

000
001
002
003
004
005
006
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
053054
055
056
057
058
059
060
061
062
063
064
065
066
067
068
069
070
071
072
073
074
075
076
077
078
079
080
081
082
083
084
085
086
087
088
089
090
091
092
093
094
095
096
097
098
099
100
101
102
103
104
105
106
107

Improving 2D Human Pose Estimation across Unseen Camera Views with Synthetic Data

Anonymous WACV Algorithms Track submission

Paper ID 442

Abstract

Human Pose Estimation is a thoroughly researched problem; however, most datasets focus on the side and front-view scenarios. We address the limitation by proposing a novel approach that tackles the challenges posed by extreme viewpoints and poses. We introduce a new method for synthetic data generation – RePoGen, RarE POses GENerator – with comprehensive control over pose and view to augment the COCO dataset. Experiments on a new dataset of real images show that adding RePoGen data to the COCO surpasses previous attempts to top-view pose estimation and significantly improves performance on the bottom-view dataset. Through an extensive ablation study on both the top and bottom view data, we elucidate the contributions of methodological choices and demonstrate improved performance. We will release the code and the datasets with the publication.

1. Introduction

The availability of large-scale, manually annotated datasets has greatly advanced research in human pose estimation from 2D monocular images. Current datasets primarily focus on camera viewpoints from what we call *an orbital view*, i.e. side, front, and back views, where challenges such as occlusion by objects or individuals are prevalent. Similarly, they focus on common poses like standing, sitting, or walking by sampling everyday activities. As a result, much of the research has been dedicated to tackling occlusion. Specialized datasets have been curated to evaluate the effectiveness of pose estimation models in scenarios involving occluded individuals.

On the other hand, the issue of unusual viewpoints has received less attention. In what we refer to as *extreme viewpoints* (top and bottom view; the complement of orbital view), the appearance of humans significantly differs from that of the orbital view. Although such views are less common in everyday activities and videos, they frequently ap-

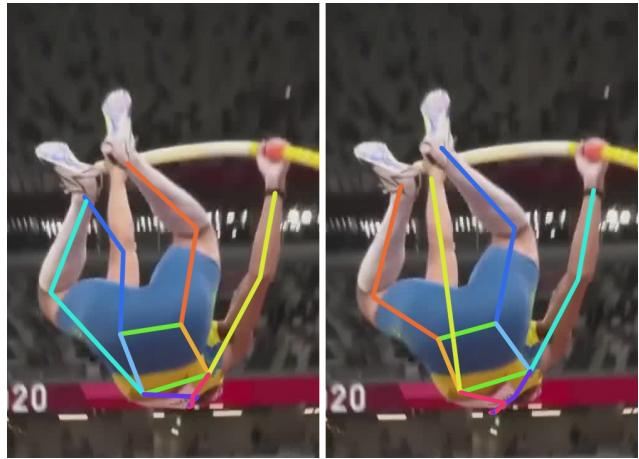


Figure 1. Pose estimation trained on COCO (left) and by our method (right). The COCO model mistakes the left and right sides and interprets the **right hand** as the **left leg** and the **right leg** as the **left hand** (color indicates the corresponding label).

pear in sports or surveillance footage. Annotating persons in extreme views poses considerable challenges as human annotators struggle to comprehend scenes unfamiliar to the human eye.

We employ an SMPL-based [23] synthetic data approach similar to previous methods [8, 30] to address the scarcity of training data. However, we distinguish ourselves by generating novel poses, even if they occasionally deviate from anatomical accuracy. We allow for the possibility of body parts, like limbs, intersecting with each other, as long as the overall pose maintains physical plausibility. Minor mesh intersections can simulate body deformations without impeding training. This novel approach allows us to generate new poses from a wider distribution than previous methods. We demonstrate that pose variability, combined with novel views, is crucial for accurate pose estimation in sports, where extreme poses and extreme views are prevalent.

108 We introduce a novel method for generating likely
109 realistic poses and utilize them to augment existing datasets,
110 thereby incorporating novel views and poses. Furthermore,
111 we demonstrate the applicability of our approach to the top
112 view, which is on par with or potentially superior to previous
113 methods. The main contributions of the paper are:
114

- 115 1. RePoGen - a new method for generating synthetic real-
116 looking images with humans.
- 117 2. The RePoGen dataset - a new dataset of synthetic im-
118 ages prioritizing rare poses and viewpoints.
- 119 3. RePo - a new manually annotated dataset of real im-
120 ages of rare poses from the top and bottom views en-
121 abling comprehensive evaluation of pose estimation
122 from unusual views.
- 123 4. We demonstrate a significant increase in the pose es-
124 timation accuracy on extreme views without harm-
125 ing COCO performance by augmenting the existing
126 COCO dataset with RePo synthetic data.

127 We will release the RePoGen code and the synthetic Re-
128 PoGen and real-world, annotated, RePo datasets. Additionally,
129 we provide enhanced annotations for the previously
130 published PoseFES dataset [32].

134 2. Related Work

135 Numerous datasets have been developed to support ad-
136 vancements in human pose estimation. Real-world datasets
137 like COCO [17] and MPII [1] offer diverse images that
138 capture human poses in everyday scenes, while the LSP
139 dataset [14] focuses on sports-related poses. To address the
140 challenge of occlusion, specialized datasets such as OCHu-
141 man [35] and CrowdPose [16] have been curated, enabling
142 the evaluation of pose estimation algorithms in occluded
143 scenarios.

144 Several models have emerged, demonstrating significant
145 advancements in accuracy and performance. These mod-
146 els primarily fall into top-down approaches, which rely on
147 bounding boxes as input for pose estimation. Among these
148 models, ViTPose [31] stands out as the current SOTA on
149 the COCO dataset leveraging the transformer architecture.
150 Similarly, models such as SWIN [20] and PSA [19] also
151 employ transformer-based architectures, although they per-
152 form slightly below ViTPose in terms of accuracy.

153 An alternative approach that garnered attention is the
154 HRNet model [27], which combines convolutional neural
155 networks with an integral part, Unbiased Data Processing
156 [10]. This combination yields excellent results and has be-
157 come a common baseline for evaluating the performance of
158 new pose estimation methods.

159 Addressing the challenges posed by occlusion and
160 crowded scenes, specialized models have been developed to
161

162 focus on these specific scenarios. For example, the I2RNet
163 [7] is a transformer-based network designed to tackle the
164 challenges of occlusion and crowd-related issues.

165 Furthermore, proper data processing techniques have
166 been proposed to enhance the performance of pose estima-
167 tion models. The DARK algorithm [33] and the UDP (Un-
168 grouped Distance Parameterization) method [10] are two
169 notable papers that highlight the importance of data pro-
170 cessing in achieving superior results.

171 To facilitate pose estimation research and development,
172 the MMPose framework [6] has emerged as a comprehen-
173 sive resource. It offers an extensive model zoo and many
174 pre-trained models, including the widely used HRNet.

175 Synthetic datasets have also played a significant role in
176 augmenting the available data and expanding the range of
177 pose variations. The THEODORE+ dataset [32] provides
178 a synthetic collection of top-view videos generated using
179 a game engine. These videos depict individuals walking
180 in a room, although they only pro 13 keypoints instead of
181 the more commonly used 17. Synthetic datasets like SUR-
182 REAL [30] and PanopTOP [8] utilize the SMPL model [21],
183 fitting it to measured 3D point clouds of real poses from
184 datasets such as Human36M [13] and Panoptic [15]. How-
185 ever, PanopTOP has limitations regarding low resolution
186 and issues with ghost hands, which should be considered.

187 The estimation of poses from extreme viewpoints is
188 another research area of interest. The WEPDFO-Pose
189 dataset [11] represents the largest dataset of top-view im-
190 ages for pose estimation. Although specialized for top-view
191 poses, it is noteworthy that most people captured in the
192 dataset are from the orbital view due to fisheye lens distor-
193 tion. Similarly, the PoseFES dataset [32], designed for eval-
194 uating top-view human pose estimation, also suffers from a
195 prevalence of orbital views caused by fisheye lens distor-
196 tion. Another dataset, ITop [9], focuses on pose estimation
197 from top-view depthmaps with no RGB images available.

198 Data augmentation is critical in addressing the scarcity
199 of annotated real-world data for human pose estimation.
200 Various methods have been introduced to tackle this chal-
201 lenge, often involving human parsing techniques for body
202 part segmentation. HumanPaste [18] and AdversarialAug-
203mentation [3] employ strategies to simulate occlusion by
204 pasting additional people or selective body parts. Similarly,
205 JointlyOD [24] and NearbyPersonOD [5] augment data by
206 introducing body parts or whole bodies to mimic occlusion
207 and crowded scenarios.

208 While these augmentation methods prove effective for
209 specific challenges, they do not directly address the prob-
210 lem of unseen viewpoints. In contrast, generating synthetic
211 data using game engines have been explored to introduce
212 variability. However, datasets created with game engines,
213 such as PoseFES [32] and LetsPF [25], often suffer from
214 limited pose variability, typically showcasing walking or a

216

narrow range of everyday activities.

217

Another avenue for synthetic data generation involves fitting the SMPL model [21] to 3D point clouds obtained from motion capture systems. For example, SURREAL [30] fits the SMPL model to the Human36M dataset, providing a pool of textures applicable to SMPL models. Similarly, PanopticTOP [8] employs the SMPL model fitted to the Panoptic dataset. However, these methods face challenges in fitting the model to point clouds, resulting in issues such as ghost hands. Furthermore, the limitations of motion capture systems make capturing extreme dynamic poses or new poses challenging. SyntheticHF [29] estimates the SMPL pose and shape from a monocular image and modifies the shape while preserving the pose, creating data resembling SURREAL and Panoptic. However, this approach has limitations due to the initial SMPL estimate, resulting in difficulties handling poses beyond its accurate capture.

234

Efforts have also been made to enhance the realism of the SMPL model. SMPL-X [23] enhances the previous model with hand poses and facial expressions. PoseNDF [28] learns a manifold of known poses, enabling the generation of random realistic poses within the manifold. Similarly, CAPE [22] introduces a clothing layer on top of existing SMPL models, aiming to narrow the domain gap between generated and real data.

242

GAN-based methods like SynthesizeAnyone [12], UnpairedPG [4], and SynthesizingIO [2] generate synthetic data by preserving the given pose or style. On the other hand, diffusion-based methods such as StableDiffusion [26] and ControlNet [34] offer promising approaches for synthetic data generation, allowing control over the rendered images. However, both approaches have limitations regarding extreme views and rare poses due to the need for more training data.

251

252

253

3. Method

254

This section provides a detailed description of our approach to enhancing an existing dataset using synthetic data generation. We developed a novel method inspired by prior works that offer enhanced control over pose parameters. Unlike previous approaches that relied on re-using point clouds from motion capture, RePoGen allows us to define a pose simplicity and generate individuals in rare poses. Although the realism of the generated poses is not guaranteed, we demonstrate that it is not a prerequisite for effective performance.

265

The proposed RePoGen pipeline is outlined in the Fig. 3. Following paragraphs present a step-by-step walkthrough of the RePoGen data generation process, highlighting the main techniques employed to achieve pose control and generate diverse synthetic data.

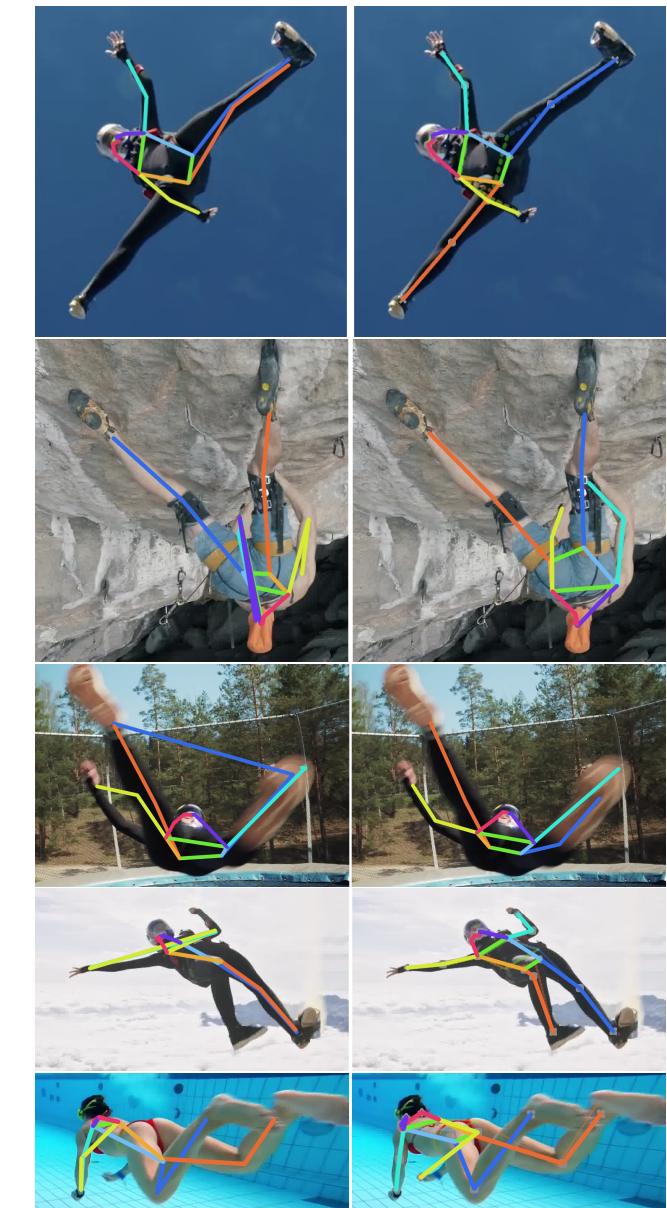


Figure 2. Examples from the RePo test set. ViTPose-s estimates when trained on COCO (left) and on RePoGen data (right). Colors as in Fig. 1 – right hand, right leg, left hand and left leg

3.1. Pose Generation

RePoGen leverages the SMPL-X model [23], which defines 21 body joints with free rotation around three axes each. In addition to the basic SMPL model [21], SMPL-X also includes joints for hands and face. The rotation angles for the face and hand joints are randomly determined, as they do not influence the 17 COCO keypoints.

We sample each body angle from an asymmetrical normal distribution, composed of two normal distributions with

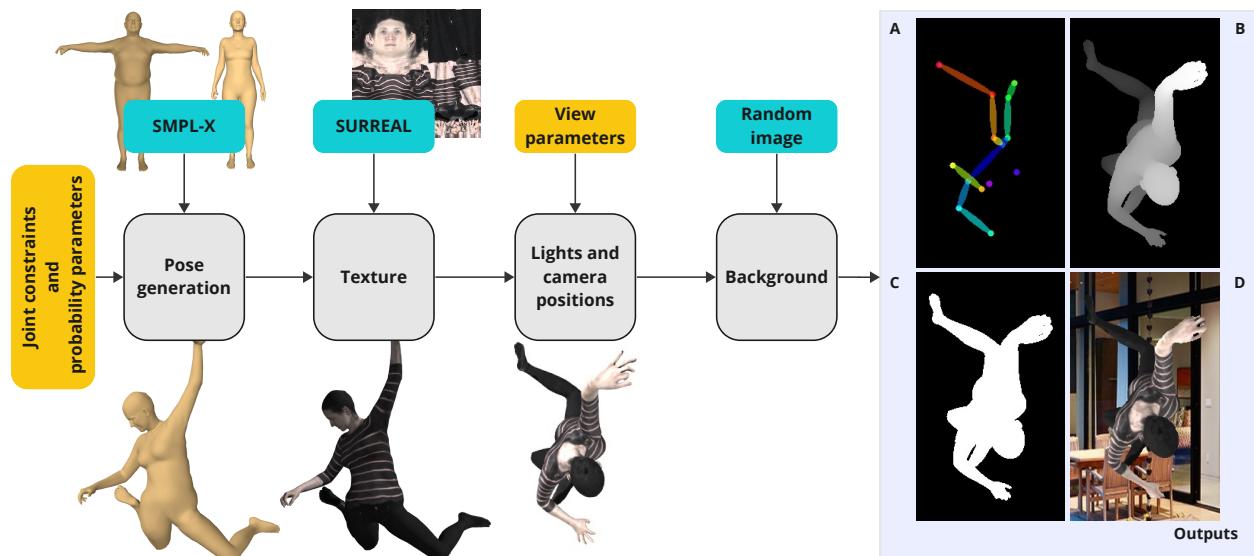
324
325
326
327
328
329
330
331
332
333
334
335
336
337
338
339
340
341
342

Figure 3. RePoGen synthetic data generation pipeline. All steps are detailed in Sec. 3. The ground truth outputs of the method are (A) 2D and 3D keypoints, (B) the depth map, (C) the mask, and (D) an RGB image.

343

344

345

346

347

348

349

350

351

352

353

354

355

356

357
358
359
360
361

Figure 4. Examples of joint angle distributions used in data generation. Baseline - a hand-crafted joint angle distribution approximating statistics of common poses. Uniform sampling of joint angles generates many extreme poses.

362
363
364
365
366
367
368
369
370
371
372
373
374
375
376
377

different variances, visualized in Fig. 4, to generate diverse poses. Each angle has its unique constraints and mean. This distribution allows us to generate pose angles centered around a standard pose, with unique and asymmetric ranges for each joint. It is a hand-crafted approximation of angle distribution in common poses.

By applying constraints on joint rotation, a substantial portion of the pose space, primarily composed of unrealistic poses, is effectively eliminated. The remaining poses are highly likely to exhibit realistic characteristics, although some instances of mesh intersection may occur. However, these small-scale mesh intersections do not pose significant issues during training, as they effectively simulate minor body deformations within the rendered images. The major advantage of our approach is the ability to generate rare

378
379
380
381
382
383
384
385
386
387
388
389
390
391
392
393
394
395
396
397
398
399400
401
402
403
404
405
406
407
408
409
410
411
412
413
414
415
416
417
418
419
420
421
422
423
424
425
426
427
428
429
430
431

poses that are not present in previous datasets. On the other hand, it is important to acknowledge the inherent limitation of the SMPL-X model, which represents the human body with only 21 joints. In comparison, the actual human body consists of over 300 joints. This discrepancy poses challenges, particularly in accurately modeling complex spine rotations.

Another advantage of the methos is the ability to control the complexity of the generated poses using a single parameter referred to as *pose simplicity* α . By scaling the distribution by a constant, we restrict the pose space, and generated poses are closer to the standard pose. Changing the standard pose mathematically means changing the mean of the composed distribution. We experimented with two standard poses - standing straight and the default SMPL pose. Additionally, we introduce the option to sample joints from a uniform distribution instead of the composed normal distribution, which produces more frequent extreme poses. The ablation study in Sec. 4.5 refers to this option as *uniform distribution*.

Last, we changed the default pose to standing straight with hands along the body instead of the default SMPL pose with hands horizontally. Both poses are visualized in the Fig. 3.

The output of this stage is a triangular mesh representing a human body in a randomly generated pose. The generated mesh is smooth and without noise, ensuring a consistent and visually coherent pose representation.

432 **3.2. Texture**
 433
 434 Once the random pose is generated, we apply a ran-
 435 domized texture to the mesh. For this purpose, we uti-
 436 lize textures provided by the SURREAL project [30] and
 437 do not differentiate between male and female textures. If
 438 no texture is applied (as examined in the ablation study in
 439 Sec. 4.5), we color the mesh to resemble natural skin tones.
 440 This approach ensures that the generated synthetic data ex-
 441 hibits variation in texture, contributing to a more realistic
 442 appearance.

443 **3.3. Lights and Camera Positions**

444
 445 In our pose generation technique, we randomly sample
 446 both light and camera positions from a surface of a unit
 447 sphere. Initially, we distribute five light sources randomly
 448 on the unit sphere, creating shadows on the texture to en-
 449 hance the realism of the generated data.

450 All distances utilized in our pose generation process are
 451 measured in the coordinates of the SMPL-X model. The
 452 SMPL unit corresponds to a length of approximately less
 453 than 1 meter. The coordinate system is visually represented
 454 in the Fig. 6, aiding in understanding the coordinate trans-
 455 formations involved in RePoGen.

456 **3.4. Random Background**

457
 458 The final component for generating visually appealing
 459 images is the background. We incorporate a random image
 460 as the background and crop the rendered scene to a 1.25
 461 multiple of the bounding box size. When selecting back-
 462 ground images, we ensure that they depict environments
 463 where people are commonly observed. However, we refrain
 464 from including discernible individuals in the background,
 465 which could confuse the network since we do not focus on
 466 crowded scenes.

467 **3.5. Ground Truth Extraction**

468
 469 The output of the pipeline includes not only the rendered
 470 RGB image but also the corresponding ground truth infor-
 471 mation. We first extract the depth map from the trian-
 472 gular mesh representation to obtain the ground truth. This
 473 depthmap is then used to generate a segmentation mask
 474 through thresholding. The segmentation mask defines the
 475 bounding box.

476
 477 However, determining the visibility of joints is a com-
 478 plex process, as the joints of the SMPL-X model are posi-
 479 tioned within the triangular mesh and are, therefore, always
 480 hidden from view in the rendered image. To address this,
 481 we define a neighborhood around each joint and consider
 482 the joint visible if at least one vertex from its respective
 483 neighborhood is visible in the image. The size of the neigh-
 484 borhood is proportional to the joint size and is determined
 485 based on the human annotation error defined in the OKS

Dataset name	# of poses	486
PoseFES Top	431	487
RePo (Bottom Val)	31	488
RePo (Bottom Test)	94	489
RePo (Bottom Seq)	62	490
RePo (Top Val)	91	491
		492
		493

Table 1. The number of annotated poses for the new datasets.

494
 495
 496
 metric from the COCO dataset. This approach allows us to
 497 estimate the visibility of the joints and accurately generate
 498 the corresponding ground truth annotations for evaluation
 499 and training purposes.

500 **4. Experiments**

501 **4.1. Implementation Details**

502
 503
 504
 To optimize computation power and time efficiency, we
 505 primarily conduct experiments using the ViTPose-s model
 506 unless otherwise specified. The training parameters align
 507 with the ViTPose model, with a batch size of 128 and a
 508 base learning rate 5e-5. We follow the training paradigm
 509 from [32] and fine-tune the model pretrained on the COCO
 510 dataset.

511
 512
 513
 514
 515
 To focus on analyzing and improving the pose estimation
 516 model, we utilize ground truth bounding boxes to crop in-
 517 dividuals from the images. This approach is chosen to miti-
 518 gate errors from detectors, particularly in extreme views.

519
 520
 521
 All synthetic images used in experiments are generated
 522 exclusively through RePoGen, with a preference for the top
 523 or bottom views. Synthetic data from orbital views are not
 524 generated as they provide no notable improvement.

525
 526
 527
 During training, the model is not exposed to any real
 528 extreme view images that are not present in the original
 529 COCO dataset. Instead, all additional data used for train-
 530 ing purposes are synthetically generated. The model used
 531 for comparison with other approaches used 3 000 images.
 532 The ablation study was done using 1 000 images.

533
 534
 535
 536
 537
 538
 539
Rotation. During training, we incorporate extensive ro-
 530 tation data augmentation of COCO and synthetic images. In
 531 experiments labeled as *w/o rotation*, we follow the standard
 532 rotation augmentation up to 40°, while in other cases, we
 533 apply a rotation up to 180°.

534 **4.2. Datasets**

535
 536
 537
 538
 539
 We created a new dataset to evaluate pose estimation
 530 from extreme views in real-world data. We conduct experi-
 531 ments on the following datasets:

532
 533
 534
 535
 536
 537
 538
 539
COCO. [17] This standard dataset is commonly used for
 531 human pose estimation. It contains approximately 250,000
 532 annotated poses from various everyday activities. However,

540 the COCO dataset includes very few images captured from
 541 extreme views.
 542

543 **PoseFES.** [32] PoseFES is a manually annotated dataset
 544 captured by a ceiling-mounted fisheye camera, serving as
 545 the solely available top-view dataset for human pose esti-
 546 mation. Although we know another dataset (WEPDTON-
 547 Pose [11]), our attempts to obtain it from the authors were
 548 unsuccessful. PoseFES consists of two sequences: one fo-
 549 cusing on two well-separated individuals, while the second
 550 involves multiple people interacting and creating challeng-
 551 ing scenarios with occlusions. We primarily utilize the first
 552 sequence for testing to align with our research focus on
 553 single-person human pose estimation. However, since this
 554 sequence predominantly contains orbital view images due
 555 to the fisheye transformation, we extracted a subset of im-
 556 ages and annotations from both sequences to create Pose-
 557 FES Top, which consists of images of individuals directly
 558 beneath the camera, representing the extreme top view.
 559

560 **Bottom Val, Test, and Seq.** Since no existing datasets
 561 specifically cater to bottom-view data, we created a new
 562 dataset called RePo (RarE POses) to evaluate our approach.
 563 The dataset consists of images extracted from various sports
 564 videos obtained from YouTube. The most common sports
 565 featured are swimming, climbing, and skydiving. The Val
 566 and Test sets possess similar structures derived from com-
 567 parable videos, while the Seq set comprises consecutive
 568 frames from one specific video of the pole vault. We em-
 569 ploy the Seq set to demonstrate that substantial rotations of
 570 the person often accompany extreme views. Examples of
 571 real images from the new dataset are in the Sec. 2.
 572

573 **Top Val.** Similar to the Bottom datasets, this dataset is
 574 collected from sports videos focusing on the top-view per-
 575 spective. It serves as a validation set during the top-view
 576 training phase. The Top Val is also part of the new RePo
 577 dataset.
 578

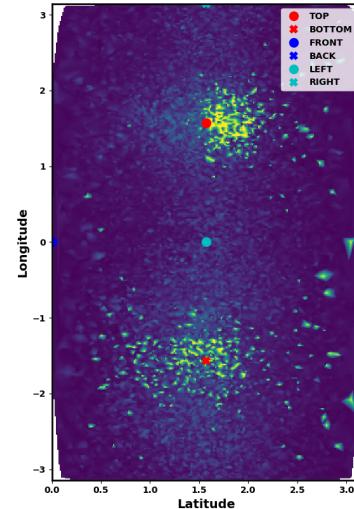
579 For further reference, a summary of the new datasets in-
 580 troduced in this work is presented in the Tab. 1.
 581

582 **Metrics.** All experiments were conducted following the
 583 COCO-style settings. The evaluation metric used was OKS-
 584 based AP (average precision), as specified in the COCO
 585 dataset [17].
 586

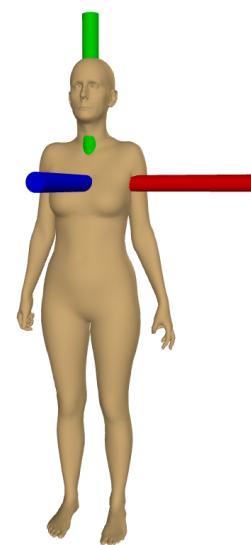
587 4.3. Viewpoint Dependency Analysis

588 RePoGen enables us to analyze the performance of state-
 589 of-the-art methods from different viewpoints. We analyze
 590 the performance in controlled settings, where individuals
 591 are well-separated and have clearly defined bounding boxes.
 592 Given the vast and complex pose space, we do not sam-
 593 ple poses systematically. Instead, we sample 4 000 random
 594 poses with uniform pose simplicity between 1.0 and 3.0 and
 595 render each one from 5 views uniformly distributed along a
 596 sphere surface resulting in 20 000 images.
 597

598 The analysis is based on the ViTPose model, which
 599



599 Figure 5. Pose estimation quality as
 600 a function of viewpoint, in spherical
 601 coordinates. Darker colors mark
 602 higher OKS (smaller error).
 603



604 Figure 6. SMPL coordi-
 605 nates: x (red), y (green),
 606 and z (blue)
 607

608 demonstrated the best performance on the COCO dataset
 609 at the time of writing. However, the results were also veri-
 610 fied on other models, namely SWIN and HRNet, as imple-
 611 mented in the MPMoP framework.
 612

613 The Fig. 5 visualizes the errors of each sample in a sph-
 614 erical coordinate system with a fixed radius, where the hori-
 615 zontal and vertical axis represents latitude and longitude,
 616 respectively. The top view is indicated by a red circle at
 617 coordinates $[\frac{\pi}{2}, \frac{\pi}{2}]$, and the bottom view is denoted by a red
 618 cross at coordinates $[\frac{\pi}{2}, -\frac{\pi}{2}]$. The front view corresponds to
 619 coordinates $[0, 0]$, located at the left edge of the image. The
 620 OKS score of each sample is indicated by the color of the
 621 point, with darker blue indicating a higher score and yellow
 622 representing a lower score.
 623

624 As expected, the findings showed that state-of-the-art
 625 methods performed poorly on extreme views. Notably, the
 626 top-back view performed worse than the top-front view,
 627 while the error distribution around the bottom view ap-
 628 peared symmetric. The spread of the error around the bot-
 629 tom view is wider. The image is not smooth because some
 630 poses with lower pose simplicity proved challenging even
 631 in orbital views.
 632

633 4.4. Comparison with baseline

634 The comparison table Tab. 2 illustrates the perfor-
 635 mance comparison between the baseline model (off-the-
 636 shelf ViTPose-s trained on the COCO dataset) and the
 637 proposed approach. We show variants with bottom-view,
 638 top-view, and mixed bottom and top-view RePoGen syn-
 639 thetic images. The results highlight a notable improvement
 640 achieved through the utilization of synthetic data and train-
 641

648	Dataset	Bottom Test	PoseFES Top
649	COCO	35.1	42.0
650	RePoGen (bottom)	61.8	52.9
651	RePoGen (top)	46.3	53.9
652	RePoGen (top+bottom)	53.9	54.1
653			
654			

655 Table 2. AP on the RePo Bottom Test set and PoseFES Top; training
656 on COCO and sets of 3 000 images from the RePoGen.
657

658	Dataset	PoseFES 1
659	COCO	75.7
660	THEODORE+	76.1 [†]
661	RePoGen (30 epochs)	77.9
662	RePoGen	79.5
663		
664		

665 Table 3. AP on the PoseFES1 set; training on COCO,
666 THEODORE+ by [32] and RePoGen dataset. The result marked
667 ([†]) taken from [32].
668

669 ing with rotation augmentation. Interestingly, incorporating
670 synthetic data from the bottom view enhances the model’s
671 performance on the bottom and top view, suggesting a simi-
672 larity between the two extreme view domains. Similarly,
673 training with synthetic data from the top-view demonstrates
674 improvements across top-view and bottom-view scenarios.
675

676 To facilitate a comprehensive comparison of RePoGen
677 with prior research, we conducted fine-tuning of the HR-
678 Net [27] model from the MPMpose [6] model zoo follow-
679 ing the same procedure as described by Yu et al. [32].
680 The performance evaluation, as presented in Tab. 3, show-
681 cases the effectiveness of RePoGen in comparison to the
682 THEODORE+ dataset and a model trained solely on the
683 COCO dataset. We observed that surpassing the prescribed
684 30-epoch fine-tuning, as mentioned in [32], led to further
685 improvements in performance. Consequently, we report re-
686 sults for the 30-epoch mark and the best-achieved perfor-
687 mance. RePoGen achieves superior results despite utilizing
688 significantly fewer data, incorporating 3000 synthetic im-
689 ages compared to 160,000 THEODORE+ images.
690

691 4.5. Ablation Study

692 We analyze and evaluate the influence of each compo-
693 nent individually, as described in the following paragraphs.
694 Throughout the ablation study, the strong rotation augmen-
695 tation is consistently applied, and unless otherwise speci-
696 fied, 1000 RePoGen are used for experimentation.
697

698 **Number of images.** The Tab. 4 provides insights into the
699 impact of adding additional images to the COCO dataset.
700 With the COCO train set already containing over 200 000
701 poses, adding 5 000 images represents approximately 2%
of the dataset, resulting in minimal impact on training time.

702	# of images	Bottom Test	Bottom Seq
703	500	54.1	86.1
704	1000	59.1	89.0
705	3000	61.8	90.5
706	5000	58.8	86.1
707			
708			

709 Table 4. AP on the Bottom dataset of RePo; training with different
710 number of RePoGen images.
711

712	RePoGen data	Bottom Test	Bottom Seq
713	baseline	59.1	89.0
714	w/o rotation	45.9	72.3
715	w/o background	56.2	85.2
716	w/o texture	59.5	88.2
717	default SMPL pose	60.4	88.4
718	uniform distribution	59.2	89.8
719			
720			

721 Table 5. Ablation study. Training without various components -
722 AP comparison on the Bottom dataset of RePo.
723

724 Remarkably, even including as few as 500 images yields no-
725 ticeable improvements. However, saturation is observed at
726 around 3 000 images, beyond which further additions may
727 have a marginal negative effect on performance probably
728 due to the overfit to the synthetic data.
729

730 **Texture and background.** The Tab. 5 validates other
731 design choices made in the pose generation technique. It
732 demonstrates the improvement in performance on the Test
733 set, which assesses the model’s ability to handle extreme
734 views. Additionally, the results on the Seq set, which in-
735 cludes extreme and adjacent views, further support the ef-
736 fectiveness of these design choices. Notably, including
737 background images contributes to a modest enhancement
738 in performance. On the other hand, adding random texture
739 does not yield significant improvements, suggesting that the
740 realism of the data may not be a crucial factor in this con-
741 text.
742

743 **Rotation.** Incorporating stronger rotation yields signifi-
744 cant performance improvements. The effect is particularly
745 pronounced in the Seq set, where the presence of views ad-
746 jacent to the extreme ones amplifies the difference even fur-
747 ther. Even without rotation, our approach outperforms the
748 off-the-shelf model, highlighting the importance of includ-
749 ing extreme view data in the training. Consequently, it is
750 advisable always to employ rotation data augmentation up
751 to 180° for applications involving pose estimation in videos
752 with extreme views.
753

754 **Default SMPL pose and uniform distribution.** The
755 impact of *uniform joint angle distribution* and the *default
756 SMPL pose* remains inconclusive. In the experimentation
757 with the Bottom datasets, training models with uniform dis-
758

tribution proved advantageous compared to the baseline. However, contrasting results were observed when training on the top view and evaluating on the PoseFES dataset. This discrepancy may be attributed to the nature of the Bottom datasets, which encompass sports activities characterized by extreme poses. In contrast, the PoseFES dataset primarily features individuals engaged in walking and standing. Similar results can be observed with the default SMPL pose. Both approaches generate poses from less usual distribution than the baseline. The observed difference in performance compared to the baseline is approximately 0.5 percentage points, indicating a relatively minor effect. Nonetheless, employing poses aligned with the target domain appears preferable for optimal results.

5. Conclusions

In conclusion, this paper presented a novel method for generating synthetic images (RePoGen) with accurate human pose ground truth by incorporating constraints on joint rotation. The view dependency of performance in SOTA methods was thoroughly analyzed, revealing substantial performance degradation in extreme views. We then trained a state-of-the-art model on the COCO dataset enhanced by RePoGen data to improve performance in extreme views. The key findings can be summarized as follows:

1. The SOTA methods perform worse in top and bottom views. The top-back view exhibited poorer results than the top-front view, likely attributed to challenges associated with face visibility.
2. Including a small number of synthetic training samples with extreme views significantly improved extreme view pose estimation.
3. Stronger rotation data augmentation proved crucial, particularly for views adjacent to extreme viewpoints. This augmentation technique is recommended especially for fisheye ceiling-mounted cameras.
4. The pose estimation performance increased when synthetic data closely resembled the poses observed in the target domain.

The next step would be utilizing the proposed model to pre-annotate a larger dataset of extreme views from sports using a human-in-the-loop approach. This process will enable further investigation into the challenges arising from extreme poses. By delving deeper into these complexities, future research endeavors can enhance the understanding and performance of pose estimation in extreme-view scenarios. Furthermore, the annotated dataset comprising almost 200 images of the bottom view and nearly 100 images of the front view, primarily sourced from sports ac-

tivities, will be made publicly available, contributing to the advancement of the field.

Potential misuse. Among other things, our method improves the pose estimation models in ceiling-mounted and surveillance cameras, and it is important to consider potential privacy implications when coupled with face recognition or action recognition systems. This paper focuses on enhancing pose estimation rather than utilizing privacy-sensitive identification models. Nevertheless, we will restrict the usage of our code in a legal way as other fields could benefit from improved extreme view pose estimation.

References

- [1] Mykhaylo Andriluka, Leonid Pishchulin, Peter Gehler, and Bernt Schiele. 2d human pose estimation: New benchmark and state of the art analysis. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2014. **2**
- [2] Guha Balakrishnan, Amy Zhao, Adrian V. Dalca, Frédéric Durand, and John V. Guttag. Synthesizing images of humans in unseen poses. *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8340–8348, 2018. **3**
- [3] Yanrui Bin, Xuan Cao, Xinya Chen, Yanhao Ge, Ying Tai, Chengjie Wang, Jilin Li, Feiyue Huang, Changxin Gao, and Nong Sang. Adversarial semantic data augmentation for human pose estimation. In *European Conference on Computer Vision*, 2020. **2**
- [4] Xu Chen, Jie Song, and Otmar Hilliges. Unpaired pose guided human image generation. *ArXiv*, abs/1901.02284, 2019. **3**
- [5] Yucheng Chen, Mingyi He, and Yuchao Dai. Nearby-person occlusion data augmentation for human pose estimation with non-extra annotations. *2021 Asia-Pacific Signal and Information Processing Association Annual Summit and Conference (APSIPA ASC)*, pages 282–287, 2021. **2**
- [6] MMPose Contributors. Openmmlab pose estimation toolbox and benchmark. <https://github.com/openmmlab/mmpose>, 2020. **2, 7**
- [7] Yiwei Ding, Wenjin Deng, Yinglin Zheng, Pengfei Liu, Meihong Wang, Xuan Cheng, Jianmin Bao, Dong Chen, and Ming Zeng. I^2 r-net: Intra- and inter-human relation network for multi-person pose estimation, 2022. **2**
- [8] Nicola Garau, Giulia Martinelli, Piotr Bródka, Niccolò Bisagno, and Nicola Conci. Panoptop: a framework for generating viewpoint-invariant human pose estimation datasets. In *2021 IEEE/CVF International Conference on Computer Vision Workshops (ICCVW)*, pages 234–242, 2021. **1, 2, 3**
- [9] Albert Haque, Boya Peng, Zelun Luo, Alexandre Alahi, Serena Yeung, and Li Fei-Fei. Towards viewpoint invariant 3d human pose estimation, 2016. **2**
- [10] Junjie Huang, Zheng Zhu, Feng Guo, and Guan Huang. The devil is in the details: Delving into unbiased data processing for human pose estimation. In *The IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2020. **2**
- [11] Linzhi Huang, Yulong Li, Hongbo Tian, Yue Yang, Xian-gang Li, Weihong Deng, and Jieping Ye. Semi-supervised

- 864 2d human pose estimation driven by position inconsistency
 865 pseudo label correction module. In *Proceedings of the*
 866 *IEEE/CVF Conference on Computer Vision and Pattern*
 867 *Recognition (CVPR)*, pages 693–703, June 2023. 2, 6
- 868 [12] Håkon Hukkelås and Frank Lindseth. Synthesizing anyone,
 869 anywhere, in any pose, 2023. 3
- 870 [13] Catalin Ionescu, Dragos Papava, Vlad Olaru, and Cristian
 871 Sminchisescu. Human3.6m: Large scale datasets and predic-
 872 tive methods for 3d human sensing in natural environments.
 873 *IEEE Transactions on Pattern Analysis and Machine Intelli-
 874 gence*, 36(7):1325–1339, jul 2014. 2
- 875 [14] Sam Johnson and Mark Everingham. Clustered pose and
 876 nonlinear appearance models for human pose estimation. In
 877 *Proceedings of the British Machine Vision Conference*, pages
 878 12.1–12.11. BMVA Press, 2010. doi:10.5244/C.24.12. 2
- 879 [15] Hanbyul Joo, Hao Liu, Lei Tan, Lin Gui, Bart C. Nabbe,
 880 I. Matthews, Takeo Kanade, Shohei Nobuhara, and Yaser
 881 Sheikh. Panoptic studio: A massively multiview system for
 882 social motion capture. *2015 IEEE International Conference on
 883 Computer Vision (ICCV)*, pages 3334–3342, 2015. 2
- 884 [16] Jiefeng Li, Can Wang, Hao Zhu, Yihuan Mao, Hao-Shu
 885 Fang, and Cewu Lu. Crowdpose: Efficient crowded scenes
 886 pose estimation and a new benchmark. *arXiv preprint
 887 arXiv:1812.00324*, 2018. 2
- 888 [17] Tsung-Yi Lin, Michael Maire, Serge J. Belongie, James
 889 Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and
 890 C. Lawrence Zitnick. Microsoft coco: Common objects in
 891 context. In *European Conference on Computer Vision*, 2014.
 892 2, 5, 6
- 893 [18] Evan Ling, De-Kai Huang, and Minhoe Hur. Humans need
 894 not label more humans: Occlusion copy & paste for occluded
 895 human instance segmentation. In *British Machine Vision
 896 Conference*, 2022. 2
- 897 [19] Huajun Liu, Fuqiang Liu, Xinyi Fan, and Dong Huang. Po-
 898 larized self-attention: Towards high-quality pixel-wise re-
 899 gression. *Arxiv Pre-Print arXiv:2107.00782*, 2021. 2
- 900 [20] Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng
 901 Zhang, Stephen Lin, and Baining Guo. Swin transformer:
 902 Hierarchical vision transformer using shifted windows. In
 903 *Proceedings of the IEEE/CVF International Conference on
 904 Computer Vision*, pages 10012–10022, 2021. 2
- 905 [21] Matthew Loper, Naureen Mahmood, Javier Romero, Ger-
 906 ard Pons-Moll, and Michael J. Black. SMPL: A skinned
 907 multi-person linear model. *ACM Trans. Graphics (Proc.
 908 SIGGRAPH Asia)*, 34(6):248:1–248:16, Oct. 2015. 2, 3
- 909 [22] Qianli Ma, Jinlong Yang, Anurag Ranjan, Sergi Pujades,
 910 Gerard Pons-Moll, Siyu Tang, and Michael J. Black. Learn-
 911 ing to dress 3D people in generative clothing. In *Computer
 912 Vision and Pattern Recognition (CVPR)*, June 2020. 3
- 913 [23] Georgios Pavlakos, Vasileios Choutas, Nima Ghorbani,
 914 Timo Bolkart, Ahmed A. A. Osman, Dimitrios Tzionas, and
 915 Michael J. Black. Expressive body capture: 3d hands, face,
 916 and body from a single image. In *Proceedings IEEE Conf.
 917 on Computer Vision and Pattern Recognition (CVPR)*, 2019.
 918 1, 3
- 919 [24] Xi Peng, Zhiqiang Tang, Fei Yang, Rogério Schmidt Feris,
 920 and Dimitris N. Metaxas. Jointly optimize data augmenta-
 921 tion and network training: Adversarial data augmentation
 922 in human pose estimation. *2018 IEEE/CVF Conference
 923 on Computer Vision and Pattern Recognition*, pages 2226–
 924 2234, 2018. 2
- 925 [25] Alina Roitberg, David Schneider, Aulia Djamel, Constantin
 926 Seibold, Simon Reiß, and Rainer Stiefelhagen. Let’s play
 927 for action: Recognizing activities of daily living by learning
 928 from life simulation video games. *2021 IEEE/RSJ Interna-
 929 tional Conference on Intelligent Robots and Systems (IROS)*,
 930 pages 8563–8569, 2021. 2
- 931 [26] Robin Rombach, Andreas Blattmann, Dominik Lorenz,
 932 Patrick Esser, and Björn Ommer. High-resolution image
 933 synthesis with latent diffusion models. In *Proceedings of the
 934 IEEE Conference on Computer Vision and Pattern Recog-
 935 nition (CVPR)*, 2022. 3
- 936 [27] Ke Sun, Bin Xiao, Dong Liu, and Jingdong Wang. Deep
 937 high-resolution representation learning for human pose es-
 938 timation. In *Proceedings of the IEEE conference on computer
 939 vision and pattern recognition*, pages 5693–5703, 2019. 2, 7
- 940 [28] Garvita Tiwari, Dimitrije Antic, Jan Eric Lenssen, Nikolaos
 941 Sarafianos, Tony Tung, and Gerard Pons-Moll. Pose-ndf:
 942 Modeling human pose manifolds with neural distance fields.
 943 In *European Conference on Computer Vision (ECCV)*, Octo-
 944 ber 2022. 3
- 945 [29] Güл Varol, Ivan Laptev, Cordelia Schmid, and Andrew Zis-
 946 serman. Synthetic humans for action recognition from un-
 947 seen viewpoints. *International Journal of Computer Vision*,
 948 129:2264 – 2287, 2019. 3
- 949 [30] Güл Varol, Javier Romero, Xavier Martin, Naureen Mah-
 950 mood, Michael J. Black, Ivan Laptev, and Cordelia Schmid.
 951 Learning from synthetic humans. In *CVPR*, 2017. 1, 2, 3, 5
- 952 [31] Yufei Xu, Jing Zhang, Qiming Zhang, and Dacheng Tao.
 953 ViTPose: Simple vision transformer baselines for human
 954 pose estimation. In *Advances in Neural Information Pro-
 955 cessing Systems*, 2022. 2
- 956 [32] Jingrui Yu, Tobias Scheck, Roman Seidel, Yukti Adya, Di-
 957 pankar Nandi, and Gangolf Hirtz. Human pose estimation
 958 in monocular omnidirectional top-view images. In *Pro-
 959 ceedings of the IEEE/CVF Conference on Computer Vision and
 960 Pattern Recognition (CVPR) Workshops*, pages 6410–6419,
 961 June 2023. 2, 5, 6, 7
- 962 [33] Feng Zhang, Xiatian Zhu, Hanbin Dai, Mao Ye, and Ce Zhu.
 963 Distribution-aware coordinate representation for human pose
 964 estimation. In *Proceedings of the IEEE/CVF Conference
 965 on Computer Vision and Pattern Recognition (CVPR)*, June
 966 2020. 2
- 967 [34] Lvmin Zhang and Maneesh Agrawala. Adding conditional
 968 control to text-to-image diffusion models, 2023. 3
- 969 [35] Song-Hai Zhang, Ruilong Li, Xin Dong, Paul L. Rosin, Zixi
 970 Cai, Han Xi, Dingcheng Yang, Hao-Zhi Huang, and Shi-Min
 971 Hu. Pose2seg: Detection free human instance segmentation,
 972 2019. 2