 **A.3 AUGMENTATION DETAILS**

788 The training pipeline first loads the input image and computes the bounding box center and  
 789 scale. It applies random horizontal flipping, half-body augmentation, and random bounding  
 790 box transformations. The image is then affine-transformed to the target input resolution  
 791 using UDP (Huang et al., 2020). Albumentations-based augmentations are then applied,  
 792 including Gaussian blur ( $p = 0.1$ ), median blur ( $p = 0.1$ ), and coarse dropout ( $p = 1.0$ , with  
 793 up to one hole of size 20%–40% of the image).

794 **A.4 EVALUATION DATASETS AND METRICS**

796 **Evaluation Datasets**

798 **COCO 2017 Keypoint Detection.** For in-domain evaluation, we use the  
 799 COCO-2017 validation set (Lin et al., 2014) annotated with 17 body keypoints,  
 800 bounding boxes, and visibility flags. Following the standard top-down evalua-  
 801 tion protocol, we generate person crops from the COCO-released detection results  
 802 (`COCO_val2017_detections_AP_H_70_person.json`), and report COCO keypoint AP/AR  
 803 on this diverse in-the-wild dataset.

804 **HumanArt.** We use HumanArt (Ju et al., 2023) as an cross domain benchmark: 50k  
 805 human-centric images across 20 scenarios (5 natural, 15 artistic—oil painting, sculpture,  
 806 cartoon, sketch, stained glass, Ukiyo-e, watercolor, etc.) with annotations for boxes and 2D  
 807 keypoints. We follow the official protocol and report keypoint AP/AR to assess robustness  
 808 under artistic domain shift.

809 **COCO-WholeBody (133-keypoint whole-body).** We train and evaluate a  
 133-keypoint variant on COCO-WholeBody (Jin et al., 2020), which shares COCO’s

train/val split. Each person has 133 keypoints (17 body, 6 foot, 68 face, 42 hand) plus boxes for person/face/left/right hand. We follow the official protocol and report Whole-Body AP and part-wise AP (body/foot/face/hand). This dataset spans diverse in-the-wild scenes and stresses fine-grained articulation, complementing COCO for structured keypoint evaluation.

**Metrics** We follow the standard COCO keypoint evaluation protocol, which is based on the Object Keypoint Similarity (OKS). For each keypoint  $i$ , the similarity is defined as

$$KS_i = \exp\left(-\frac{d_i^2}{2s^2k_i^2}\right),$$

where  $d_i$  denotes the Euclidean distance between predicted and ground-truth keypoints,  $s$  is the object scale (square root of the segmentation area), and  $k_i$  is a per-keypoint constant controlling falloff. The OKS for an instance is the average  $KS_i$  over visible keypoints:

$$OKS = \frac{\sum_i KS_i \cdot \delta(v_i > 0)}{\sum_i \delta(v_i > 0)},$$

where  $v_i$  is the visibility flag. Using OKS as the matching criterion, COCO computes Average Precision (AP) as the mean precision over OKS thresholds  $[0.50 : 0.05 : 0.95]$ , and Average Recall (AR) analogously as the mean recall across the same thresholds.

## A.5 DETAILS OF COCO-OOD

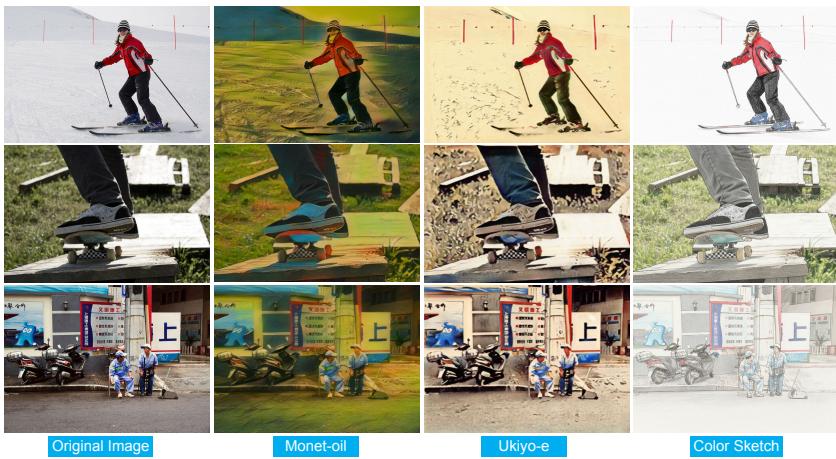


Figure 9: **COCO-OOD visualizations.** We create an OOD split of the COCO validation set by stylizing images with CycleGAN, StyTr2, and Nano-Banana into a Monet-like oil-painting domain, enabling robustness evaluation under appearance shifts.

To complement HumanArt and enable out-of-distribution evaluation under matched content and annotations, we construct COCO-OOD (Fig. 9) by applying artistic style transfer to the original COCO images while preserving human bounding boxes and keypoints. We generate three stylistic variants, each capturing different types of artistic domain shifts. Importantly, for fair comparison and to avoid introducing priors from large-scale pretrained diffusion models, we intentionally adopt the earlier CycleGAN (Zhu et al., 2017) and StyTr2 (Deng et al., 2022) framework rather than more recent style transfer approaches.

For the Monet oil-painting variant, we employ the StyTr2 framework, which produces high-fidelity oil-painting textures and color palettes while maintaining the overall scene geometry. For the Ukiyo-e variant, we adopt the official CycleGAN implementation to translate natural photographs into the Ukiyo-e domain. Finally, we also explore a Color Sketch variant generated using the Nano-Banana model (Nano Banana). While Nano-Banana preserves global shape in most cases, its stylization can occasionally introduce slight pixel-level misalignment with respect to the source images. Because such spatial deviations may compromise compatibility with COCO’s keypoint annotations, we use the Color Sketch variant only as

a supplemental resource for latent-space t-SNE analysis and exclude it from quantitative evaluation. By drawing from multiple style-transfer frameworks, COCO-OOD reduces the bias introduced by any single model and produces stylistic shifts with diverse characteristics. This diversity leads to a more balanced and comprehensive benchmark. Overall, COCO-OOD maintains the geometric content of COCO while producing substantial appearance shifts across distinct artistic domains, providing a controlled and annotation-preserving testbed for assessing cross-domain robustness in pose estimation.

## B DETAILS OF QUANTITATIVE COMPARISON AND ADDITIONAL QUALITATIVE COMPARISON

### B.1 FULL QUANTITATIVE COMPARISON ON COCO

SDPose achieves 81.3 AP / 85.2 AR on the COCO validation set (Table 3) with only 40 training epochs using a 0.95B SD-v2 backbone. It matches the accuracy of Sapiens (82.1–82.2 AP) despite requiring 5× fewer epochs and a smaller backbone, and surpasses GenLoc (+3.7 AP, +4.5 AR). SDPose also outperforms ViTPose++ (+1.9 AP), which relies on multiple auxiliary datasets, while being trained solely on COCO.

Table 3: Quantitative comparison on the COCO validation set.

Model	Input Size	AP	AR
SimpleBaseline (Xiao et al., 2018)	256×192	73.5	79.0
HRNet (Sun et al., 2019)	384×288	76.3	81.2
HRFormer (Yuan et al., 2021)	256×192	77.2	82.0
ViTPose-S (Xu et al., 2022)	256×192	73.8	79.2
ViTPose-B (Xu et al., 2022)	256×192	75.8	81.1
ViTPose-L (Xu et al., 2022)	256×192	78.3	83.5
ViTPose-H (Xu et al., 2022)	256×192	79.1	84.1
ViTPose++-S (Xu et al., 2023)	256×192	75.8	81.0
ViTPose++-B (Xu et al., 2023)	256×192	77.0	82.6
ViTPose++-L (Xu et al., 2023)	256×192	78.6	84.1
ViTPose++-H (Xu et al., 2023)	256×192	79.4	84.8
Sapiens-0.3B (Khirodkar et al., 2024)	1024×768	79.6	83.6
Sapiens-0.6B (Khirodkar et al., 2024)	1024×768	81.2	84.9
Sapiens-1B (Khirodkar et al., 2024)	1024×768	82.1	85.9
Sapiens-2B (Khirodkar et al., 2024)	1024×768	82.2	86.0
<b>SDPose (Ours)</b>	1024×768	<b>81.3</b>	<b>85.2</b>

### B.2 FULL QUANTITATIVE COMPARISON ON HUMANART

On HumanArt (Table 4), SDPose sets a new state of the art with 71.2 AP / 73.9 AR, surpassing large-scale foundation baselines under the same COCO-only training: +1.6 AP over Sapiens-2B (69.6 AP) and +4.7 AP over ViTPose-H (66.5 AP), with consistent AR gains. Compared with traditional baselines such as RTMPose (Jiang et al., 2023) and HRNet (Sun et al., 2019), SDPose further delivers substantial improvements in AP, exceeding them by more than 14 points.

### B.3 FULL QUANTITATIVE COMPARISON ON COCO WHOLEBODY

As shown in Table 5, SDPose achieves 71.5 AP / 78.4 AR on the COCO-WholeBody validation set. This result is highly competitive with the large-scale Sapiens (Khirodkar et al., 2024) models: while Sapiens-2B reaches 74.4 AP with over 2B parameters and long training

918  
 919 Table 4: **Quantitative Comparison on the HumanArt validation set.** Models trained  
 920 on COCO, evaluated with GT bounding boxes. SDPose achieves new state-of-the-art per-  
 921 formance.

Model	AP	AP <sup>50</sup>	AP <sup>75</sup>	AR	AR <sup>50</sup>
RTMPose-T (Jiang et al., 2023)	44.4	72.5	45.3	48.8	75.0
RTMPose-S (Jiang et al., 2023)	48.0	73.9	49.8	52.1	76.3
RTMPose-M (Jiang et al., 2023)	53.2	76.5	56.3	57.1	78.9
RTMPose-L (Jiang et al., 2023)	56.4	78.9	60.2	59.9	80.8
ViTPose-S (Xu et al., 2022)	50.7	75.8	53.1	55.1	78.0
ViTPose-B (Xu et al., 2022)	55.5	78.2	59.0	59.9	80.9
ViTPose-L (Xu et al., 2022)	63.7	83.8	68.9	67.7	85.9
ViTPose-H (Xu et al., 2022)	66.5	86.0	71.5	70.1	87.1
HRNet-W32 (Sun et al., 2019)	53.3	77.1	56.2	57.4	79.2
HRNet-W48 (Sun et al., 2019)	55.7	78.2	59.3	59.5	80.4
Sapiens-1B (Khirodkar et al., 2024)	64.3	82.1	67.9	67.4	83.7
Sapiens-2B (Khirodkar et al., 2024)	69.6	85.3	73.3	72.2	86.8
<b>SDPose (Ours)</b>	<b>71.2</b>	<b>87.3</b>	<b>76.3</b>	<b>73.9</b>	<b>88.6</b>

937  
 938  
 939  
 940 Table 5: **Quantitative comparison on the COCO-WholeBody validation set.**  
 941

942 (a) Body, Feet, Face

Model	Body AP	Body AR	Feet AP	Feet AR	Face AP	Face AR
HRNet (Sun et al., 2019)	70.1	77.3	58.6	69.2	72.7	78.3
VitPose+-L (Xu et al., 2023)	75.3	-	77.1	-	63.0	-
VitPose+-H (Xu et al., 2023)	75.9	-	77.9	-	63.6	-
RTMPose-x (Jiang et al., 2023)	71.4	78.4	69.2	81.0	88.8	92.2
DWPose-l (Yang et al., 2023)	72.2	78.9	70.4	81.7	88.7	92.1
Sapiens-0.3B (Khirodkar et al., 2024)	66.4	73.4	67.3	78.4	87.1	91.2
Sapiens-0.6B (Khirodkar et al., 2024)	74.3	80.2	79.4	87.0	89.5	92.9
Sapiens-1B (Khirodkar et al., 2024)	77.4	82.9	83.0	89.8	90.7	93.6
Sapiens-2B (Khirodkar et al., 2024)	<b>79.2</b>	<b>84.6</b>	<b>84.1</b>	<b>90.9</b>	<b>91.2</b>	<b>93.8</b>
<b>SDPose (Ours)</b>	77.9	83.4	81.5	88.7	88.5	92.2

952  
 953 (b) Hands and Whole-body

Model	Hand AP	Hand AR	Whole AP	Whole AR
HRNet (Sun et al., 2019)	51.6	60.4	58.6	67.4
VitPose+-L (Xu et al., 2023)	54.2	-	60.6	-
VitPose+-H (Xu et al., 2023)	54.7	-	61.2	-
RTMPose-x (Jiang et al., 2023)	59.0	68.5	65.3	73.3
DWPose-l (Yang et al., 2023)	62.1	71.0	66.5	74.3
Sapiens-0.3B (Khirodkar et al., 2024)	58.1	67.1	62.0	69.4
Sapiens-0.6B (Khirodkar et al., 2024)	65.4	74.0	69.5	76.3
Sapiens-1B (Khirodkar et al., 2024)	69.2	77.1	72.7	79.2
Sapiens-2B (Khirodkar et al., 2024)	<b>70.4</b>	<b>78.1</b>	<b>74.4</b>	<b>81.0</b>
<b>SDPose (Ours)</b>	65.2	74.0	71.5	78.4

966 schedules, SDPose attains comparable accuracy with a smaller 0.95B backbone trained for  
 967 only 42 epochs. In terms of sub-part analysis, our method closely matches Sapiens-2B on  
 968 body (77.9 vs. 79.2) and feet (81.5 vs. 84.1), while maintaining strong performance on face  
 969 (88.5 vs. 91.2). The largest gap appears on hand keypoints (65.2 vs. 70.4), reflecting the  
 970 intrinsic difficulty of fine-grained articulation. Nevertheless, compared with classical base-  
 971 lines such as HRNet (Sun et al., 2019), RTMPose (Jiang et al., 2023), or DWPose (Yang  
 et al., 2023), SDPose shows substantial gains of +6–12 AP across whole-body evaluation.

972  
 973 **Table 6: Quantitative Comparison on the COCO-OOD-Monet Wholebody vali-**  
 974 **dation set.**

975  
 976 (a) **Body, Feet, Face**  
 977 

Model	Body AP	Body AR	Feet AP	Feet AR	Face AP	Face AR
Sapiens-1B	52.1	58.6	55.9	66.2	57.6	63.0
Sapiens-2B	59.9	65.8	63.8	72.4	58.4	64.2
SDPose (Ours)	<b>60.0</b>	<b>66.1</b>	<b>62.5</b>	<b>72.0</b>	<b>64.2</b>	<b>69.9</b>

  
 981 (b) **Hands and Whole-body**  
 982 

Model	L-Hand AP	L-Hand AR	R-Hand AP	R-Hand AR	Whole AP	Whole AR
Sapiens-1B	43.1	52.0	41.5	50.6	38.7	46.8
Sapiens-2B	48.2	56.8	46.8	55.3	44.4	53.0
SDPose (Ours)	46.3	55.2	44.9	54.4	<b>46.6</b>	<b>54.8</b>

983  
 984 **B.4 FULL QUANTITATIVE COMPARISON ON COCO-OOD WHOLEBODY**

985 Detailed whole-body pose estimation results on COCO-OOD Monet are reported in Ta-  
 986 ble 6. Breaking down by body part, SDPose matches or surpasses Sapiens on the most  
 987 stable regions: body (60.0 AP vs. 59.9) and feet (62.5 vs. 63.8), and delivers a notable  
 988 margin on face landmarks (+5.8 AP, 64.2 vs. 58.4), highlighting the reliability of SD-  
 989 Pose’s features under appearance shifts. For hands, which are the most challenging due to  
 990 fine-grained articulation and limited resolution, SDPose attains 46.3/44.9 AP on left/right  
 991 hands, respectively, remaining competitive but still slightly behind Sapiens-2B (48.2/46.8).  
 992 Nevertheless, the overall whole-body AP/AR of 46.6/54.8 establishes SDPose as the most  
 993 robust framework under OOD whole-body evaluation. Detailed whole-body pose estimation  
 994 results on COCO-OOD Ukiyo-e are presented in Table 7. Breaking down by region, SDPose  
 995 remains competitive with Sapiens-2B on the more structurally stable parts: body (61.2 AP  
 996 vs. 62.2) and feet (64.7 vs. 67.3), and shows a clear advantage on face landmarks (+4.8 AP,  
 997 64.7 vs. 59.9), reflecting the stronger style-invariant geometry encoded by diffusion U-Net  
 998 features under artistic shifts. For hands—the most challenging subset due to fine-scale pose  
 999 variation—SDPose achieves 45.9/44.1 AP on left/right hands, remaining close to Sapiens-  
 1000 2B (49.8/46.7) despite the latter’s larger model capacity. Importantly, SDPose attains the  
 1001 highest whole-body AP/AR (47.7/56.4), demonstrating consistently stronger robustness  
 1002 than both Sapiens-1B and Sapiens-2B when evaluated under Ukiyo-e style transformations.

1003 **Table 7: Quantitative Comparison on the COCO-OOD-Ukiyoe Wholebody vali-**  
 1004 **dation set.**

1005  
 1006 (a) **Body, Feet, Face**  
 1007 

Model	Body AP	Body AR	Feet AP	Feet AR	Face AP	Face AR
Sapiens-1B	54.0	60.9	57.9	69.9	58.7	65.3
Sapiens-2B	<b>62.2</b>	<b>68.2</b>	<b>67.3</b>	<b>76.5</b>	59.9	66.9
SDPose (Ours)	61.2	67.5	64.7	75.0	<b>64.7</b>	<b>71.2</b>

  
 1016 (b) **Hands and Whole-body**  
 1017 

Model	L-Hand AP	L-Hand AR	R-Hand AP	R-Hand AR	Whole AP	Whole AR
Sapiens-1B	43.7	52.8	40.9	50.6	40.5	49.4
Sapiens-2B	<b>49.8</b>	<b>58.7</b>	<b>46.7</b>	<b>55.7</b>	46.6	55.8
SDPose (Ours)	45.9	55.7	44.1	54.3	<b>47.7</b>	<b>56.4</b>

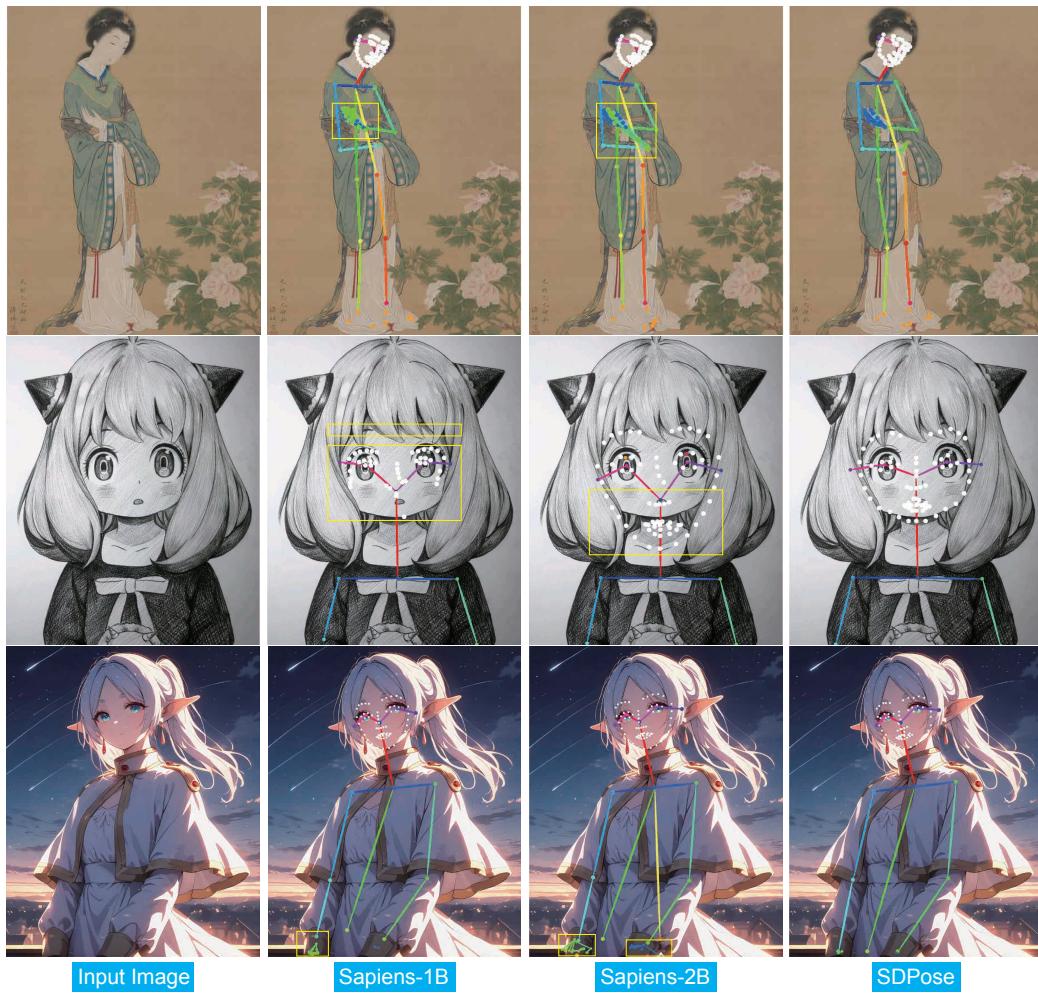
1026  
1027 B.5 QUALITATIVE COMPARISON FOR WHOLE-BODY POSE ESTIMATION  
1028

Figure 10: **Comparison on Stylized Paintings: Sapiens WholeBody vs. SDPose WholeBody.** All erroneous predictions are highlighted with yellow boxes. SDPose yields fewer false positives and notably better facial keypoint localization.