

Data Augmentation in Human-Centric Vision

Wentao Jiang^{1*}, Yige Zhang¹, Shaozhong Zheng¹, Si Liu¹, Shuicheng Yan²

¹Beihang University, Xueyuan Road No.37, Beijing, 100191, China.

² Skywork AI, 7 Straits View, Singapore, 018936, Singapore.

*Corresponding author(s). E-mail(s): jiangwentao@buaa.edu.cn;

Contributing authors: 19241002@buaa.edu.cn; 19241094@buaa.edu.cn; liusi@buaa.edu.cn;
shuicheng.yan@kunlun-inc.com;

Abstract

This survey presents a comprehensive analysis of data augmentation techniques in human-centric vision tasks, a first of its kind in the field. It delves into a wide range of research areas including person ReID, human parsing, human pose estimation, and pedestrian detection, addressing the significant challenges posed by overfitting and limited training data in these domains. Our work categorizes data augmentation methods into two main types: data generation and data perturbation. Data generation covers techniques like graphic engine-based generation, generative model-based generation, and data recombination, while data perturbation is divided into image-level and human-level perturbations. Each method is tailored to the unique requirements of human-centric tasks, with some applicable across multiple areas. Our contributions include an extensive literature review, providing deep insights into the influence of these augmentation techniques in human-centric vision and highlighting the nuances of each method. We also discuss open issues and future directions, such as the integration of advanced generative models like Latent Diffusion Models, for creating more realistic and diverse training data. This survey not only encapsulates the current state of data augmentation in human-centric vision but also charts a course for future research, aiming to develop more robust, accurate, and efficient human-centric vision systems.

Keywords: Data Augmentation, Human-Centric Vision

1 Introduction

Human-centric perception remains a core focus within the computer vision and machine learning communities, encompassing a wide array of research tasks and applications such as person ReID [1–7], human parsing [8–14], human pose estimation [15–22], and pedestrian detection [23–29]. Despite significant advancements, these algorithms are inherently data-intensive and frequently encounter overfitting issues [30–36], where models excel with training data but falter on unseen test data. This challenge intensifies

when access to large datasets is limited, often due to privacy concerns or the need for labor-intensive and costly human annotation tasks [37].

Addressing this issue, data augmentation emerges as a practical solution, particularly in the context of high collection and annotation costs. While previous surveys [38, 39] have explored data augmentation across various computer vision tasks, they often overlook the unique aspects of human-centric vision tasks. This survey seeks to fill this gap by extensively summarizing data augmentation works specific to human-centric vision

tasks, such as person ReID, human parsing, human pose estimation, and pedestrian detection. These tasks, while distinct, share commonalities in augmentation methods that leverage human-specific features, with some techniques being applicable across multiple human-centric tasks.

We aim to comprehensively survey data augmentation methods in human-centric vision, including several works previously discussed in surveys [38, 39], while primarily focusing on methods not extensively covered before, especially those concerning human body augmentation. Data augmentation techniques in this domain can be broadly categorized into two types: data generation and data perturbation. Data generation involves creating or expanding datasets by adding new examples through various techniques, such as collecting additional samples, applying transformations, or introducing variations to bolster model training and generalization. This category encompasses graphic engine-based generation, generative model-based generation, and data recombination.

Complementing data generation, data perturbation is frequently employed, subdivided into image-level and human-level perturbations. Image-level perturbation involves applying transformations to the entire image, such as rotations, flips, zooms, or adjustments in brightness and contrast, to artificially enhance the training dataset’s diversity. On the other hand, human-level perturbation introduces alterations at the level of individual samples, aiming to increase a model’s adaptability to varied instances by simulating the diversity encountered in real scenes.

Furthermore, data augmentation methods can be differentiated based on their application in specific human-centric tasks, including person ReID, human parsing, human pose estimation, and pedestrian detection. Each of these tasks employs distinct augmentation methods that perform both data perturbation and data generation, which will be elaborated upon in this survey.

Our contribution can be summarized as:

- We are the first to conduct a comprehensive survey of data augmentation methods tailored for human-centric vision tasks, highlighting the unique characteristics of these methods in relation to human-centric tasks.

- We provide a comprehensive literature review on data augmentation methods for human-centric vision tasks, summarizing and categorizing methods from various perspectives. This provides an in-depth understanding of crucial factors influencing augmentation techniques in human-centric vision.
- We present a thorough discussion about open issues and potential future directions based on our investigation. Our comprehensive studies uncover the pros/cons of current methods and bring new observations and insights to the community.

Outline. The remainder of this paper is organized as follows. Section 2 introduces the existing review papers classified by the type of augmentation method. Sections 3 describe augmentation methods categorized by the applied human-centric tasks. Finally, Section 4 concludes the paper and discusses several promising future research directions.

2 Categorized by Data Augmentation Method

The proposed taxonomy classifies data augmentation in human-centric vision into two main branches: data perturbation and data augmentations, as presented in Table 1. The former indicates methods that perturb the original example for data augmentation, while the latter refers to methods that generate training new examples for data augmentation. The specifics of each data augmentation method are thoroughly discussed in subsequent sections.

2.1 Data Perturbation

Data perturbation aims to augment data using the existing original example. It can be classified as image-level data perturbation that applies transformation, erasing, and mixing in the whole image and human-level data perturbation that only alters the human instances.

2.1.1 Image-level Perturbation

Image-level data augmentation method involves applying various transformations to the image to artificially increase the diversity of the training

Categories	Sub-categories	Subsub-categories	Methods
Data perturbation	Image-level perturbation	Global perturbation	Scaling and rotation: (X Peng 2018 [40])
			Style transfer: (CamStyle 2018 [41])(C Michaelis 2019 [42]) (Z Zhong 2018 [43])(Z lin 2021 [44])
			Noise injection: (Wang 2021 [45])
		Region-level perturbation	Information dropping: (Z Zhong 2020 [46])(J Huang 2020 [47]) (W Sun 2020 [48])(Pedhunter 2020 [49])
			Grayscale patch: (Y Gong 2021 [50])
			Patch stylized: (S Cygert 2020 [51])
	Human-level perturbation	Human-level occlusion generation	Keypoint masking: (L Ke 2018 [52])
			Copy-paste for complex scenario synthesis: (Y Bin 2020 [53])
			Nearby-person occlusion: (Y Chen 2021 [54])
		Human body perturbation	Altering pedestrian shapes: (Zh Chen 2019 [55])
			2D human pose transformation: (PoseTrans 2022 [56])
			3D human pose transformation: (Li 2021 [57])(Z Xin 2022 [58]) (L Huang 2022 [59])(PoseGU 2023 [60])
Data generation	Graphic Engine-based		(MixedPeds 2017 [61])(W Chen 2017 [62])(Varol 2018 [63]) (Bo lu 2022 [64])(SynPoses 2022 [65])(J Nilsson 2014 [66]) (D Mehta 2018 [67])
	Generative Model-based		(A Siarohin 2018 [68])(X Zhang 2021 [69])(PAC-GAN 2020 [70]) (S Liu 2020 [71])(FD-GAN 2018 [72])(L Zhang 2021 [73]) (V Uc-Cetina 2023 [74])(Z Yang 2023 [75])(J Liu 2018 [76]) (R Zhi 2021 [77])(D Wu 2018 [78])(Q Wu 2021 [79])
	Data recombination	Image-level recombination	Background replacing: (N McLaughlin 2015 [80])(Dai 2022 [81]) (T Kikuchi 2017 [82])(L Chen 2017 [83])(M Tian 2018 [84])
			Copy-paste: (D Dwibedi 2017 [85])(CL Li 2021 [86]) (Instaboost 2019 [87])(J Deng 2022 [88]) (G Ghiasi 2020 [89])(T Remez 2018 [90])
		Human-level recombination	2D image: (F Chen 2020 [91])(K Han 2023 [92])(X Jia 2022 [93])
			3D pose: (PoseAug 2022 [94])

Table 1 Categorized by data augmentation methods.

dataset. These transformations may include rotations, flips, zooms, or brightness changes in image-level or region-level. By augmenting the dataset with these modified images, the model becomes more robust and better generalizes to variations in the input data, enhancing its performance in tasks related to human-centric vision.

Global Perturbation. Global perturbation methods modify the overall characteristics of an image, introducing variations such as noise addition [45] or style transfer [41], as shown in Figure 1. These techniques aim to change the

global appearance of the image, providing the model with exposure to different visual patterns and scenarios during training.

- **Style Transfer [41–44]:** In the realm of global perturbation methods, one notable technique involves the stylization of training images. This method enhances robustness against various types of corruptions, severities, and across different datasets. It is achieved by blending the content of an image with the style elements of another, leading to a significant improvement in the model’s ability to generalize across diverse

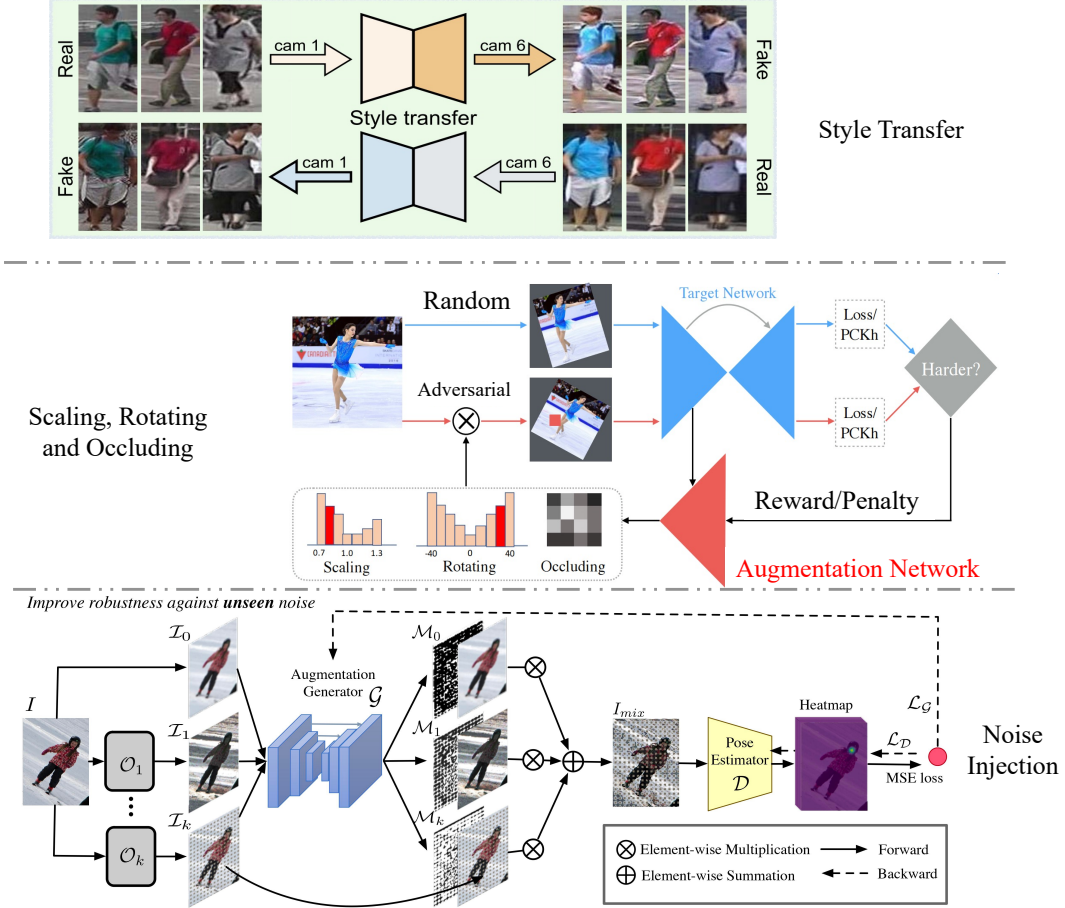


Fig. 1 Examples of global perturbation. The figure contains representative works of Style Transfer [41], Scaling Rotating and Occluding [40] and Noise Injection [45].

visual scenarios. This technique underscores the value of style variability in training data for enhancing model robustness.

- **Scaling, Rotating, and Occluding [40]:** Another approach in global perturbations is the use of an augmentation network designed to create adversarial distributions. From these distributions, specific augmentation operations such as scaling, rotating, and occluding are sampled to generate novel data points. This method represents a strategic shift towards adversarial robustness, where the model is trained on data points that are systematically varied to challenge and thereby improve the model’s resilience and adaptability to real-world variations in visual data.
- **Noise Injection [45]:** AdvMix [95] exemplifies a noise injection approach to improve the robustness of human pose estimation models

against data corruptions. This method, adaptable across various human pose estimation frameworks, employs a combination of adversarial augmentation to introduce challenging corrupted images and knowledge distillation to preserve clean pose information. This strategy effectively trains models to withstand a range of data inconsistencies, enhancing their real-world applicability.

Global perturbation methods in data augmentation present notable advantages and drawbacks. On the positive side, their uniform application across all data instances simplifies operations, contributing to ease of implementation and potentially enhancing the model’s overall robustness. This simplicity is particularly advantageous in

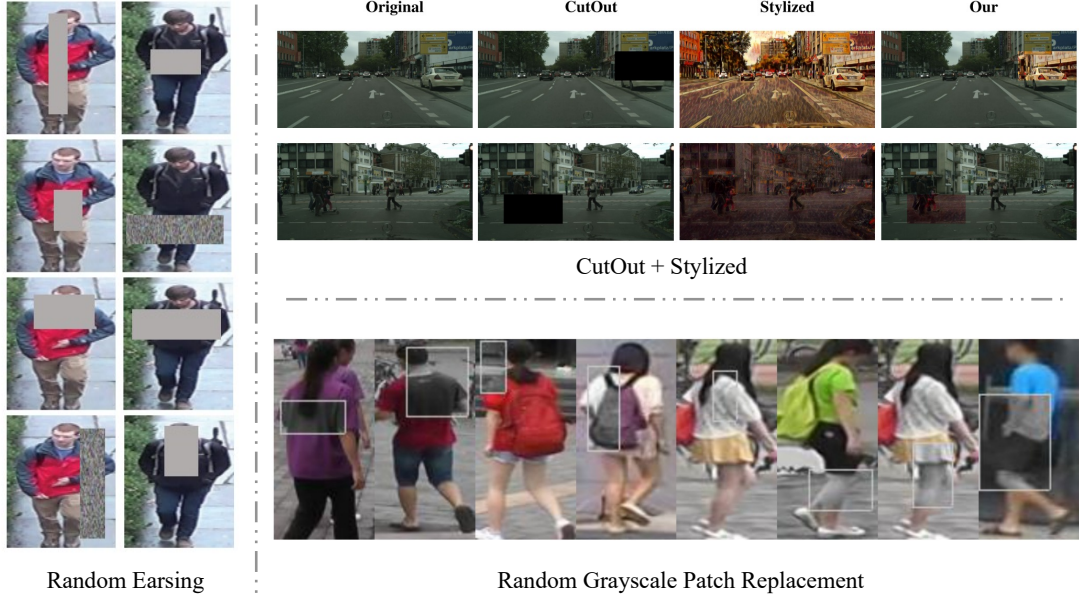


Fig. 2 Examples of region-level perturbation. The figure contains representative works of Random Erasing [46], CutOut and Stylized [51] and Random Grayscale Patch Replacement [50].

ensuring consistent treatment of diverse data, simplifying the training pipeline, and making the model more adaptable to varying scenarios.

However, these advantages come with certain caveats. Excessive modifications through global perturbation may distort original features, impacting the model’s ability to discern meaningful patterns. The risk of overemphasis on artificial variations poses a challenge, potentially diminishing the model’s efficacy in recognizing significant features representative of the natural data distribution. Furthermore, the selection of appropriate perturbation techniques requires careful consideration. Inappropriate transformations may introduce unrealistic scenarios, leading to suboptimal model performance by deviating from real-world data characteristics.

Region-level Perturbation. Augmentation methods, as shown in Figure 2, involving region-level information alteration focus on introducing variability by randomly manipulating specific areas of an image. Techniques like random erasing [46] or adding random grayscale patches [50] selectively modify portions of the image, contributing to a more diverse training set and improving the model’s ability to handle localized variations.

- **Information Dropping:** This is a category of data augmentation techniques that includes methods like random erasing [46], Cutout, and Grid Mask [47]. These methods enhance model training by deliberately removing or obscuring parts of an image. random erasing randomly selects and erases a rectangular region in an image, Cutout removes specific patches, and Grid Mask applies a grid-like pattern to mask parts of the image. All these techniques introduce varying levels of occlusion or information loss, effectively simulating real-world scenarios where parts of a scene may be obscured or missing.
- **Cutout with Stylized Patch Integration [51]:** The cutout method, inspired by the concept of patch Gaussian augmentation, adds a patch of a stylized image to the original image at a randomly sampled location. This technique allows for variation in the size of the patch, introducing a unique form of augmentation. By selectively applying stylized patches, Cutout improves the model’s accuracy and its robustness. This method simulates scenarios where part of the human subject or the background may undergo unusual visual changes, challenging the model to maintain performance under diverse conditions.

- **Random Grayscale Patch Replacement [50]:** This method involves selectively altering regions within an image by replacing them with corresponding grayscale patches. The process, encompassing Local Grayscale Patch Replacement (LGPR) and Global Grayscale Patch Replacement (GGPR), introduces varying levels of grayscale into the training images. By doing so, Random Grayscale Patch Replacement provides a unique form of region-level perturbation that enhances the model’s ability to process images with diverse color profiles. This is especially beneficial in human-centric vision applications where variations in lighting or color can significantly affect image perception and subsequent analysis.

Region-level perturbation methods offer the advantage of simulating real-world scenarios with occlusion and information loss, contributing to a more robust and adaptable model. The simplicity of operations involved in these techniques is another strength, making them easily applicable across diverse datasets. However, these methods exhibit certain drawbacks. While they aim to replicate real-world occlusion and information loss, there is a risk that the introduced occlusion may not authentically represent the complexity of real-world scenarios. This lack of realism in information loss poses challenges, as the model may not effectively learn to handle genuine occlusion situations. Excessive or unrealistic occlusion, stemming from overzealous removal or obscuring of image regions, can impede the model’s ability to recognize crucial features. Therefore, careful management of occlusion levels is imperative to strike a balance between introducing variability and preserving informative content.

In summary, region-level perturbation methods provide a simplified means of introducing variability through simulated occlusion and information loss. Despite their operational simplicity, the challenge lies in ensuring that the introduced occlusion authentically represents real-world complexities, avoiding the risk of hindering the model’s capacity to recognize essential features due to excessive or unrealistic information loss.

2.1.2 Human-level Perturbation

The main feature of human-level data perturbation is to introduce changes at the level of a single

human (instance) to increase the adaptability of the model to different instances. This may include making unique transformations to each instance to simulate the variety in the actual scene.

Human-level Occlusion Generation.

Human-level occlusion generation, as shown in Figure 3, refers to the intentional introduction of obstructions or coverings over specific instances, such as human subjects, within images. This technique is applied during the training phase of machine learning models to improve their ability to handle scenarios where parts of an object or person are obscured or occluded in real-world situations.

- **Keypoint Masking for Occlusion Simulation [52]:** The keypoint masking technique is an innovative approach to simulate occlusion scenarios in training data. It involves two primary methods: the first method covers a body keypoint with a background patch to mimic occlusion, aiding in the model’s occlusion recovery learning. The second method places body keypoint patches onto nearby background areas to simulate multiple keypoints, a scenario commonly encountered in multi-person environments. This augmentation leverages ground-truth keypoint annotations and is instrumental in training models to recognize and interpret occluded keypoints, a critical challenge in human-centric vision tasks like human pose estimation.
- **Copy-Paste for Complex Scenario Synthesis [53]:** Utilizing human parsing to segment training images into various body parts, the Copy-Paste method constructs a semantic part pool based on these segments. This pool allows for the random sampling and placement of body parts onto images, synthesizing complex cases such as symmetric appearances, occlusions, and interactions with nearby individuals. By recombining body parts with different semantic granularities, this method effectively augments training data with realistic and challenging scenarios. This technique is particularly valuable in enhancing the robustness of models to complex, real-world human-centric scenarios, where occlusion and interaction between multiple persons are common.
- **Nearby-Person Occlusion Generation [54]:** This method focuses on creating training images



Fig. 3 Examples of human-level occlusion generation. The figure contains representative works of Keypoint masking [52], Copy-Paste [53] and Nearby-person occlusion [54].

that feature occlusions caused by the proximity of other individuals. Starting with a foreground human body pool generated from rough masks and keypoint annotations, this approach involves randomly sampling and placing a human body crop over another in a training image. This simulates nearby-person occlusion, a frequent occurrence in crowded or group settings. Such augmentation is crucial for training models to accurately detect and interpret human figures in dense, cluttered environments, a common challenge in applications like pedestrian detection and crowd analysis.

Human-level occlusion generation methods navigate real-world scenarios featuring occluded objects or individuals. By intentionally introducing obstructions during the training process, these techniques enrich the dataset with diverse and challenging situations, ultimately contributing to the improvement of model robustness. The simulation of realistic occlusion scenarios enables models to develop a more nuanced understanding of object interactions and improves their

adaptability to complex visual environments. In comparison to image-level perturbation methods, human-level occlusion generation techniques provide a more nuanced and realistic simulation of occlusion scenarios, particularly in the context of human figures. This finer granularity in simulating occlusions contributes to improved model adaptation to real-world scenarios with human occlusion, enhancing overall performance in complex visual environments.

However, the intricacies involved in handling the complex structures and relationships of human body parts may introduce increased computational overhead during the generation of occlusions.

Human body perturbation. Human body perturbation involves introducing small variations or disturbances to the positions of individual joints in a 2D/3D skeletal representation, as shown in Figure 4. These techniques help improve the robustness and generalization capabilities of the model by exposing it to a wider range of variations in pose configurations.

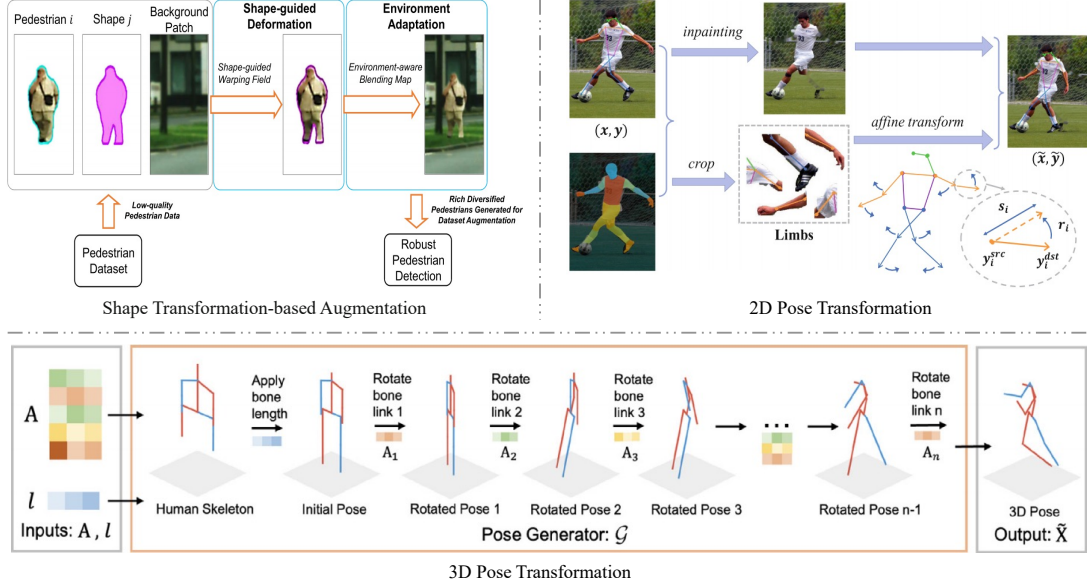


Fig. 4 Examples of human body perturbation. The figure contains representative works of Deformable Shape Augmentation [55], 2D Pose Transformation [56] and 3D Pose Transformation [60].

- **Deformable Shape Augmentation [55]:** This category of data augmentation focuses on altering pedestrian shapes within images. Techniques typically involve warping fields that deform the pedestrian’s shape and blending algorithms to adapt these changes to diverse environmental contexts. Such approaches are invaluable for enhancing the realism and variability of pedestrian datasets, where quality and diversity in real pedestrian images might be limited. By creating a range of shape deformations and integrating them into varying backgrounds, this method effectively prepares models for more accurate pedestrian detection and tracking in dynamically changing real-world scenarios.
- **2D Human Pose Transformation:** Methods in this category, including approaches like PoseTrans [56], utilize various transformation techniques, such as affine transformations, to modify the pose of 2D human figures. These methods often start with limb erasure using human parsing results, followed by selective transformation of each limb. The transformations can include scaling, rotation, and translation, generating a pool of diverse poses. Some techniques may also incorporate Generative Adversarial Networks (GANs) or pose evaluators to ensure the naturalness of generated poses. These methods are crucial for augmenting datasets in 2D

human pose estimation tasks, allowing models to learn from and adapt to a wide array of human poses, improving their accuracy and robustness in real-world applications.

- **3D Human Pose Transformation:** In the realm of 3D human pose estimation, data augmentation methods like PoseGU [60] focus on generating diverse 3D human poses. These techniques typically utilize skeletal models where rotation transformations and bone length adjustments are applied to generate various human poses. The goal is to create a rich dataset from minimal seed samples, significantly enhancing the diversity of poses available for training. This approach is particularly beneficial in scenarios where 3D pose data is scarce or limited to specific types of movements. By providing a broad spectrum of 3D poses, these methods greatly aid in the development of more accurate and versatile 3D human pose estimation models.

Human Body Perturbation methods offer a notable advantage in their ability to intricately simulate the shapes and poses of real-world human bodies, contributing to the generation of more realistic sample data. This precision in mimicking the complexities of human forms is particularly valuable in tasks like pedestrian detection and human pose estimation, where diverse and

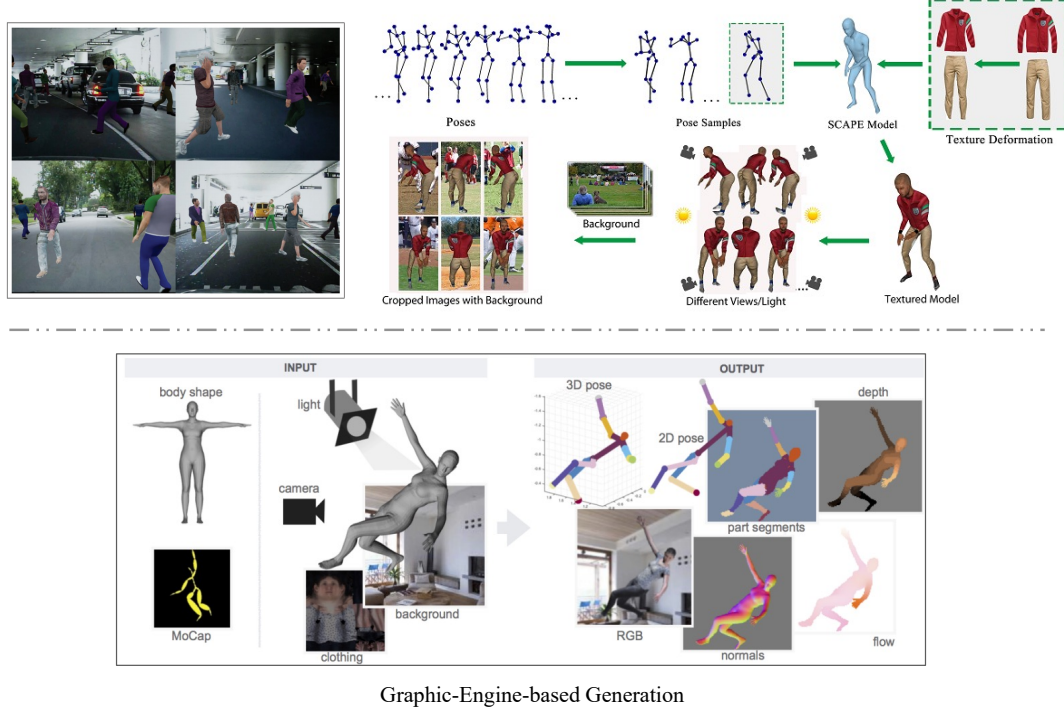


Fig. 5 Examples of graphic engine-based generation. The figure contains representative works of MixedPeds [61], Synthetic Human on Real Background [62] and Synthetic Humans [63].

authentic training samples are essential for model robustness. However, these advantages come with associated drawbacks. The computational complexity linked to techniques such as Deformable Shape Augmentation, especially in the context of 3D human pose estimation, introduces a significant demand for computational resources. Consequently, the implementation of these methods may require substantial computing power and storage capacity, potentially limiting their feasibility in resource-constrained environments. Moreover, ensuring the naturalness and realism of the generated poses remains a critical consideration. In 3D human pose estimation, where the challenge lies in capturing the nuances of three-dimensional movements, the need for natural and believable poses is paramount.

In summary, while Human Body Perturbation methods excel in creating fine-grained simulations of human shapes and poses, their integration comes with computational challenges, particularly in 3D human pose estimation contexts. Striking a balance between computational demands and natural pose generation is crucial for maximizing

the benefits of these methods in enhancing model robustness and generalization capabilities.

2.2 Data Generation

2.2.1 Graphic Engine-based Generation

As shown in Figure 5, synthetic data generation involves creating artificial datasets to augment training sets, employing graphics engine-based methods to produce realistic synthetic examples, and enhancing model robustness and performance.

- MixedPeds [61] automatically extract the vanishing point from the dataset to calibrate the virtual camera and extract the pedestrians' scales to generate a Spawn Probability Map, which guides the algorithm to place the pedestrians at appropriate locations. An HSV color model is used to generate clothing colors according to the color of the pedestrians. Putting synthetic human agents in the unannotated images to use these augmented images to train a Pedestrian Detector.

- **Synthetic Human on Real Background** [62] presents a fully automatic, scalable approach that samples the human pose space for guiding the synthesis procedure. The 3D pose space is sampled and the samples are used for deforming SCAPE models. Meanwhile, various clothes textures are mapped onto the human models. The deformed textured models are rendered using a variety of viewpoints and light sources and finally composited over real image backgrounds.
- **Synthetic Humans** [63] generates RGB images together with 2D/3D poses, surface normals, optical flow, depth images, and body-part segmentation maps for rendered people. A 3D human body model is posed using motion capture data and a frame is rendered using a background image, a texture map of the body, lighting, and a camera position. These ingredients are randomly sampled to increase the diversity of the data.

Graphic engine-based generation in human-centric vision data augmentation starts with extracting key features such as pedestrian scale ratios, vanishing points, or sampling the 3D human pose space. These features guide the creation and customization of human models, often enhanced with varied clothing textures and rendered from multiple viewpoints and lighting conditions. The resulting synthetic humans are strategically placed within real-world scenes, leading to a seamless and realistic blend of virtual and real elements. Such datasets, rich in detail with features like RGB images, 2D/3D poses, and surface normals, are invaluable for a variety of applications, including pedestrian detection and human pose estimation. This approach effectively bridges the gap between simulated and natural environments, significantly enhancing model training and performance in complex, real-world scenarios.

However, the use of Graphic Engine-Based Generation comes with computational demands, especially in tasks like rendering from multiple perspectives and lighting conditions. Additionally, ensuring the naturalness and authenticity of the generated scenes remains a challenge, as overly synthetic or unrealistic elements may impact the model’s ability to generalize effectively. In summary, while Graphic Engine-Based Generation enhances realism and dataset richness, careful consideration is needed to balance computational

demands and maintain the authenticity required for effective model training and real-world application.

2.2.2 Generative Model-based Generation

The generative model-based generation method, as shown in Figure 6, involves creating synthetic data by generating new instances or variations of existing instances within the dataset, providing the model with a richer training set. This kind of method contributes to the model’s adaptability and performance in recognizing and classifying instances under various conditions, ultimately improving its ability to handle complex and diverse datasets. Generative Adversarial Networks are always adopted for data augmentation. GAN-based methods [70–72, 74, 75] for data augmentation in human-centric vision primarily focus on enhancing the diversity and realism of training datasets through the generation of synthetic human poses. These methods leverage pose transfer GANs, often combined with modules for similarity measurement or hard example mining, to create new instances of human images in various poses. By extracting skeletal poses and pairing them with different human appearances, these techniques enable the generation of augmented data that is crucial for tasks such as person re-identification (ReID). The generated images not only enrich the pose variation in the dataset but also improve the model’s ability to recognize and classify human figures across a wide range of conditions. This approach significantly contributes to the adaptability and performance of models in complex and diverse human-centric vision scenarios, making it a valuable tool in the realm of advanced data augmentation.

- **Pose Transferring** [76] proposes a pose-transferrable person ReID framework that utilizes pose-transferred sample augmentations to enhance ReID model training. On one hand, novel training samples with rich pose variations are generated via transferring pose instances from the MARS dataset, and they are added to the target dataset to facilitate robust training. On the other hand, in addition to the conventional discriminator of GAN (i.e., to distinguish between REAL/FAKE samples), a novel

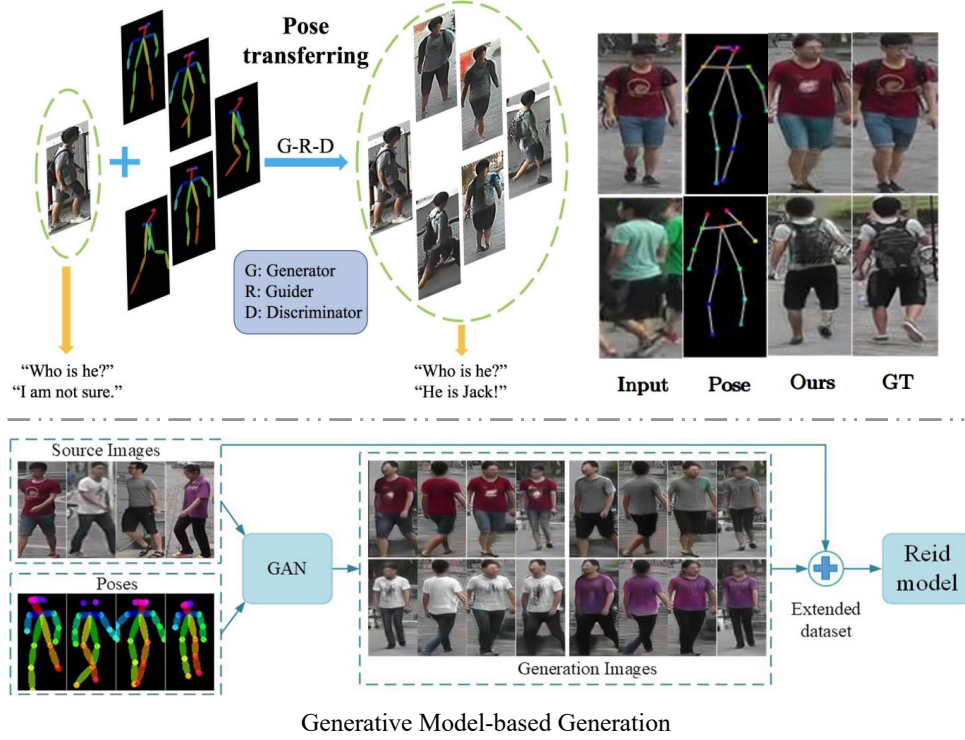


Fig. 6 Examples of generative model-based generation. The figure contains representative works of Pose Transferring [76], Pose Variation Adaptation [73] and Pose Variation Aware Data Augmentation [96]

guider sub-network that encourages the generated sample (i.e., with novel pose) towards better satisfying the ReID loss (i.e., cross-entropy ReID loss, triplet ReID loss).

- Pose Variation Adaptation [73] proposes a pose variation adaptation method for person ReID. It can reduce the probability of deep learning network over-fitting. Specifically, this method introduced a pose transfer generative adversarial network with a similarity measurement module. With the learned pose transfer model, training images can be transferred to any given pose, and with the original images, forming an augmented training dataset.
- Pose Variation Aware Data Augmentation [96] proposes a pose transfer generative adversarial network (PTGAN). PTGAN introduces a similarity measurement module to synthesize realistic person images that are conditional on the pose, and with the original images, form an augmented training dataset.

Generative Adversarial Networks excel in generating realistic and diverse samples. By employing a generator and discriminator in an adversarial training setup, GANs can produce high-quality, visually appealing outputs. GANs are known for their ability to capture complex data distributions and generate images with intricate details. However, GANs have certain drawbacks, including mode collapse, where the generator may focus on a limited subset of modes in the data distribution, and training instability, which can make it challenging to achieve convergence.

Diffusion Models. In the in-depth exploration of the future directions of diffusion models, a particularly noteworthy avenue is the application of advanced generative models, especially pre-trained Latent Diffusion Models, to enhance human-centric vision datasets. Diffusion models stand out with their distinctive paradigm for generative modeling, revolving around the iterative introduction of noise into samples to simulate the generative process of data distribution. The application of diffusion models has demonstrated significant potential in generating high-quality samples

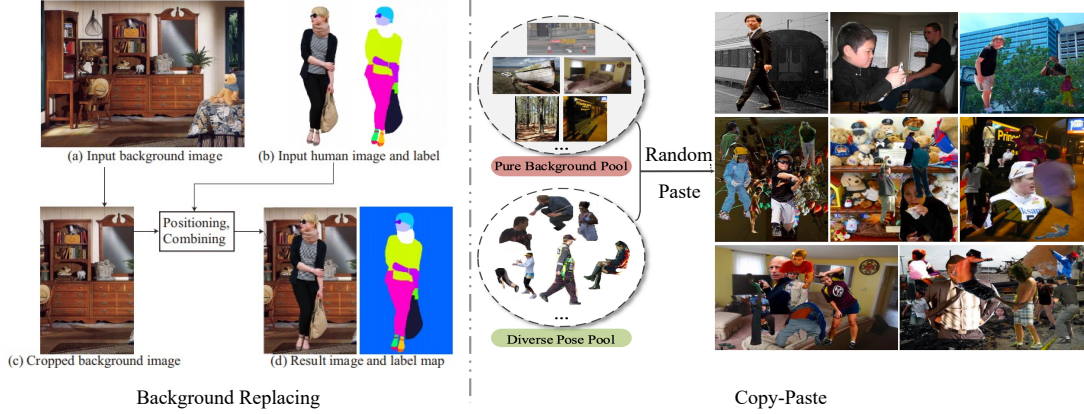


Fig. 7 Examples of image-level recombination. The figure contains representative works of Background Replacing [82] and Copy-Paste [81].

and mitigating issues associated with mode collapse. Research findings suggest that the utilization of diffusion models holds promise for achieving superior results in generating visual data, contributing to the production of more diverse and realistic images for human-centric visual tasks. Notably, the unique aspect of this approach lies in its deviation from the traditional dependence on a discriminator during the training process, potentially simplifying the entire training pipeline. This simplification could offer an effective and feasible pathway for the development and optimization of future generative models. In summary, within the field of generative models, the potential of latent diffusion models in guiding future research directions presents a significant outlook, introducing new possibilities for artificial intelligence applications in visual processing tasks.

2.2.3 Data Recombination

Image-level Recombination. Image-level recombination, as shown in Figure 7, indicates the combination and rearrangement of elements from original images, such as foreground and background recombination, background replacement, or domain mixing.

- **Background Replacing [80–82]:** This method focuses on enhancing the diversity of backgrounds in human-centric images. It involves integrating a human pose estimation network to transfer features across domains and replacing original backgrounds with varied scenes from large-scale scenery datasets. Such augmentation

is essential for improving model performance in environments with complex backgrounds, enhancing the robustness of human parsing networks in varied settings.

- **Copy-paste [85–90]:** Copy-paste augmentation represents a versatile and widely-employed technique in human-centric vision, fundamentally involving the extraction of human figures from one image and their integration into another. This method varies in its application: some approaches, like InstaBoost [87], perform copy-paste within the original image to create variations, while others construct separate pools of pure backgrounds and diverse human figures for recombination. The latter involves randomly selecting human subjects and compositing them into different backgrounds at various locations, enhancing the diversity and complexity of the training data. Some methods also take into account the spatial context and positioning, ensuring that the pasted figures align realistically with the new backgrounds. This technique is particularly effective for training models in tasks such as object detection, person re-identification, and scene understanding, as it introduces a wide range of human appearances and contextual scenarios, enriching the dataset with realistic and challenging examples.

One of the primary advantages of image-level recombination is its efficiency in rapidly generating a large volume of augmented samples, thereby enriching the dataset. This process requires relatively fewer computational resources compared to

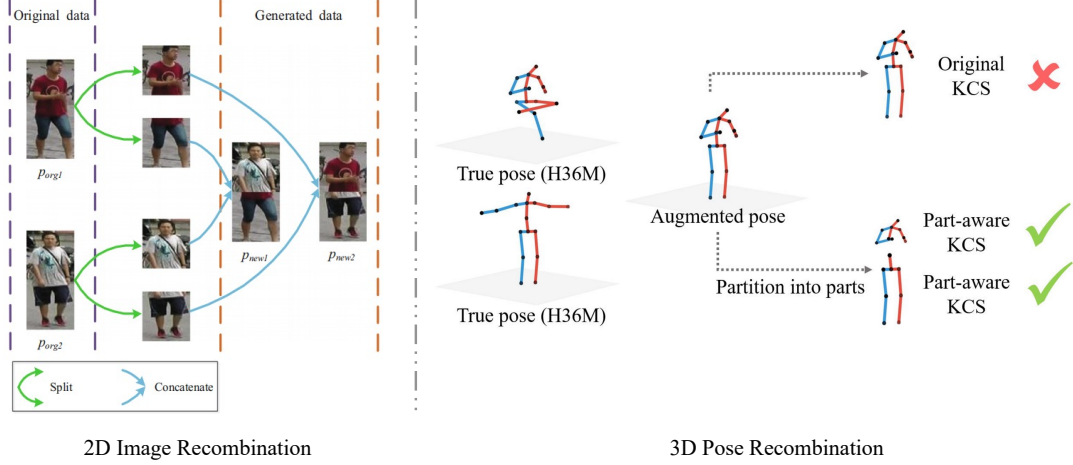


Fig. 8 Examples of human-level recombination. The figure contains representative works of 2D Image Recombination [91] and 3D Image Recombination [94].

certain complex augmentation methods. The ability to create diverse training samples promptly contributes to improved model robustness and generalization. Moreover, the simplicity and speed of image-level recombination make it particularly suitable for scenarios with resource constraints, enabling the augmentation of datasets at scale without excessive computational demands.

However, the generated data may exhibit a trade-off in terms of realism. The artificial nature of the recombined images might introduce a lack of authenticity, potentially leading the model to learn incorrect or unrealistic features. This limitation is crucial, especially when the goal is to train models for real-world scenarios. The challenge lies in striking a balance between the efficiency gained through rapid generation and the fidelity required for effective learning. As such, careful consideration and validation are necessary to ensure that the recombined data aligns with the actual distribution of real-world images, minimizing the risk of model learning spurious patterns that may not generalize well to authentic scenarios.

Human-level Recombination. The human-level recombination method, as shown in Figure 8, involves manipulating image data by selectively extracting and relocating human parts within an image or across images to generate synthetic examples.

- **2D Human Image Recombination:** This technique [91] involves splitting 2D human images

into distinct upper and lower parts and then creatively recombining these segments to generate new synthetic examples. The method effectively increases the diversity of human appearances and postures within the training dataset, introducing variations that are essential for robust model training. By reassembling different combinations of human features, this augmentation strategy enhances the ability of models to accurately recognize and classify individuals or detect pedestrians in varied scenarios.

- **3D Human Pose Recombination:** PoseAug [94] exemplifies a human-level manipulation approach in data augmentation, focusing on enhancing pose diversity for 3D human pose estimation. It introduces a novel pose augmentor capable of adjusting various geometric factors such as posture, body size, viewpoint, and position through differentiable operations. This capability allows for the augmentor to be jointly optimized with the 3D human pose estimator, using estimation errors as feedback to generate more diverse and challenging poses in an online manner. A key feature of PoseAug is its ability to split and recombine the upper and lower parts of the human body in 3D poses, creating new, synthetic examples.

One notable advantage of human-level recombination lies in its capability to introduce nuanced and realistic variations into the dataset. By selectively extracting and relocating human parts, this method can simulate a wide range of scenarios,

Tasks		Categories	Methods
Person ReID	Image-level perturbation		Global perturbation: (Z Zhong 2019 [41])(Z Zhong 2018 [43]) (Z Lin 2021 [44])
			Region-level perturbation: (Y Gong 2021 [50])(W Sun 2020 [48])
	Data recombination		Image-level recombination: (L Chen 2017 [83])(N McLaughlin 2015 [80]) (M Tian 2018 [84])
			Human-level recombination: (K Han 2023 [92])(X Jia 2022 [93]) (F Chen 2020 [91])
	Generative model-based generation		(J Liu 2018 [76])(L Zhang 2021 [73])(L Zhang 2022 [96]) (D Wu 2018 [78])(PAC-GAN 2020 [70])(V Uc-Cetina 2023 [74]) (Z Yang 2023 [75])(Q Wu 2021 [79])
Human Pose Estimation	2D	Image-level perturbation	Global perturbation: (X Peng 2018 [40])(Wang 2021 [45])
			Region-level perturbation: (J Huang 2020 [47])
		Data recombination	Image-level recombination: (Instaboost 2019 [87])(Dai 2022 [81])
			Human-level perturbation
		Human body perturbation: (W Jiang 2022 [56])	
	3D	Human-level perturbation	Human body perturbation: (Li 2021 [57])(Z Xin 2022 [58]) (L Huang 2022 [59])(PoseGU 2023 [60])
		Data recombination	Human-level recombination: (Gong 2022 [94])
Graphic engine-based generation		(Rogez 2016 [97])(Mehta 2017 [67])(W Chen 2017 [62])(Varol 2018 [63])	
Human Parsing		Data recombination	Image-level recombination: (T Kikuchi 2017 [82])(T Remez 2018 [90]) (G Ghiasi 2020 [89])(Instaboost 2019 [87])
Pedestrian Detection	Image-level perturbation		Region-level perturbation: (Z Zhong 2020 [46])(S Cygert 2020 [51]) (Pedhunter 2020 [49])
			Global perturbation: (C Michaelis 2019 [42])
	Data recombination	Image-level recombination: (D Dwibedi 2017 [85])(CL Li 2021 [86]) (J Deng 2022 [88])	
	Human-level perturbation	Human body perturbation: (Z Chen 2019 [55])	
	Graphic engine-based generation	(J Nilsson 2014 [66])(MixedPeds 2017 [61])(SynPoses 2022 [65])	
	Generative model-based generation	(Bo Lu 2022 [64])(R Zhi 2021 [77])(X Zhang 2020 [69])(S Liu 2020 [71])	

Table 2 Categorized by human-centric vision tasks.

such as occlusions, interactions, and variations in body poses. This diversity contributes to the creation of a more comprehensive training dataset, enhancing the model’s adaptability to complex, real-world conditions.

On the other side, the realism of the generated synthetic examples may be subject to the precision and appropriateness of the recombination process. Inaccuracies in part extraction or improper

relocation could result in unrealistic images that deviate from the natural distribution of real-world data. This challenge poses a risk of the model learning patterns or features that may not be representative of authentic scenarios. Consequently, careful validation and refinement of the human-level recombination process are essential to ensure that the generated samples maintain authenticity

and contribute positively to the model’s robustness. Balancing the introduction of diversity with the preservation of realism is crucial for maximizing the effectiveness of human-level recombination in data augmentation.

3 Categorized by Human-Centric Vision Tasks

In Table 2, we classify data augmentation in human-centric vision into two main branches: data perturbation and data augmentations. The former indicates methods that perturb the original example for data augmentation, while the latter refers to methods that generate training new examples for data augmentation. The specifics of each data augmentation method are thoroughly discussed in subsequent sections.

3.1 Person ReID

Person Re-Identification (ReID) in computer vision is a challenging task that involves recognizing and matching individuals across different camera views. This task becomes crucial in surveillance and security applications, where the goal is to track individuals’ movements without compromising personal privacy. Advanced ReID systems utilize deep learning techniques to analyze features like clothing, gait, and even subtle physical characteristics, striving to achieve high accuracy even in crowded or dynamic environments. The main challenge in ReID lies in handling variations in lighting, pose, and occlusion, making robust feature extraction and matching essential for effective identification.

The data augmentation methods in person re-identification (ReID) tasks include image-level perturbation, image/human-level recombination, and generative model-based data generation. Image-level perturbation techniques, such as camera style transfer [41, 43] and grayscale patch augmentation [50], enable ReID models to adapt to different camera styles and environmental conditions by altering color information and encouraging them to focus on structural features. Data recombination augmentation involves background replacing [62, 80, 84] at the image level and human instance recombination at the human level [91],

Methods	mAP (↑)	Rank-1 (↑)
Baselines		
SVDNet [98]	62.1	82.3
Data augmentation methods		
DeformGAN [68]	61.3	80.6
LRSO [99]	66.1	84.0
Random erasing [46]	71.3	87.0
CamStyle [43]	71.6	89.5
PN-GAN [100]	72.6	89.4
FD-GAN [101]	77.7	90.5
DG-Net [102]	86.0	94.8

Table 3 Person Re-identification Performance Comparison of Methods with Data Augmentation on Market1501 Dataset. We compared data augmentation methods based on SVDNet with ResNet-50 as the backbone.

which generates new backgrounds and human appearances, respectively, increasing the diversity of training data. Generative model-based data generation methods employ GANs [6, 68, 70, 74, 75, 79] to synthesize human images with varied clothing and textures while maintaining consistent identity. This significantly enriches the dataset for each identity and improves the ReID model’s ability to generalize across different appearances and poses. These augmentation strategies collectively enhance the robustness and generalizability of ReID systems, enabling them to perform effectively in diverse surveillance and security contexts.

Dataset. In the context of Person Re-identification (ReID) tasks, prominent datasets such as Market-1501 [3], DukeMTMC-reID [103], and CUHK03 [104] are commonly employed. In this study, our primary focus lies on the evaluation of the efficacy of diverse data augmentation methods using the Market-1501 dataset. Market-1501 furnishes an environment simulating real-world surveillance scenarios, thereby endowing algorithms with enhanced robustness for practical applications. The dataset, originating from multiple surveillance cameras, encompasses over 1500 identities and a total of 32,000 images. Notably, it incorporates challenging factors such as disparate backgrounds, varied poses, and occlusions, rendering it an ideal benchmark for assessing the performance of Person ReID algorithms. The chosen dataset, with its realistic surveillance setting and diverse challenges, serves as a pertinent backdrop for evaluating the impact of augmentation strategies on the robustness and generalization capabilities of Person ReID models.

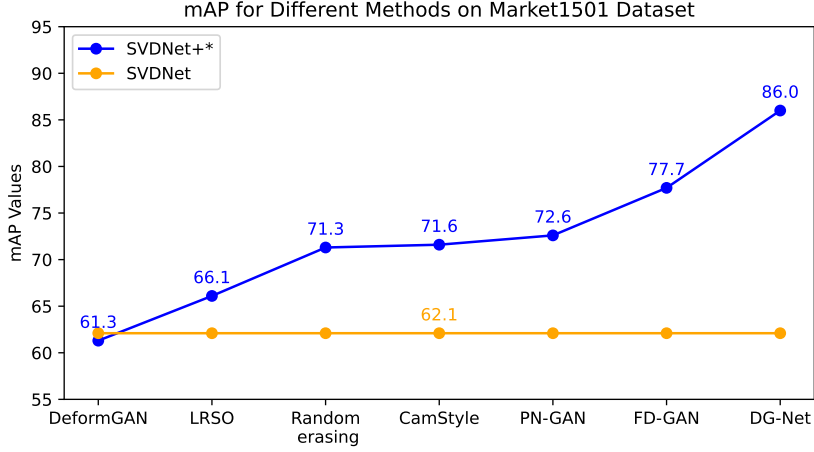


Fig. 9 mAP of different methods on Market1501 Dataset. The orange line represents the baseline SVDNet, while the blue line represents the baseline combined with various data augmentation methods.

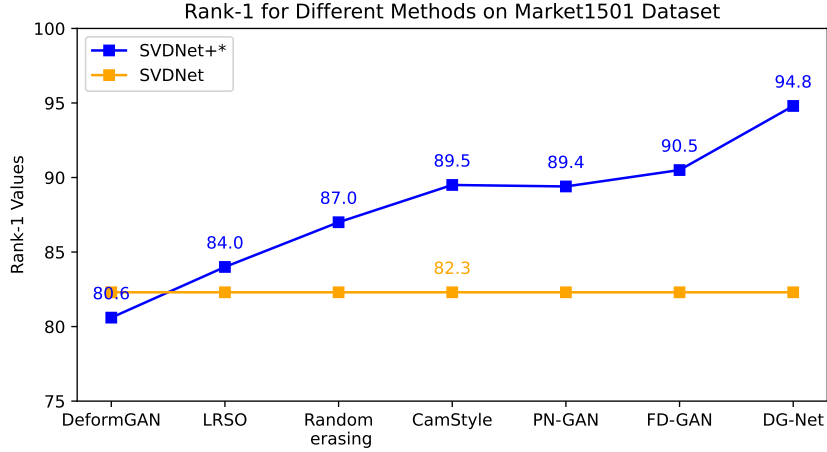


Fig. 10 Most of the data augmentation methods outperformed the baseline in the rank-1 metric. The data augmentation method of generating a human body image is more effective than the image-level perturbation.

Person re-identification experimental analysis. We collect and compare the performance of various methods with data augmentation for person re-identification on the Market1501 dataset using mAP (mean Average Precision) and Rank-1 accuracy as evaluation metrics in Table 3. The data augmentation methods show good performance, as shown in Figure 9 and 10. All methods use SVDNet as the baseline and ResNet-50 as the backbone. To assess the effectiveness of data augmentation methods, we employed SVDNet [98] as the baseline. In the data augmentation methods, based on the evaluation metrics,

DG-Net [102] demonstrates the highest performance among all the methods, followed by FD-GAN [101], PN-GAN [100], and CamStyle [41]. These methods exhibit better performance compared to the other data augmentation methods and the baselines. The results showcase the effectiveness of data augmentation for improving the performance of human-centric vision tasks such as person re-ID.

3.2 Human Pose Estimation

Human pose estimation involves detecting the positions and orientations of human joints in

images or videos, aiming to understand human body language and actions. This technology has far-reaching implications, particularly in sports analytics, physical therapy, and entertainment. By accurately tracking joint positions, human pose estimation algorithms enable the analysis of body movements, offering feedback for performance improvement in athletes or monitoring rehabilitation progress in patients. The major challenge in human pose estimation is achieving accurate joint detection in real-time, especially in complex environments where occlusions or rapid movements occur.

3.2.1 2D Human Pose Estimation

For 2D human pose estimation tasks, the data augmentation methods focus on image-level perturbation and human-level perturbation. Image-level perturbation involves direct modifications to the image. Techniques such as learnable scaling and rotating transformations, as applied in Peng [40], adjust the orientation and size of images to simulate different viewing angles and distances. Wang [45] proposes the injection of noise into images, enhancing the model’s robustness against variations in image quality and real-world disturbances. These perturbations at the image level ensure that the human pose estimation models are not only accurate but also adaptable to a wide range of imaging conditions.

Human-level perturbation introduces more variance specific to humans. Methods like those presented in Ke [52] and Chen [54] create occlusions in keypoint areas, simulating real-world scenarios where humans are partially obscured by objects or other people. Another innovative approach involves cutting and pasting human limbs from different images (as shown in Bin [53]), replicating scenarios of overlapping and interacting human figures. These occlusions and limb manipulations closely mimic complex, crowded environments, providing a more comprehensive training ground for human pose estimation models. An advanced human body perturbation method, as described in Jiang [56], involves the direct transformation of limbs in the original image. This technique modifies the position, rotation, and size of specific body parts, offering a direct and effective means of simulating a wide range of human poses and movements. Such direct

manipulation of the human figure itself presents a unique challenge for human pose estimation models, pushing them to accurately detect joints and limb orientations even under substantial alterations.

Dataset. For the task of 2D human pose estimation, commonly utilized datasets include MS-COCO [111] and MPII Human Pose Dataset [112]. In the ensuing discussion, we aim to assess the effectiveness of various data augmentation methods specifically in the realm of 2D human pose estimation, with a primary focus on the MS-COCO dataset. The MS-COCO dataset comprises over 200,000 images, spanning 80 different object categories, encompassing entities such as humans, animals, vehicles, furniture, and more. With detailed annotations for the task of human pose estimation, each human instance is meticulously labeled, associating keypoint markers. The MS-COCO dataset holds significant value for advancing research in 2D human pose estimation algorithms.

2D human pose estimation experimental analysis. Table 4 shows the comparison of the 2D human pose estimation performance of different methods on the MS-COCO validation set. The evaluation metrics used are AP (average precision), AP50, AP75, and AR (average recall). All methods are compared against the same baseline using HRNet-W32 as the backbone. Various Baseline-based data augmentation methods have seen improvements, as shown in Figure 11. Data-augmented models tend to improve in various ways, such as better detecting human pose in dense crowds, as shown in Figure 12. The baseline method, HRNet-W32 [105], achieves a moderate level of performance across the evaluation metrics. The data augmentation methods, including Cutout [106], GridMask [107], Photometric Distortion [108], AdvMix [45], InstaBoost [87], ASDA [53], and PoseTrans [56], show improvements in performance compared to the baseline method. It appears that some data augmentation methods, such as PoseTrans [56], demonstrate higher performance than others, as indicated by higher AP, AP50, AP75, and AR scores. Overall, Table 4 suggests that employing data augmentation techniques can enhance the performance of 2D human pose estimation models on the MS-COCO dataset. Experimenting with different

Method	AP (\uparrow)	AP50 (\uparrow)	AP75 (\uparrow)	AR (\uparrow)
Baseline				
HRNet-W32 [105]	74.4	90.5	81.9	79.8
Data augmentation methods				
+Cutout* [106]	74.5	90.5	81.7	78.8
+GridMask [107]	74.7	90.6	82.0	80.1
+Photometric Distortion [108]	74.6	90.3	81.9	80.0
+AdvMix [45]	74.7	-	-	-
+InstaBoost [87]	74.7	90.5	82.0	80.1
+ASDA [53]	75.2	91.0	82.4	80.4
+PoseTrans [56]	75.5	91.0	82.9	80.7

Table 4 2D human pose estimation performance comparison of methods with data augmentation on MS-COCO val set. Results marked with ‘*’ are using CascadeRCNN bounding boxes. We compared data augmentation methods based on the HRNet-W32 backbone.

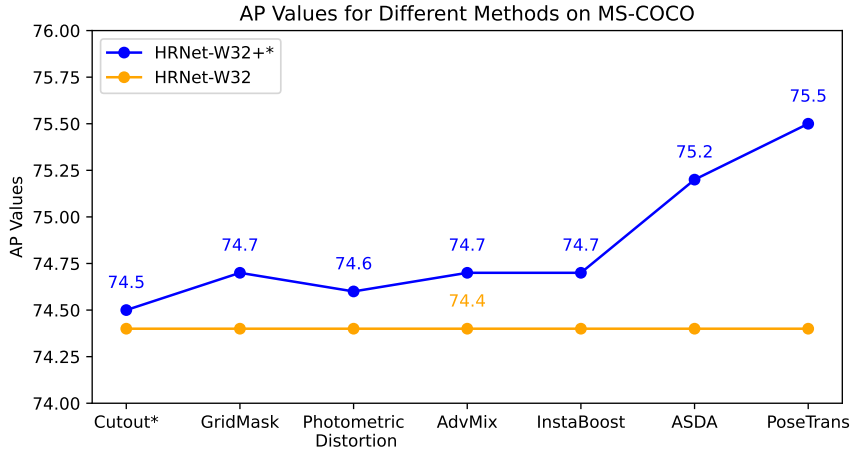


Fig. 11 AP of different methods on MS-COCO Dataset. The orange line represents the baseline HRNet-W32, while the blue line represents the baseline combined with various data augmentation methods. Human-level perturbation methods like PoseTrans can usually get better results.

augmentation methods, such as PoseTrans, may yield better results in terms of accuracy.

3.2.2 3D Human Pose Estimation

In the field of 3D human pose estimation, data augmentation techniques involve various methods, including those generated in graphics engines such as graphic engine-based generation, and human-level perturbation recombination. Compared with 2D human pose estimation, 3D pose ground truth is far more difficult to obtain. Most of the existing 3D human pose estimation datasets are collected indoors, which results in poor generalization in real-world applications. Thus, data augmentation is a very important technique for 3D human pose estimation.

Method	MPJPE (\downarrow)
Non-data augmentation methods	
SemGCN [113]	57.60
Sharma [114]	58.00
Moon [115]	54.40
VPose [116]	52.70
Data augmentation methods	
Li [57]	50.90
VPose + PoseAug [94]	50.20
VPose + DH-AUG [59]	49.81

Table 5 3D human pose estimation performance comparison of methods with Data Augmentation on H36M dataset. For 3D human pose estimation tasks, there is no standardized backbone for data augmentation methods.

Graphic engine-based generation is a specific data augmentation strategy for obtaining 3D



Fig. 12 Comparison of data augmentation method and baseline performance on the CrowdHuman Dataset [109]. Red circles emphasize the inaccurate keypoints predicted by HigherHRNet [110], while green circles demonstrate predictions of data augmentation method Full-DG [81].

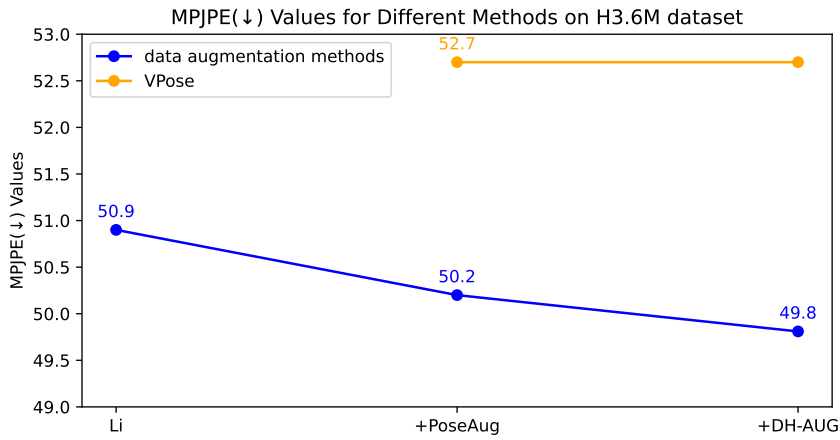


Fig. 13 MPJPE of different methods on H3.6M Dataset. Compared to the baseline, the results are significantly better after using the data augmentation methods.

ground truth. This kind of method employs simulator or rendering methods to synthesize 3D human instances and paste them into the diverse real background. It helps create training examples with almost zero cost and greatly increases the number of training set. Generative model-based generation is also adopted, Rogez et al. [97] introduce an image-based synthesis engine that artificially augments a dataset of real images with 2D human pose annotations using 3D motion capture data.

Apart from generation, human-level perturbation aims to perform the transformation in the limb of 3D pose by rotating and scaling. Human-level recombination for 3D human pose estimation aims to split the 3D human pose into upper/lower parts and then recombine them.

Dataset. For the task of 3D human pose estimation, commonly employed datasets include Human3.6M [124] and MPI-INF-3DHP Dataset [123]. In the subsequent experimental analysis, we assess the impact of data augmentation methods in the domain of 3D human pose estimation on both the H3.6M and 3DHP datasets. Both datasets, H3.6M and 3DHP provide a substantial collection of images captured from diverse environments, accompanied by meticulous 3D keypoint annotations. This furnishes a solid foundation for evaluating the effectiveness of data augmentation methods in the realm of 3D human pose estimation.

3D human pose estimation experimental analysis. Based on Tables 5 and 6, we

Methods	CE	MPJPE (\downarrow)	PCK (\uparrow)	AUC (\uparrow)
Non-data augmentation methods				
Multi Person [117]		122.20	75.20	37.80
OriNet [118]		89.40	81.80	45.20
LCN [119]	✓	-	74.00	36.70
HMR [120]	✓	113.20	77.10	40.70
SRNet [121]	✓	-	77.60	43.80
RepNet [122]	✓	92.50	81.80	54.80
VPose [116]	✓	86.60	-	-
Data augmentation methods				
VNect [67]		124.70	76.60	40.40
Mehta [123]		117.60	76.50	40.80
Li [57]	✓	99.70	81.20	46.10
VPose+PoseAug [94]	✓	73.00	88.60	57.30
VPose+DH-AUG [59]	✓	71.17	89.45	57.93

Table 6 3D human pose estimation performance comparison of methods with Data Augmentation on 3DHP dataset. CE means evaluation across datasets. For 3D human pose estimation tasks, there is no standardized backbone for data augmentation methods.

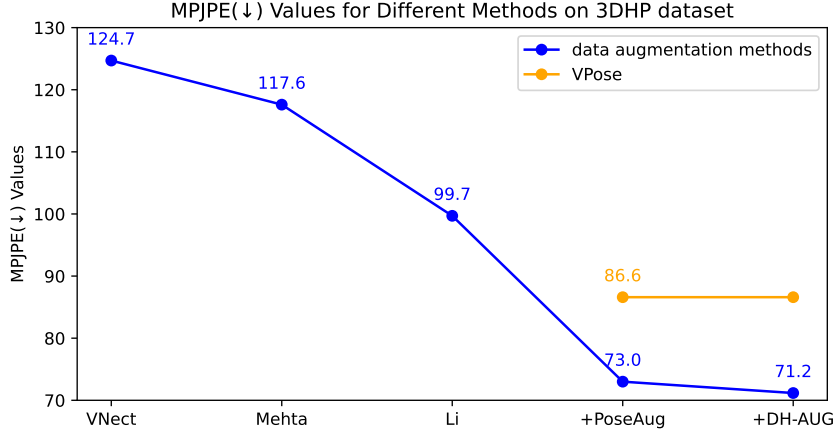


Fig. 14 MPJPE of different methods on 3DHP Dataset. The baseline has achieved good results on this dataset, and data augmentation based on the baseline can also make significant progress.

can analyze the performance of different methods for 3D human pose estimation. As shown in Figure 13 and 14, we can see that with data augmentation methods, we can get better results. Figure 15 shows more accurate 3d human pose estimation with augmented data. The evaluation metrics include MPJPE (Mean Per Joint Position Error), PCK (Percentage of Correct Keypoints), and AUC (Area Under the Curve). In Table 5, we can see that the data augmentation methods, Li, VPose + PoseAug [94], and VPose + DH-AUG [59] show improvements in MPJPE compared to the non-data augmentation

methods. In Table 6, which evaluates methods on the 3DHP dataset, among the data augmentation methods, Li, VPose + PoseAug [94], and VPose + DH-AUG [59] demonstrate improvements in MPJPE compared to the non-data augmentation methods. Based on the provided information, we can conclude that VPose [116] is one of the better-performing methods for 3D human pose estimation in both datasets. VPose + DH-AUG also demonstrates good results, especially on the 3DHP dataset.

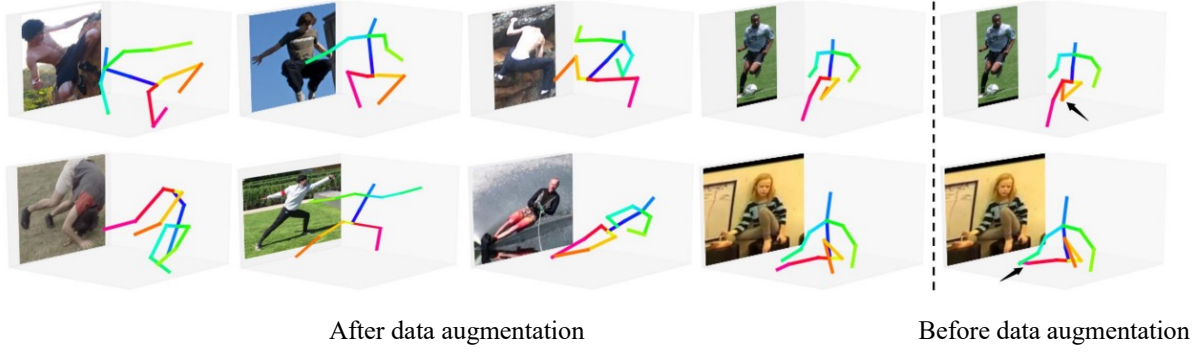


Fig. 15 Example 3D pose estimations from LSP, MPII, 3DHP and 3DPW. Results of PoseAug [94] are shown in the left four columns. The rightmost column shows the results of Baseline—VPoser [116] trained without PoseAug. Errors are highlighted by black arrows.

3.3 Pedestrian Detection

Pedestrian detection is the task of identifying and locating people within images or video frames, predominantly used in autonomous vehicle systems and urban surveillance. This task is crucial for ensuring pedestrian safety, as accurate and rapid detection of pedestrians enables vehicles or monitoring systems to react appropriately to avoid accidents. Modern pedestrian detection systems leverage deep learning models to distinguish pedestrians from various backgrounds and under different lighting conditions. The primary challenge here is to minimize false positives and negatives, ensuring reliable detection in diverse and often unpredictable urban settings.

In the pursuit of enhancing pedestrian detection models, a diverse array of data augmentation methods proves instrumental. Leveraging image-level perturbations, such as style transfer [42], cutout [49, 51], and random erasing, introduces essential variability, allowing models to adapt to different visual conditions. Meanwhile, human-level perturbations focus on altering pedestrian shapes [55]. It can involve techniques such as geometric transformations, shape warping, or body part swapping. By altering the pedestrian’s shape, the model learns to recognize pedestrians across different body proportions, poses, and articulations. Data recombination methods, including copy-paste techniques [85, 86, 88], contribute to dataset diversity by rearranging image components. Furthermore, the integration of virtual data through graphics engine-based generation [61, 64–66] introduces controlled variations, expanding

the model’s exposure to diverse scenes and environmental conditions. In the realm of generative approaches, utilizing generative models like GANs [69, 71, 77] enables the generation of synthetic pedestrian images. This not only broadens the dataset but also empowers models to discern pedestrians across a spectrum of appearances, enhancing their generalization capabilities.

In essence, these augmentation strategies collectively contribute to a holistic training approach, enriching the dataset with realistic variations and bolstering the model’s adaptability to the intricacies of real-world pedestrian scenarios.

Dataset. Commonly used datasets for Pedestrian Detection tasks include CrowdHuman [109] and the INRIA Person Dataset [130]. To assess the effectiveness of various data augmentation methods in the context of Pedestrian Detection, we opt to conduct our evaluation on the CrowdHuman dataset. Comprising approximately 15,000 images, each densely populated with pedestrian instances, the CrowdHuman dataset presents a diverse range of challenges, including different perspectives, various occlusions, and complex scenarios involving intersecting pedestrians. The dataset offers rich annotation information, facilitating the evaluation of algorithmic performance. This choice of dataset allows for a comprehensive validation of data augmentation methods in the domain of Pedestrian Detection.

Pedestrian detection experimental analysis. Table 7 compares the performance of different data augmentation methods on pedestrian detection accuracy using Faster R-CNN [125] and

Methods	MR^{-2} (\downarrow)	AP@0.5 (\uparrow)	AP@0.5:0.95 (\uparrow)	JI (\uparrow)
on Faster R-CNN [125]				
Baseline	50.42	84.95	-	-
Mosaic [108]	43.71	85.21	52.66	78.35
RandAug [126]	42.17	87.48	53.19	80.40
SAutoAug [127]	42.13	87.64	53.35	80.39
SimCP [89]	41.88	87.36	53.36	79.53
CrowdAug [128]	40.21	88.61	54.88	81.41
on RetinaNet [129]				
Baseline	63.33	80.83	-	-
Mosaic [108]	52.53	82.95	48.87	75.60
RandAug [126]	50.25	83.94	49.77	76.58
SAutoAug [127]	50.21	84.02	49.85	76.80
SimCP [89]	50.01	84.12	50.05	77.02
CrowdAug [128]	47.35	85.29	51.84	77.79

Table 7 Pedestrian detection performance comparison of methods with Data Augmentation on CrowdHuman val set. Results are in percentage (%). We compared data augmentation methods based on the Faster R-CNN and RetinaNet backbones.

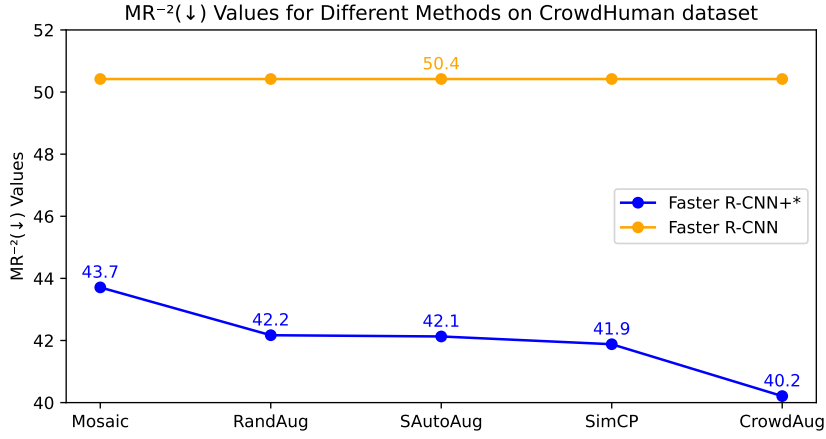


Fig. 16 MR^{-2} of different data augmentation methods with baseline Faster R-CNN on CrowdHuman Dataset. All data augmentation methods show significant performance improvements over baseline.

RetinaNet [129] models on the CrowdHuman validation dataset. Figure 16 and 17 clearly demonstrates the superiority of data augmentation methods over the baseline in detecting pedestrians, as evidenced by their ability to detect pedestrians that were missed by the baseline, as shown in Figure 18. The evaluation metrics used in Table 7 include MR^{-2} (Mean Recall at 2 False Positives per Image), AP@0.5 (Average Precision at 0.5 Intersection over Union threshold), AP@0.5:0.95 (Average Precision at 0.5 to 0.95 Intersection over Union threshold range), and JI (Jaccard Index). For both Faster R-CNN [125] and RetinaNet [129] models, we observe that the baseline model has

the highest MR^{-2} value, indicating a high percentage of recall at 2 false positives per image, but lower AP and JI values. On the other hand, all data augmentation methods improve the AP and JI values, indicating better precision and overlap between predicted and ground truth bounding boxes. Among the data augmentation methods, CrowdAug [128] consistently outperforms other methods in terms of all evaluation metrics for both Faster R-CNN [125] and RetinaNet [129] models. SimCP [89] and SAutoAug [127] follow closely

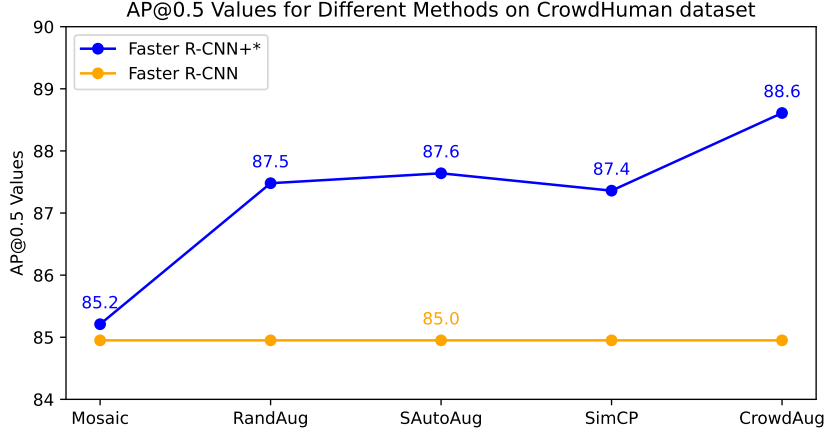


Fig. 17 AP@0.5 of different data augmentation methods with baseline Faster R-CNN on CrowdHuman Dataset.



Fig. 18 Comparison of qualitative results about whether applying STDA [55] for data augmentation. Red boxes are ground-truths. Green dotted boxes are detection results.

behind CrowdAug [128], demonstrating competitive performance across evaluation metrics. RandAug [126] and Mosaic [108] also show improvements in AP and JI values compared to the baseline, although the improvement is relatively small. In conclusion, Table 7 highlights the effectiveness of data augmentation methods in improving the accuracy of pedestrian detection models on the CrowdHuman validation dataset, with CrowdAug [128], SimCP [89], and SAutoAug [127] being the most effective methods.

3.4 Human Parsing

Human parsing refers to the process of segmenting a human image into multiple parts or regions, typically labeling each pixel with a category like head, arms, torso, or legs. This task is vital in applications such as augmented reality, fashion analysis, and advanced human-computer interactions. By understanding the spatial arrangement of different body parts, human parsing algorithms can provide detailed insights into human posture and attire, enabling personalized recommendations in fashion retail or accurate gesture recognition in interactive systems. The complexity of human parsing

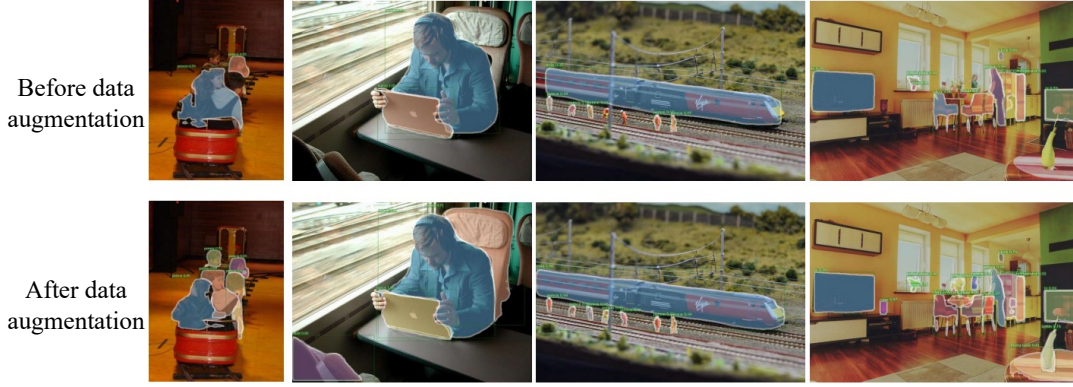


Fig. 19 Instance segmentation result of vanilla Mask R-CNN [131] vs. Mask R-CNN trained with InstaBoost [87] (bottom). InstaBoost guarantees finer instance segmentation results.

arises from the diversity in human poses, clothing styles, and body shapes, requiring sophisticated algorithms that can generalize well across varied scenarios.

In human parsing tasks, data augmentation methods similar to rotation and scale adjustments, which are common in image-level perturbation, can also be applied. However, specific data augmentation methods for human parsing are primarily image-level recombination, which includes image background replacement and copy-paste techniques. These methods are crucial for creating diverse training samples with varied backgrounds and compositions, leading to improved model performance and generalization. Through data augmentation methods, we can segment the people in the image more accurately, as shown in Figure 19. The scarcity of dedicated data augmentation methods for human parsing tasks has resulted in limited experimentation. The absence of standardized datasets across related studies contributes to the challenge of comparing results. The use of different datasets by various research articles introduces inconsistencies and impedes direct result comparisons.

Dataset. The objective of Human Parsing tasks is to segment human body images into distinct semantic parts such as hair, clothing, and skin, thereby facilitating data augmentation. Datasets like the LIP (Look into Person) Dataset [132] and COCO [111] are commonly used for training and evaluating Human Parsing algorithms. The LIP Dataset encompasses over 50,000 human body images, spanning diverse

scenes and poses, and provides detailed annotations covering 50 different body parts. Similarly, the COCO dataset encompasses a multitude of object categories, including the human body, with a substantial number of annotated images suitable for human parsing tasks.

4 Future Work

The future of leveraging data augmentation techniques in human-centric computer vision appears to be a promising avenue to address the challenges posed by overfitting or lack of training data in deep convolutional neural networks. As these networks require extensive datasets to learn effectively, data augmentation can play a critical role in artificially expanding the dataset size and introducing a diverse range of variations in training samples. This process helps in simulating real-world scenarios where changes in lighting, orientation, and background are common. Looking forward, advancements in data augmentation are likely to focus on generating more realistic and complex augmentations that closely mimic real-world variations. This could involve integrating advanced generative models, like Generative Adversarial Networks (GANs) [134, 135], and Latent Diffusion Models (LDMs) [136, 137] to create lifelike, varied training samples. Moreover, the development of task-specific augmentation techniques, specially tailored for specific human-centric tasks such as human pose estimation or person ReID, will be crucial. These specialized augmentations would account for human-specific

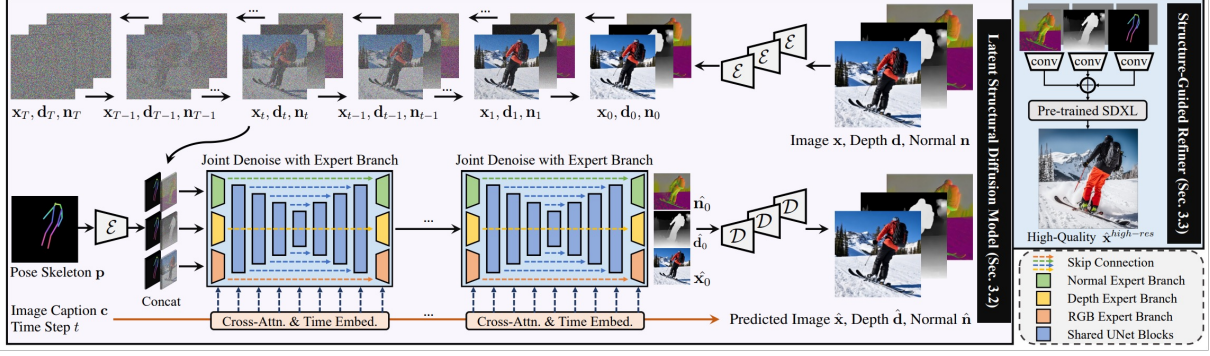


Fig. 20 Visualization of using diffusion model [133] for data generation using pose and depth as conditions.

characteristics and movements, thereby enhancing the robustness and accuracy of these tasks.

4.1 Data Generation with Diffusion Models

Among the above future directions, the most promising one is to leverage the current powerful generative models, especially the pre-trained Latent Diffusion Models (e.g., Stable Diffusion [136]) to generate human data for the human-centric vision for data augmentation. Latent Diffusion Models (LDMs) like Stable Diffusion operate by gradually transforming a random noise distribution into a coherent image representation in a latent space. This process is guided by learned priors, allowing the generation of high-quality, diverse images. In the context of human-centric computer vision, these models can be adapted to generate human images with specific attributes, such as poses or appearances.

Person Re-Identification (ReID): In person ReID, the model must recognize individuals across different scenes and camera angles. A controllable diffusion model can be used to generate images of the same person in various outfits, poses, and lighting conditions, as well as from different camera perspectives. By inputting specific attributes or features (like clothing color, type, or individual physical characteristics) into the model, it can produce a diverse set of images that simulate different scenarios in which a person might be captured by surveillance systems. This helps in training ReID models to be more robust in identifying individuals despite changes in appearance or context.

Pose-Guided Synthesis: Human pose estimation requires accurately identifying the position and orientation of various body parts. The diffusion model can be used to generate images of humans in a multitude of poses, from common to rare or challenging ones, in different environments, as shown in Figure 20. Using models like ControlNet [139] or HyperHuman [133], which demonstrates proficiency in pose-guided image synthesis, we can generate human images in various poses. This is particularly useful for tasks like human pose estimation, where diverse, accurately posed human figures can enrich the training data. By inputting desired pose parameters, the model can synthesize human figures that match these poses, providing a rich and varied dataset for training human pose estimation models. By training on these augmented datasets, human pose estimation models can achieve greater accuracy and flexibility in real-world applications, where human poses can vary greatly.

Human Parsing: For tasks like human parsing, LDMs can be employed to create images based on parsing maps. These maps delineate different human body parts or clothing items, allowing the generation of images with a wide range of appearances and configurations. This approach ensures diversity not only in poses but also in the representation of clothing and body types. The model can be controlled to create complex scenarios where clothing and body parts are partially occluded or overlapping, a common challenge in real-world images. By training on these augmented datasets, human parsing models can learn to more effectively distinguish and segment various parts of the human body under diverse conditions.

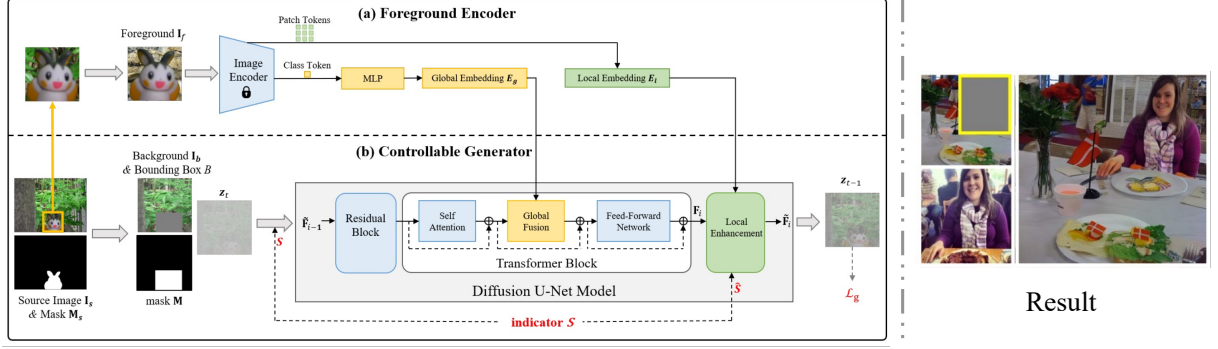


Fig. 21 Visualization of using diffusion model [138] for data recombination.

Pedestrian Detection: In pedestrian detection, the challenge lies in identifying individuals in a variety of urban settings and conditions. The diffusion model can create images of pedestrians in different urban landscapes, under various weather conditions, and during different times of the day. The model can also simulate challenging scenarios, such as crowded scenes or pedestrians partially obscured by objects like vehicles or street furniture. Training on these augmented datasets enables pedestrian detection models to become more adept at recognizing pedestrians in complex and dynamic urban environments.

The generated data can be seamlessly integrated into existing training pipelines. By augmenting the dataset with a wide range of synthetic images, the model’s ability to generalize and perform accurately on real-world data is significantly enhanced. This approach also helps in addressing data scarcity issues, especially in domains where collecting extensive real human datasets is challenging or ethically problematic.

In summary, leveraging pre-trained Latent Diffusion Models for data augmentation in human-centric vision tasks holds immense potential. By leveraging a controllable diffusion model, each of these tasks can benefit from a rich, diverse, and tailored dataset that addresses specific challenges inherent to that task. This approach not only enhances the robustness and accuracy of models in these areas but also significantly contributes to the advancement of human-centric computer vision technologies as a whole.

4.2 Data Perturbation and Recombination with Diffusion Models

Apart from generating the augmented images directly, data perturbation and recombination can leverage the current existing data to create more diverse data. However, the previous methods that perform perturbation and recombination mainly use simple copy-paste and inpainting techniques, which produce artifacts that may harm the model training. With the current powerful generative models like diffusion models, we can enable realistic and plausible data perturbation, and recombination, including background/foreground composition, human switching, and perturbation.

Background/Foreground Composition [140, 141]: The diffusion model, as shown in Figure 21, can be used to seamlessly integrate foreground human subjects into a variety of background scenes. This involves more than just superimposing figures onto backgrounds; it requires an understanding of lighting, perspective, and environmental context to make the composition realistic. In tasks like pedestrian detection or person ReID, this technique helps create scenarios where subjects are placed in diverse environments, under different lighting conditions, and from various camera angles, greatly enhancing the diversity of training data.

Human Switching [142, 143]: This involves replacing one human subject in an image with another while maintaining the integrity and realism of the original scene. The diffusion model can recognize and adapt to the original image’s context, such as lighting, pose, and interaction with

the environment, ensuring a realistic switch. For person ReID, this technique can be particularly useful, as it allows the creation of varied instances of the same individual in different settings or different individuals in the same setting, aiding the model in learning to focus on identifying features of persons rather than the context.

Human Perturbation [144]: Here, the model introduces subtle changes to human subjects or their surroundings. This can include altering aspects like clothing textures, colors, or even background elements. The key is to do this in a way that maintains the overall realism of the image. In human parsing and human pose estimation, such perturbations can create a more robust dataset by introducing minor variations that a model might encounter in real-world scenarios, improving its accuracy and generalization ability.

5 Conclusion

In this survey, we have presented an in-depth analysis of data augmentation techniques in the context of human-centric vision tasks. Our exploration distinctly categorizes these techniques into data generation and perturbation, offering a clear framework for understanding their application in person ReID, human parsing, human pose estimation, and pedestrian detection. This work stands out as the first comprehensive survey specifically addressing data augmentation for human-centric vision, providing a structured overview and critical insights into the nuances of these methods. Future research directions point towards the integration of advanced generative models like Latent Diffusion Models for creating more realistic and diverse training data. This approach shows promise in enhancing model performance across various human-centric tasks by generating tailored, contextually appropriate augmented data. Such advancements are expected to significantly mitigate challenges such as overfitting and data scarcity, marking a substantial step forward in the field. Overall, this survey not only synthesizes the current state of data augmentation in human-centric vision but also paves the way for novel methodologies and applications. The insights provided here will guide future efforts in developing more robust, precise, and efficient human-centric vision systems.

6 Declarations

6.1 Availability of data and material

This work is a comprehensive survey focused on reviewing existing data augmentation methods within human-centric vision research. Given the nature of this study, it primarily involves the analysis and synthesis of findings from previously published papers. As such, this survey does not involve the generation of new datasets or the creation of original materials that would necessitate release or archiving.

6.2 Competing interests

The authors declare the following affiliations: Wentao Jiang, Yige Zhang, Shaozhong Zheng, and Si Liu are affiliated with Beihang University. Shuicheng Yan is affiliated with Skywork AI. Beyond these affiliations, the authors declare no competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

6.3 Funding

This research is supported in part by National Key R&D Program of China (2022ZD0115502), National Natural Science Foundation of China (NO. 62122010, U23B2010), Zhejiang Provincial Natural Science Foundation of China under Grant No. LDT23F02022F02, Key Research and Development Program of Zhejiang Province under Grant 2022C01082.

6.4 Authors' contributions

The survey paper was a collaborative effort, where Wentao Jiang, Yige Zhang, and Shaozhong Zheng, as students, were instrumental in conducting the survey, synthesizing findings, and drafting the manuscript. Professors Si Liu and Shuicheng Yan significantly contributed to guiding the research direction, framework, and providing critical revisions to ensure the manuscript's intellectual integrity. All authors have approved the final manuscript and agreed to be accountable for all aspects of the work.

References

- [1] Ye, M., Shen, J., Lin, G., Xiang, T.,

- Shao, L., Hoi, S.C.: Deep learning for person re-identification: A survey and outlook. *IEEE transactions on pattern analysis and machine intelligence* **44**(6), 2872–2893 (2021)
- [2] Zheng, L., Yang, Y., Hauptmann, A.G.: Person re-identification: Past, present and future. *arXiv preprint arXiv:1610.02984* (2016)
- [3] Zheng, L., Shen, L., Tian, L., Wang, S., Wang, J., Tian, Q.: Scalable person re-identification: A benchmark. In: *Proceedings of the IEEE International Conference on Computer Vision*, pp. 1116–1124 (2015)
- [4] Wang, H., Du, H., Zhao, Y., Yan, J.: A comprehensive overview of person re-identification approaches. *Ieee Access* **8**, 45556–45583 (2020)
- [5] Zheng, L., Zhang, H., Sun, S., Chandraker, M., Yang, Y., Tian, Q.: Person re-identification in the wild. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 1367–1376 (2017)
- [6] Wu, D., Zheng, S.-J., Zhang, X.-P., Yuan, C.-A., Cheng, F., Zhao, Y., Lin, Y.-J., Zhao, Z.-Q., Jiang, Y.-L., Huang, D.-S.: Deep learning-based methods for person re-identification: A comprehensive review. *Neurocomputing* **337**, 354–371 (2019)
- [7] Zhang, X., Luo, H., Fan, X., Xiang, W., Sun, Y., Xiao, Q., Jiang, W., Zhang, C., Sun, J.: Alignedreid: Surpassing human-level performance in person re-identification. *arXiv preprint arXiv:1711.08184* (2017)
- [8] Li, P., Xu, Y., Wei, Y., Yang, Y.: Self-correction for human parsing. *IEEE Transactions on Pattern Analysis and Machine Intelligence* **44**(6), 3260–3271 (2020)
- [9] Ruan, T., Liu, T., Huang, Z., Wei, Y., Wei, S., Zhao, Y.: Devil in the details: Towards accurate single and multiple human parsing. In: *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 33, pp. 4814–4821 (2019)
- [10] Gong, K., Liang, X., Li, Y., Chen, Y., Yang, M., Lin, L.: Instance-level human parsing via part grouping network. In: *Proceedings of the European Conference on Computer Vision (ECCV)*, pp. 770–785 (2018)
- [11] Li, J., Zhao, J., Wei, Y., Lang, C., Li, Y., Sim, T., Yan, S., Feng, J.: Multiple-human parsing in the wild. *arXiv preprint arXiv:1705.07206* (2017)
- [12] Gong, K., Gao, Y., Liang, X., Shen, X., Wang, M., Lin, L.: Graphonomy: Universal human parsing via graph transfer learning. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 7450–7459 (2019)
- [13] Wang, W., Zhu, H., Dai, J., Pang, Y., Shen, J., Shao, L.: Hierarchical human parsing with typed part-relation reasoning. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 8929–8939 (2020)
- [14] Gong, K., Liang, X., Zhang, D., Shen, X., Lin, L.: Look into person: Self-supervised structure-sensitive learning and a new benchmark for human parsing. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 932–940 (2017)
- [15] Fang, H.-S., Xie, S., Tai, Y.-W., Lu, C.: Rmpe: Regional multi-person pose estimation. In: *Proceedings of the IEEE International Conference on Computer Vision*, pp. 2334–2343 (2017)
- [16] Chen, C.-H., Ramanan, D.: 3d human pose estimation= 2d pose estimation+ matching. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 7035–7043 (2017)
- [17] Martinez, J., Hossain, R., Romero, J., Little, J.J.: A simple yet effective baseline for 3d human pose estimation. In: *Proceedings of the IEEE International Conference on Computer Vision*, pp. 2640–2649 (2017)

- [18] Andriluka, M., Iqbal, U., Insafutdinov, E., Pishchulin, L., Milan, A., Gall, J., Schiele, B.: Posetrack: A benchmark for human pose estimation and tracking. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp. 5167–5176 (2018)
- [19] Papandreou, G., Zhu, T., Kanazawa, N., Toshev, A., Tompson, J., Bregler, C., Murphy, K.: Towards accurate multi-person pose estimation in the wild. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp. 4903–4911 (2017)
- [20] Belagiannis, V., Zisserman, A.: Recurrent human pose estimation. In: 2017 12th IEEE International Conference on Automatic Face & Gesture Recognition (FG 2017), pp. 468–475 (2017). IEEE
- [21] Zheng, C., Wu, W., Chen, C., Yang, T., Zhu, S., Shen, J., Kehtarnavaz, N., Shah, M.: Deep learning-based human pose estimation: A survey. *ACM Computing Surveys* **56**(1), 1–37 (2023)
- [22] Zheng, C., Zhu, S., Mendieta, M., Yang, T., Chen, C., Ding, Z.: 3d human pose estimation with spatial and temporal transformers. In: Proceedings of the IEEE/CVF International Conference on Computer Vision, pp. 11656–11665 (2021)
- [23] Liu, W., Liao, S., Ren, W., Hu, W., Yu, Y.: High-level semantic feature detection: A new perspective for pedestrian detection. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 5187–5196 (2019)
- [24] Gawande, U., Hajari, K., Golhar, Y.: Pedestrian detection and tracking in video surveillance system: issues, comprehensive review, and challenges. *Recent Trends in Computational Intelligence*, 1–24 (2020)
- [25] Mao, J., Xiao, T., Jiang, Y., Cao, Z.: What can help pedestrian detection? In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp. 3127–3136 (2017)
- [26] Zhang, S., Benenson, R., Schiele, B.: Citypersons: A diverse dataset for pedestrian detection. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp. 3213–3221 (2017)
- [27] Lan, W., Dang, J., Wang, Y., Wang, S.: Pedestrian detection based on yolo network model. In: 2018 IEEE International Conference on Mechatronics and Automation (ICMA), pp. 1547–1551 (2018). IEEE
- [28] Iftikhar, S., Zhang, Z., Asim, M., Muthanna, A., Koucheryavy, A., Abd El-Latif, A.A.: Deep learning-based pedestrian detection in autonomous vehicles: Substantial issues and challenges. *Electronics* **11**(21), 3551 (2022)
- [29] Roszyk, K., Nowicki, M.R., Skrzypczyński, P.: Adopting the yolov4 architecture for low-latency multispectral pedestrian detection in autonomous driving. *Sensors* **22**(3), 1082 (2022)
- [30] Ying, X.: An overview of overfitting and its solutions. In: *Journal of Physics: Conference Series*, vol. 1168, p. 022022 (2019). IOP Publishing
- [31] Bejani, M.M., Ghatee, M.: A systematic review on overfitting control in shallow and deep neural networks. *Artificial Intelligence Review*, 1–48 (2021)
- [32] Bartlett, P.L., Long, P.M., Lugosi, G., Tsigler, A.: Benign overfitting in linear regression. *Proceedings of the National Academy of Sciences* **117**(48), 30063–30070 (2020)
- [33] Chen, T., Zhang, Z., Liu, S., Chang, S., Wang, Z.: Robust overfitting may be mitigated by properly learned smoothening. In: *International Conference on Learning Representations* (2020)
- [34] Zhang, Z., Dong, M., Ota, K., Zhang, Y., Ren, Y.: Lbcf: A link-based collaborative filtering for overfitting problem in recommender system. *IEEE Transactions on Computational Social Systems* **8**(6), 1450–1464 (2021)

- [35] Hosseini, M., Powell, M., Collins, J., Callahan-Flintoft, C., Jones, W., Bowman, H., Wyble, B.: I tried a bunch of things: The dangers of unexpected overfitting in classification of brain data. *Neuroscience & Biobehavioral Reviews* **119**, 456–467 (2020)
- [36] Zhang, Z.-Y., Sheng, X.-R., Zhang, Y., Jiang, B., Han, S., Deng, H., Zheng, B.: Towards understanding the overfitting phenomenon of deep click-through rate prediction models. *arXiv preprint arXiv:2209.06053* (2022)
- [37] Ashktorab, Z., Desmond, M., Andres, J., Muller, M., Joshi, N.N., Brachman, M., Sharma, A., Brimijoin, K., Pan, Q., Wolf, C.T., *et al.*: Ai-assisted human labeling: Batching for efficiency without overreliance. *Proceedings of the ACM on Human-Computer Interaction* **5**(CSCW1), 1–27 (2021)
- [38] Shorten, C., Khoshgoftaar, T.M.: A survey on image data augmentation for deep learning. *Journal of big data* **6**(1), 1–48 (2019)
- [39] Kumar, T., Turab, M., Raj, K., Mileo, A., Brennan, R., Bendechache, M.: Advanced data augmentation approaches: A comprehensive survey and future directions. *arXiv preprint arXiv:2301.02830* (2023)
- [40] Peng, X., Tang, Z., Yang, F., Feris, R.S., Metaxas, D.: Jointly optimize data augmentation and network training: Adversarial data augmentation in human pose estimation. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 2226–2234 (2018)
- [41] Zhong, Z., Zheng, L., Zheng, Z., Li, S., Yang, Y.: Camstyle: A novel data augmentation method for person re-identification. *IEEE Transactions on Image Processing* **28**(3), 1176–1190 (2018)
- [42] Michaelis, C., Mitzkus, B., Geirhos, R., Rusak, E., Bringmann, O., Ecker, A.S., Bethge, M., Brendel, W.: Benchmarking robustness in object detection: Autonomous driving when winter is coming. *arXiv preprint arXiv:1907.07484* (2019)
- [43] Zhong, Z., Zheng, L., Zheng, Z., Li, S., Yang, Y.: Camera style adaptation for person re-identification. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 5157–5166 (2018)
- [44] Lin, Z., Liu, C., Qi, W., Chan, S.-C.: A color/illuminance aware data augmentation and style adaptation approach to person re-identification. *IEEE Access* **9**, 115826–115838 (2021)
- [45] Wang, J., Jin, S., Liu, W., Liu, W., Qian, C., Luo, P.: When human pose estimation meets robustness: Adversarial algorithms and benchmarks. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 11855–11864 (2021)
- [46] Zhong, Z., Zheng, L., Kang, G., Li, S., Yang, Y.: Random erasing data augmentation. In: *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 34, pp. 13001–13008 (2020)
- [47] Huang, J., Zhu, Z., Huang, G., Du, D.: Aid: Pushing the performance boundary of human pose estimation with information dropping augmentation. *arXiv preprint arXiv:2008.07139* (2020)
- [48] Sun, W., Zhang, X., Zhang, X., Zhang, G., Ge, N.: Triplet erasing-based data augmentation for person re-identification. *International Journal of Sensor Networks* **34**(4), 226–235 (2020)
- [49] Chi, C., Zhang, S., Xing, J., Lei, Z., Li, S.Z., Zou, X.: Pedhunter: Occlusion robust pedestrian detector in crowded scenes. In: *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 34, pp. 10639–10646 (2020)
- [50] Gong, Y., Zeng, Z., Chen, L., Luo, Y.-X., Weng, B., Ye, F.: A person re-identification data augmentation method with adversarial defense effect. *ArXiv* **abs/2101.08783**

- (2021)
- [51] Cygert, S., Czyżewski, A.: Toward robust pedestrian detection with data augmentation. *IEEE Access* **8**, 136674–136683 (2020)
 - [52] Ke, L., Chang, M.-C., Qi, H., Lyu, S.: Multi-scale structure-aware network for human pose estimation. In: *Proceedings of the European Conference on Computer Vision (ECCV)*, pp. 713–728 (2018)
 - [53] Bin, Y., Cao, X., Chen, X., Ge, Y., Tai, Y., Wang, C., Li, J., Huang, F., Gao, C., Sang, N.: Adversarial semantic data augmentation for human pose estimation. In: *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XIX* **16**, pp. 606–622 (2020). Springer
 - [54] Chen, Y., He, M., Dai, Y.: Nearby-person occlusion data augmentation for human pose estimation with non-extra annotations. In: *2021 Asia-Pacific Signal and Information Processing Association Annual Summit and Conference (APSIPA ASC)*, pp. 282–287 (2021). IEEE
 - [55] Chen, Z., Ouyang, W., Liu, T., Tao, D.: A shape transformation-based dataset augmentation framework for pedestrian detection. *International Journal of Computer Vision* **129**(4), 1121–1138 (2021)
 - [56] Jiang, W., Jin, S., Liu, W., Qian, C., Luo, P., Liu, S.: Posetrans: A simple yet effective pose transformation augmentation for human pose estimation. In: *European Conference on Computer Vision*, pp. 643–659 (2022). Springer
 - [57] Li, S., Ke, L., Pratama, K., Tai, Y.-W., Tang, C.-K., Cheng, K.-T.: Cascaded deep monocular 3d human pose estimation with evolutionary training data. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 6173–6183 (2020)
 - [58] Xin, Z., Muqing, W., Min, Z.: 3d pose estimation by synthesizing motions. In: *2022 IEEE 8th International Conference on Computer and Communications (ICCC)*, pp. 2092–2096 (2022). IEEE
 - [59] Huang, L., Liang, J., Deng, W.: Dh-aug: Dh forward kinematics model driven augmentation for 3d human pose estimation. In: *European Conference on Computer Vision*, pp. 436–453 (2022). Springer
 - [60] Guan, S., Lu, H., Zhu, L., Fang, G.: Posegu: 3d human pose estimation with novel human pose generator and unbiased learning. *Computer Vision and Image Understanding* **233**, 103715 (2023)
 - [61] Cheung, E., Wong, A., Bera, A., Manocha, D.: Mixedpeds: Pedestrian detection in unannotated videos using synthetically generated human-agents for training. In: *AAAI Conference on Artificial Intelligence* (2017)
 - [62] Chen, W., Wang, H., Li, Y., Su, H., Wang, Z., Tu, C., Lischinski, D., Cohen-Or, D., Chen, B.: Synthesizing training images for boosting human 3d pose estimation. In: *2016 Fourth International Conference on 3D Vision (3DV)*, pp. 479–488 (2016). IEEE
 - [63] Varol, G., Romero, J., Martin, X., Mahmood, N., Black, M.J., Laptev, I., Schmid, C.: Learning from synthetic humans. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 109–117 (2017)
 - [64] Lu, B., Huang, M., Li, X., Nie, Y., Miao, Q., Lv, Y.: Pedestrian detection for autonomous vehicles using virtual-to-real augmentation. *2022 China Automation Congress (CAC)*, 3652–3657 (2022)
 - [65] Nie, Y., Lu, B., Chen, Q., Miao, Q., Lv, Y.: Synposes: Generating virtual dataset for pedestrian detection in corner cases. *IEEE Journal of Radio Frequency Identification* **6**, 801–804 (2022)
 - [66] Nilsson, J., Andersson, P., Gu, I.Y.-H., Fredriksson, J.: Pedestrian detection using augmented training data. In: *2014 22nd*

- International Conference on Pattern Recognition, pp. 4548–4553 (2014). IEEE
- [67] Mehta, D., Sridhar, S., Sotnychenko, O., Rhodin, H., Shafiei, M., Seidel, H.-P., Xu, W., Casas, D., Theobalt, C.: Vnect: Real-time 3d human pose estimation with a single rgb camera. *Acm transactions on graphics (tog)* **36**(4), 1–14 (2017)
 - [68] Siarohin, A., Sangineto, E., Lathuiliere, S., Sebe, N.: Deformable gans for pose-based human image generation. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 3408–3416 (2018)
 - [69] Zhang, X., Wang, Z., Liu, D., Lin, Q., Ling, Q.: Deep adversarial data augmentation for extremely low data regimes. *IEEE Transactions on Circuits and Systems for Video Technology* **31**(1), 15–28 (2020)
 - [70] Zhang, C., Zhu, L., Zhang, S., Yu, W.: Pac-gan: an effective pose augmentation scheme for unsupervised cross-view person re-identification. *Neurocomputing* **387**, 22–39 (2020)
 - [71] Liu, S., Guo, H., Hu, J.-G., Zhao, X., Zhao, C., Wang, T., Zhu, Y., Wang, J., Tang, M.: A novel data augmentation scheme for pedestrian detection with attribute preserving gan. *Neurocomputing* **401**, 123–132 (2020)
 - [72] Ge, Y., Li, Z., Zhao, H., Yin, G., Yi, S., Wang, X., Li, H.: Fd-gan: Pose-guided feature distilling gan for robust person re-identification. In: *Neural Information Processing Systems* (2018)
 - [73] Zhang, L., Jiang, N., Xu, Y., Diao, Q., Zhou, Z., Wu, W.: Pose variation adaptation for person re-identification. *2020 25th International Conference on Pattern Recognition (ICPR)*, 6996–7003 (2021)
 - [74] Uc-Cetina, V., Alvarez-Gonzalez, L., Martin-Gonzalez, A.: A review on generative adversarial networks for data augmentation in person re-identification systems. *arXiv preprint arXiv:2302.09119* (2023)
 - [75] Yang, Z., Shao, J., Yang, Y.: An improved cyclegan for data augmentation in person re-identification. *Big Data Research* **34**, 100409 (2023)
 - [76] Liu, J., Ni, B., Yan, Y., Zhou, P., Cheng, S., Hu, J.: Pose transferrable person re-identification. *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 4099–4108 (2018)
 - [77] Zhi, R., Guo, Z., Zhang, W., Wang, B., Kaiser, V., Wiederer, J., Flohr, F.B.: Pose-guided person image synthesis for data augmentation in pedestrian detection. *2021 IEEE Intelligent Vehicles Symposium (IV)*, 1493–1500 (2021)
 - [78] Wu, D., Zhang, K., Cheng, F., Zhao, Y., Liu, Q., Yuan, C.-A., Huang, D.-S.: Random occlusion-recovery for person re-identification. *arXiv preprint arXiv:1809.09970* (2018)
 - [79] Wu, Q., Dai, P., Chen, P., Huang, Y.: Deep adversarial data augmentation with attribute guided for person re-identification. *Signal, Image and Video Processing* **15**, 655–662 (2021)
 - [80] McLaughlin, N., Rincón, J.M., Miller, P.C.: Data-augmentation for reducing dataset bias in person re-identification. *2015 12th IEEE International Conference on Advanced Video and Signal Based Surveillance (AVSS)*, 1–6 (2015)
 - [81] Dai, Y., Wang, X., Gao, L., Song, J., Zheng, F., Shen, H.T.: Overcoming data deficiency for multi-person pose estimation. *IEEE Transactions on Neural Networks and Learning Systems* (2023)
 - [82] Kikuchi, T., Endo, Y., Kanamori, Y., Hashimoto, T., Mitani, J.: Transferring pose and augmenting background for deep human-image parsing and its applications. *Computational Visual Media* **4**, 43–54 (2018)

- [83] Chen, L., Yang, H., Wu, S., Gao, Z.: Data generation for improving person re-identification. *Proceedings of the 25th ACM international conference on Multimedia* (2017)
- [84] Tian, M., Yi, S., Li, H., Li, S., Zhang, X., Shi, J., Yan, J., Wang, X.: Eliminating background-bias for robust person re-identification. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 5794–5803 (2018)
- [85] Dwibedi, D., Misra, I., Hebert, M.: Cut, paste and learn: Surprisingly easy synthesis for instance detection. In: *Proceedings of the IEEE International Conference on Computer Vision*, pp. 1301–1310 (2017)
- [86] Li, C.-L., Sohn, K., Yoon, J., Pfister, T.: Cutpaste: Self-supervised learning for anomaly detection and localization. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 9664–9674 (2021)
- [87] Fang, H.-S., Sun, J., Wang, R., Gou, M., Li, Y.-L., Lu, C.: Instaboost: Boosting instance segmentation via probability map guided copy-pasting. In: *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 682–691 (2019)
- [88] Deng, J., Fan, D., Qiu, X., Zhou, F.: Improving crowded object detection via copy-paste. In: *AAAI Conference on Artificial Intelligence* (2022)
- [89] Ghiasi, G., Cui, Y., Srinivas, A., Qian, R., Lin, T.-Y., Cubuk, E.D., Le, Q.V., Zoph, B.: Simple copy-paste is a strong data augmentation method for instance segmentation. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 2918–2928 (2021)
- [90] Remez, T., Huang, J., Brown, M.A.: Learning to segment via cut-and-paste. In: *European Conference on Computer Vision* (2018)
- [91] Chen, F., Wang, N., Tang, J., Liang, D., Feng, H.: Self-supervised data augmentation for person re-identification. *Neurocomputing* **415**, 48–59 (2020)
- [92] Han, K., Gong, S., Huang, Y., Wang, L., Tan, T.: Clothing-change feature augmentation for person re-identification. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 22066–22075 (2023)
- [93] Jia, X., Zhong, X., Ye, M., Liu, W., Huang, W.: Complementary data augmentation for cloth-changing person re-identification. *IEEE Transactions on Image Processing* **31**, 4227–4239 (2022)
- [94] Gong, K., Zhang, J., Feng, J.: Poseaug: A differentiable pose augmentation framework for 3d human pose estimation. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 8575–8584 (2021)
- [95] Huang, Y., Fang, K., Huang, X., Yang, J.: Advmix: Data augmentation for accurate scene text spotting. In: *2021 IEEE International Conference on Image Processing (ICIP)*, pp. 954–958 (2021). IEEE
- [96] Zhang, L., Jiang, N., Diao, Q., Zhou, Z., Wu, W.: Person re-identification with pose variation aware data augmentation. *Neural Computing and Applications* **34**, 11817–11830 (2022)
- [97] Rogez, G., Schmid, C.: Mocap-guided data augmentation for 3d pose estimation in the wild. *Advances in neural information processing systems* **29** (2016)
- [98] Sun, Y., Zheng, L., Deng, W., Wang, S.: Svdnet for pedestrian retrieval. *2017 IEEE International Conference on Computer Vision (ICCV)*, 3820–3828 (2017)
- [99] Zheng, Z., Zheng, L., Yang, Y.: Unlabeled samples generated by gan improve the person re-identification baseline in vitro. In: *Proceedings of the IEEE International Conference on Computer Vision*, pp. 3754–3762 (2017)

- [100] Qian, X., Fu, Y., Xiang, T., Wang, W., Qiu, J., Wu, Y., Jiang, Y.-G., Xue, X.: Pose-normalized image generation for person re-identification. In: Proceedings of the European Conference on Computer Vision (ECCV), pp. 650–667 (2018)
- [101] Ge, Y., Li, Z., Zhao, H., Yin, G., Yi, S., Wang, X., et al.: Fd-gan: Pose-guided feature distilling gan for robust person re-identification. *Advances in neural information processing systems* **31** (2018)
- [102] Zheng, Z., Yang, X., Yu, Z., Zheng, L., Yang, Y., Kautz, J.: Joint discriminative and generative learning for person re-identification. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 2138–2147 (2019)
- [103] Ristani, E., Solera, F., Zou, R., Cucchiara, R., Tomasi, C.: Performance measures and a data set for multi-target, multi-camera tracking. In: European Conference on Computer Vision, pp. 17–35 (2016). Springer
- [104] Li, W., Zhao, R., Xiao, T., Wang, X.: Deepreid: Deep filter pairing neural network for person re-identification. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp. 152–159 (2014)
- [105] Toshev, A., Szegedy, C.: Deeppose: Human pose estimation via deep neural networks. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp. 1653–1660 (2014)
- [106] DeVries, T., Taylor, G.W.: Improved regularization of convolutional neural networks with cutout. *arXiv preprint arXiv:1708.04552* (2017)
- [107] Chen, P., Liu, S., Zhao, H., Jia, J.: Grid-mask data augmentation. *arXiv preprint arXiv:2001.04086* (2020)
- [108] Bochkovskiy, A., Wang, C.-Y., Liao, H.-Y.M.: Yolov4: Optimal speed and accuracy of object detection. *arXiv preprint arXiv:2004.10934* (2020)
- [109] Shao, S., Zhao, Z., Li, B., Xiao, T., Yu, G., Zhang, X., Sun, J.: Crowddhuman: A benchmark for detecting human in a crowd. *arXiv preprint arXiv:1805.00123* (2018)
- [110] Cheng, B., Xiao, B., Wang, J., Shi, H., Huang, T.S., Zhang, L.: Higherhrnet: Scale-aware representation learning for bottom-up human pose estimation. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 5386–5395 (2020)
- [111] Lin, T.-Y., Maire, M., Belongie, S., Hays, J., Perona, P., Ramanan, D., Dollár, P., Zitnick, C.L.: Microsoft coco: Common objects in context. In: Computer Vision—ECCV 2014: 13th European Conference, Zurich, Switzerland, September 6–12, 2014, Proceedings, Part V 13, pp. 740–755 (2014). Springer
- [112] Andriluka, M., Pishchulin, L., Gehler, P., Schiele, B.: 2d human pose estimation: New benchmark and state of the art analysis. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp. 3686–3693 (2014)
- [113] Zhao, L., Peng, X., Tian, Y., Kapadia, M., Metaxas, D.N.: Semantic graph convolutional networks for 3d human pose regression. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 3425–3435 (2019)
- [114] Sharma, S., Varigonda, P.T., Bindal, P., Sharma, A., Jain, A.: Monocular 3d human pose estimation by generation and ordinal ranking. In: Proceedings of the IEEE/CVF International Conference on Computer Vision, pp. 2325–2334 (2019)
- [115] Moon, G., Chang, J.Y., Lee, K.M.: Camera distance-aware top-down approach for 3d multi-person pose estimation from a single rgb image. In: Proceedings of the IEEE/CVF International Conference on Computer Vision, pp. 10133–10142 (2019)
- [116] Pavllo, D., Feichtenhofer, C., Grangier, D., Auli, M.: 3d human pose estimation in

- video with temporal convolutions and semi-supervised training. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 7753–7762 (2019)
- [117] Chu, X., Yang, W., Ouyang, W., Ma, C., Yuille, A.L., Wang, X.: Multi-context attention for human pose estimation. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 1831–1840 (2017)
- [118] Luo, C., Chu, X., Yuille, A.: Orinet: A fully convolutional network for 3d human pose estimation. *arXiv preprint arXiv:1811.04989* (2018)
- [119] Ci, H., Wang, C., Ma, X., Wang, Y.: Optimizing network structure for 3d human pose estimation. In: *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 2262–2271 (2019)
- [120] Kanazawa, A., Black, M.J., Jacobs, D.W., Malik, J.: End-to-end recovery of human shape and pose. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 7122–7131 (2018)
- [121] Zeng, A., Sun, X., Huang, F., Liu, M., Xu, Q., Lin, S.: Srnet: Improving generalization in 3d human pose estimation with a split-and-recombine approach. In: *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XIV 16*, pp. 507–523 (2020). Springer
- [122] Wandt, B., Rosenhahn, B.: Repnet: Weakly supervised training of an adversarial reprojection network for 3d human pose estimation. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 7782–7791 (2019)
- [123] Mehta, D., Rhodin, H., Casas, D., Fua, P., Sotnychenko, O., Xu, W., Theobalt, C.: Monocular 3d human pose estimation in the wild using improved cnn supervision. In: *2017 International Conference on 3D Vision (3DV)*, pp. 506–516 (2017). IEEE
- [124] Ionescu, C., Papava, D., Olaru, V., Sminchisescu, C.: Human3.6m: Large scale datasets and predictive methods for 3d human sensing in natural environments. *IEEE transactions on pattern analysis and machine intelligence* **36**(7), 1325–1339 (2013)
- [125] Ren, S., He, K., Girshick, R.B., Sun, J.: Faster r-cnn: Towards real-time object detection with region proposal networks. *IEEE Transactions on Pattern Analysis and Machine Intelligence* **39**, 1137–1149 (2015)
- [126] Cubuk, E.D., Zoph, B., Mane, D., Vasudevan, V., Le, Q.V.: Autoaugment: Learning augmentation policies from data. *arXiv preprint arXiv:1805.09501* (2018)
- [127] Chen, Y., Li, Y., Kong, T., Qi, L., Chu, R., Li, L., Jia, J.: Scale-aware automatic augmentation for object detection. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 9563–9572 (2021)
- [128] Deng, J., Fan, D., Qiu, X., Zhou, F.: Improving crowded object detection via copy-paste. In: *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 37, pp. 497–505 (2023)
- [129] Lin, T.-Y., Goyal, P., Girshick, R., He, K., Dollár, P.: Focal loss for dense object detection. In: *Proceedings of the IEEE International Conference on Computer Vision*, pp. 2980–2988 (2017)
- [130] Dalal, N., Triggs, B.: Histograms of oriented gradients for human detection. In: *2005 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR’05)*, vol. 1, pp. 886–893 (2005). Ieee
- [131] He, K., Gkioxari, G., Dollár, P., Girshick, R.: Mask r-cnn. In: *Proceedings of the IEEE International Conference on Computer Vision*, pp. 2961–2969 (2017)
- [132] Liang, X., Gong, K., Shen, X., Lin, L.: Look into person: Joint body parsing & pose estimation network and a new benchmark.

- IEEE transactions on pattern analysis and machine intelligence **41**(4), 871–885 (2018)
- [133] Liu, X., Ren, J., Siarohin, A., Skorokhodov, I., Li, Y., Lin, D., Liu, X., Liu, Z., Tulyakov, S.: Hyperhuman: Hyper-realistic human generation with latent structural diffusion. arXiv preprint arXiv:2310.08579 (2023)
- [134] Creswell, A., White, T., Dumoulin, V., Arulkumaran, K., Sengupta, B., Bharath, A.A.: Generative adversarial networks: An overview. IEEE signal processing magazine **35**(1), 53–65 (2018)
- [135] Goodfellow, I., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., Courville, A., Bengio, Y.: Generative adversarial networks. Communications of the ACM **63**(11), 139–144 (2020)
- [136] Rombach, R., Blattmann, A., Lorenz, D., Esser, P., Ommer, B.: High-resolution image synthesis with latent diffusion models. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 10684–10695 (2022)
- [137] Pinaya, W.H., Tudosiu, P.-D., Dafflon, J., Da Costa, P.F., Fernandez, V., Nachev, P., Ourselin, S., Cardoso, M.J.: Brain imaging generation with latent diffusion models. In: MICCAI Workshop on Deep Generative Models, pp. 117–126 (2022). Springer
- [138] Zhang, B., Duan, Y., Lan, J., Hong, Y., Zhu, H., Wang, W., Niu, L.: Controlcom: Controllable image composition using diffusion model. ArXiv **abs/2308.10040** (2023)
- [139] Zhang, L., Rao, A., Agrawala, M.: Adding conditional control to text-to-image diffusion models. In: Proceedings of the IEEE/CVF International Conference on Computer Vision, pp. 3836–3847 (2023)
- [140] Zhang, B., Duan, Y., Lan, J., Hong, Y., Zhu, H., Wang, W., Niu, L.: Controlcom: Controllable image composition using diffusion model. arXiv preprint arXiv:2308.10040 (2023)
- [141] Zhang, X., Zhao, W., Lu, X., Chien, J.: Text2layer: Layered image generation using latent diffusion model. arXiv preprint arXiv:2307.09781 (2023)
- [142] Schmitz, F., Voss, A.: Decomposing task-switching costs with the diffusion model. Journal of Experimental Psychology: Human Perception and Performance **38**(1), 222 (2012)
- [143] Ging-Jehli, N.R., Ratcliff, R.: Effects of aging in a task-switch paradigm with the diffusion decision model. Psychology and aging **35**(6), 850 (2020)
- [144] Zhang, H., Feng, G.: Enhanced example diffusion model via style perturbation. Symmetry **15**(5), 1074 (2023)