

Coarse-to-Fine Cross-modality Generation for Enhancing Vehicle Re-Identification with High-Fidelity Synthetic Data

Leyang Jin¹, Wei Ji¹, Tat-Seng Chua¹, Zhedong Zheng²

Abstract—Due to the critical issues of privacy and partial occlusion, license plate information is not always available in vehicle recognition systems. Consequently, researchers have increasingly turned towards vehicle re-identification (reID) techniques to bridge the gap between cross-view camera systems. Despite the growing interest, one major challenge persists: the scarcity of authentic, large-scale training datasets. To address this challenge, this paper introduces a coarse-to-fine generation pipeline designed to synthesize high-fidelity vehicle data, thereby facilitating subsequent vehicle representation learning. Specifically, the proposed approach consists of three stages: Prompt Processing, Diffusion Fine-tuning, and Semantic Filtering. First, we collect detailed prompts from vehicle websites and companies with fine-grained vehicle prototype attributes. Next, we leverage the prior knowledge of these automotive prototypes to fine-tune diffusion models. Finally, to ensure the quality of the synthesized data, we employ pre-trained vision-language models to filter out substandard images. Building upon the high-quality data generated by this pipeline, we validate the effectiveness using vanilla models. Extensive experimental evaluations demonstrate that our approach achieves competitive accuracy on public benchmarks such as VeRi-776 and CityFlowV2, and is compatible with various model architectures.

I. INTRODUCTION

Vehicle re-identification (reID) aims to match images of a target vehicle across multiple cameras, thus having increasing demands on the deployment of autonomous vehicles [1] and intelligent traffic systems [2]. Considering the minor intra-class difference between different car prototypes, vehicle reID is usually regarded as a fine-grained representation learning task [3], [4]. However, due to the privacy concerns [5] and annotation difficulties in multiple-sensor systems [6], [7], we face the scarcity of realistic data. Therefore, recent researches [8], [9], [10] have resorted to generating authentic data for vehicle reID to break the bottleneck. However, generating large-scale training data for vehicle reID remains challenging, considering high-fidelity images to capture subtle inter-class discrepancy and intra-class consistency. As shown in Table II, directly applying the general generation models even compromises the training and decreases model performance.

Existing efforts on vehicle reID data generation can be divided into two directions: 1) Graphics-engine-based methods, such as PAMTRI [8] and VehicleX [9]. They employ 3D CAD models to generate vehicle images. While these methods have made significant strides, they still face challenges.

¹Leyang Jin, Wei Ji and Tat-seng Chua are with School of Computing, National University of Singapore, Singapore 117417 e0792447@u.nus.edu, {jiwei, dcscts}@nus.edu.sg

²Zhedong Zheng is with the FST and ICI, University of Macau, China 999078 zhedongzheng@um.edu.mo



Fig. 1: We compare our Vehicle-Diff dataset to existing synthetic datasets. The second and third rows of datasets are based on 3D engines (PAMTRI [8] and VehicleX [9]), while PTGAN [11] and VehicleGAN [10] adopt the data-driven structure, *i.e.*, Generative Adversarial Networks [12]. We could observe that the proposed method is with a closer visual appearance compared to the real dataset, *i.e.*, VeRi-776. Besides, the generated images by the proposed method are associated with text captions, allowing for cross-modality knowledge to guide generation.

There is a notable domain gap between rendered 3D CAD vehicle images and actual real-world images. Additionally, the process of generating the VehicleX dataset relies heavily on a large amount of labeled vehicle re-identification data, which is costly and raises privacy concerns. Similarly, synthetic data from PAMTRI needs to be combined with fully labeled re-identification datasets. 2) Data-driven methods, such as generative adversarial networks (GANs) [12]. For instance, PTGAN [11] and VehicleGAN [10] mainly explore GANs to synthesize novel vehicle views. Although these methods can generate vehicle images with relatively good visual quality, they under-explore the cross-modality guidance and thus the fine-grained attributes of the same vehicle are often inconsistent, compromising the training process of the vehicle re-identification task.

To address the aforementioned challenges, we propose Vehicle-Diff, a new pipeline designed to synthesize large-

scale training data for vehicle re-identification, facilitating the representation learning. In particular, the pipeline consists of three primary stages: prompt processing, diffusion model tuning, and semantic filtering. We first collect and process the prompt for vehicles with a focus on the vehicle attribute. To harness the pre-trained inherent knowledge of car prototypes, we employ carefully crafted prompts. Then, we fine-tune the diffusion model using only 1% of unlabeled target data during the generation stage. It enables the diffusion model to adapt to the target vehicle domain at both the content and stylistic levels. In the subsequent filtering stage, we apply sophisticated post-processing techniques to enhance the semantic alignment of the generated data. Our pipeline is scalable and adaptable to multiple downstream scenarios, reducing labeling costs and privacy concerns. As shown in Fig. 1, the generated vehicle images are much closer to the real-world data. Finally, we construct a new labeled vehicle re-identification dataset, called Vehicle-Diff, comprising 149,472 images of 4,940 distinct vehicles. The efficacy of Vehicle-Diff is substantiated through comparative evaluations with synthetic datasets produced by existing approaches. In summary, our paper makes the following contributions:

- A new coarse-to-fine cross-modality generation pipeline by prompting the diffusion model to craft a synthetic vehicle re-identification dataset tailored to a downstream scene, with only about 1% unlabeled images in the original dataset. To the best of our knowledge, our work is among the early attempts for large-scale training data generation with attributes for vehicle re-identification.
- Extensive experiments have validated that our pipeline can minimize the gap between synthetic and real data, facilitating the subsequential reID model learning. The proposed method has achieved competitive performance, *e.g.*, 83.79 mAP on the VeRi-776 dataset.

II. RELATED WORK

Vehicle re-identification. Vehicle re-identification (reID) retrieves vehicles of interest from a database of images collected by traffic cameras. Previous studies [13], [14], [15], [16] rely on supervised learning and have had significant success. However, supervised training based on well-annotated datasets suffers from the high cost of annotation, as well as privacy concerns when collecting and labeling re-identification data. Some studies [17], [18] utilize unsupervised learning to reduce annotation costs. However, a substantial amount of real data is still required for general vehicle reID tasks [16], and attribute annotations [19], [20] are still preferred. In contrast, we developed a multi-modality data synthesis approach that reduces both the need for real data and annotations.

Synthetic datasets for vehicle re-identification task. Synthetic data are frequently utilized to address privacy concerns as well as the high annotation costs associated with creating a re-identification dataset [21], [22]. Some previous works [23], [24], [25] have used 3D engines to create characters for re-identification scenarios. Similar ideas are used in vehicle re-identification [8], [9]. However, assets generated by 3D

engines suffer from the intrinsic domain gap between virtual and real scenes. In addition, hand-crafting 3D assets such as persons and vehicles is time-consuming. Some methods apply GAN [12] for data augmentation. For example, VehicleGAN [10] designs the reconstruction pipeline based on the idea of AutoReconstruction and puts vehicles in the same pose. PTGAN [11] synthesizes novel views of a vehicle based on given pose information. However, they still have the following limitations. To begin, a large labeled dataset is required for effective model training. Second, the quality and patterns of the augmented data are typically constrained by the original dataset.

Text-to-image diffusion models. Diffusion models [26], [27] have recently been regarded as promising generative models. In particular, text-to-image diffusion models are able to generate images following the description of text prompts. Recent text-to-image models such as Stable Diffusion [28], Stable Diffusion XL [29], and Midjourney [30], which are based on diffusion model principles, have achieved astonishing results in text-to-image generation. Based on the great power of text-to-image diffusion models, some methods like [31], [32] utilize diffusion models such as [33] to generate synthetic data for image classification tasks. Despite the exciting visual results and some applications of text-to-image diffusion models, the potential application of these models for the vehicle re-identification task remains unexplored. In this paper, we evaluate multiple state-of-the-art text-to-image models and harness the optimal model for the downstream vehicle re-identification task.

III. METHOD

The overview of Vehicle-Diff is shown in Fig. 2. Vehicle-Diff aims to generate high-fidelity data in a coarse-to-fine manner to boost the training of reID networks, and it contains three stages: (1) prompt processing; (2) diffusion fine-tuning; and (3) semantic filtering. Specifically, the developed prompt processing (§III-A) develops a prompt library and explicitly provides vehicle attributes, *e.g.*, vehicle models and colors, for the subsequent image generation. In the second stage (§III-B), Vehicle-Diff fine-tunes the diffusion model using unlabeled vehicle images, enabling the model to better adapt to vehicle image generation. In the third stage (§III-C), Vehicle-Diff coarsely generates vehicle images with different IDs using a well-developed prompt library and a fine-tuned diffusion model, and then further filters the synthesized images with off-the-shelf detection and cross-modality alignment models.

A. Prompt Processing

The prompt processing stage aims to construct discriminative vehicle attribute prompts to guide image generation, thus enhancing inter-class consistency and intra-class diversity. We first filter the noisy online information to collect vehicle attributes, *i.e.*, brand, production year, and body style, for different car models from an online car information website¹.

¹<https://www.autoevolution.com/>

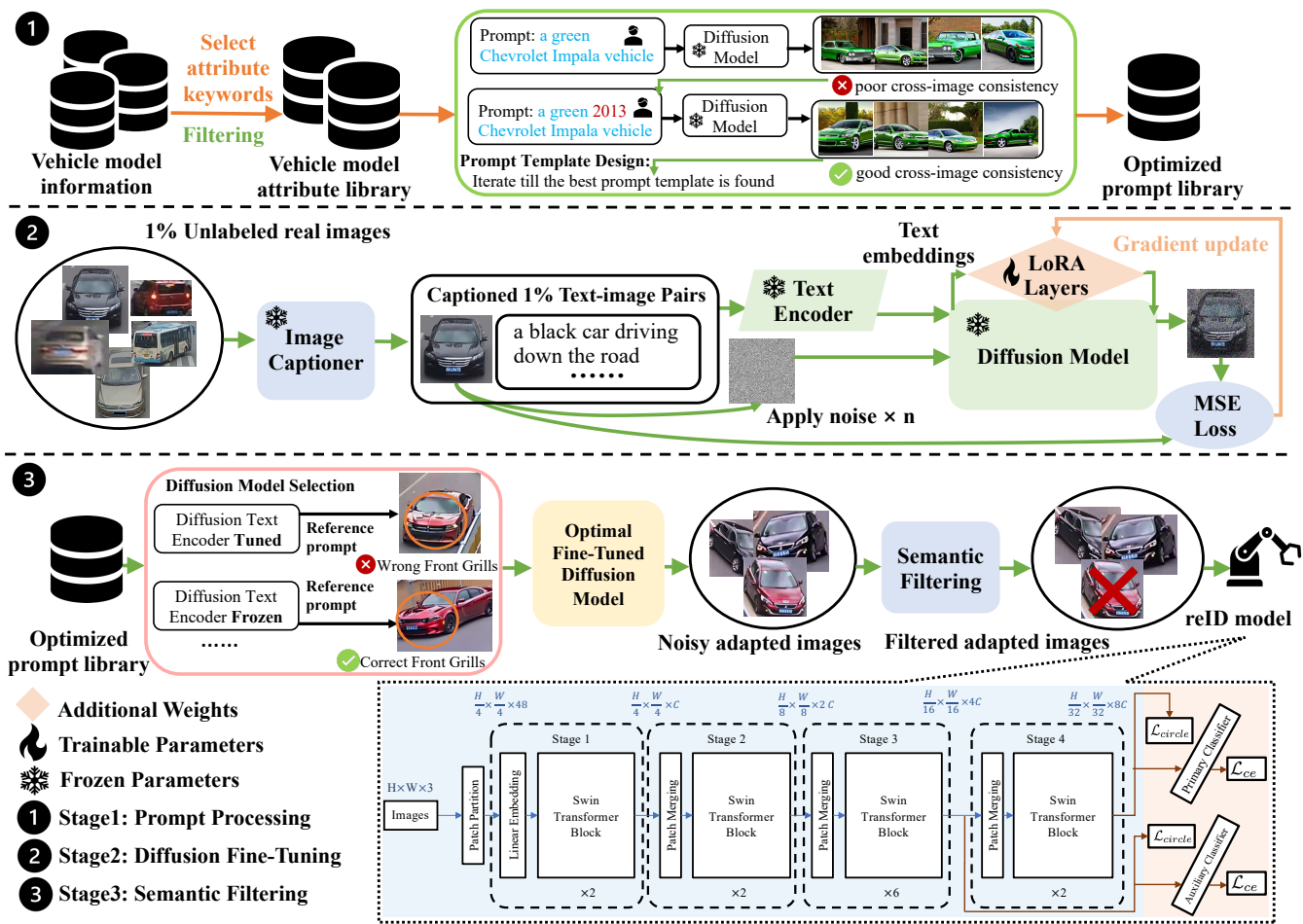


Fig. 2: An overview of our coarse-to-fine cross-modality pipeline Vehicle-Diff. It has three stages: Prompt Processing, Diffusion Fine-tuning, and Semantic Filtering. (1) We first scrape and filter vehicle model information from online vehicle websites. Given the diffusion model, we then select the prompt template according to the visual quality. (2) In the second stage, we leverage the off-the-shelf image captioner to generate the pseudo caption. It is worth noting that the proposed pipeline only requests a few unlabeled real images from the downstream dataset. After the data preparation, we fine-tune the diffusion model via Mean Squared Error (MSE) loss. (3) In the third stage, using the refined prompts, we choose the most effective diffusion model by comparing visual quality, such as consistency. Then, we create synthetic data for the vehicle re-identification task. We use the cross-modality model to filter out semantically misaligned data. Finally, we feed the high-fidelity data to train the reID model via cross-entropy loss [34], [35] and circle loss [36].

It is worth noting that color is an important attribute, and we will use it again in the third stage for semantic filtering. Moreover, inspired by alternating optimization [37] and human-diffusion interaction [38], [39], we also develop a prompt template to improve the quality of the generated images. Specifically, we adjusted one component of the prompt template based on feedback from the diffusion model. The final prompt template is designed as “a [color] [production year] [brand] [car model] [body style] driving down the road.” Please also check the bottom of Fig. 1, where we show several examples of the prompt template, as well as the resulting images.

B. Diffusion Fine-tuning

Vehicle-Diff leverages a text-to-image diffusion model to generate vehicle images according to prompts. How-

ever, a pre-trained diffusion model still struggles to adapt well to the real-world vehicle images, resulting in a domain gap between synthesized images and those in vehicle reID datasets. Therefore, we further fine-tune the diffusion model to mitigate the domain discrepancy while retaining its generation capability. As shown in Fig. 2 (Stage 2), we illustrate the step-by-step fine-tuning stage from the data preparation to the model optimization. To be specific, we first deploy an image captioner, *i.e.*, BLIP-2 [40], to predict text prompts for unlabeled vehicle images, and then employ the generated image-text pairs to fine-tune the text-to-image diffusion model. We incorporate additional weights [41] in the decoder part, while keeping the pre-trained weights unchanged. Therefore, the additional weights could adapt the final visual style, while maintaining the generative capability. The optimization objective is the mean squared error (MSE)

loss. It is worth noting that, our Vehicle-Diff could be trained with **only a few (1%) unlabeled images of the vehicle dataset for fine-tuning**, *i.e.*, 378 images for VeRi-776 and 527 images for CityFlowV2, while previous methods either require large-scale datasets (GAN-based methods [11], [10]) or rely on labeled images (graphics-engine-based methods [8], [9]). Moreover, different from these methods, Vehicle-Diff harnesses the generative power of diffusion models, enabling to generate more realistic images, as shown in Fig. 1. Similarly, we fine-tune multiple candidate diffusion models in preparation for the next stage, which involves selecting the optimal diffusion model.

C. Semantic Filtering

We first sample approximately 10 prompts from the optimized prompt library to evaluate and select the optimal fine-tuned diffusion model. With a similar idea to our prompt template design, the selection of the fine-tuned model is informed by a qualitative assessment of the images generated by each candidate model. Fig. 2 (Stage 3) provides illustrative examples of fine-tuned models evaluated alongside the corresponding generated imagery. Through this evaluation, we opt for the fine-tuned diffusion model that maintains the text encoder in a frozen state. We then feed our designed prompts into the optimal fine-tuned diffusion model, which generates synthetic images automatically. Because of the limitations of text-to-image generation models in producing fine-grained and controllable outputs, directly using generated images is insufficient for training vehicle re-identification networks due to the following two major challenges, *i.e.*, multiple objects and semantic misalignment. We only need portions of the images that include the high-quality vehicle. Diffusion models can generate low-quality images, such as those with multiple vehicles, fragmented vehicles, or no vehicle at all. To address this issue, we employ the off-the-shelf detection model, *i.e.*, YOLOv5x6 [42] trained on high-resolution images of 1280×1280 , for vehicle detection and cropping. We configure the detection model to detect only vehicle categories, specifically focusing on cars and trucks. The number of bounding box per image is limited to one, focusing on the most prominent vehicle in each scene. We only keep the image with high vehicle confidence by setting a threshold, and exclude any small vehicles with heights or widths smaller than or equal to 250 pixels. After cropping, we have the vehicle in the center of the image, and we further screen out noisy images with semantic misalignment, such as vehicles with incorrect colors. In particular, we employ a cross-modal vision-language model, *i.e.*, CLIP [43], to extract the feature for both text and image modalities. We then remove semantic misaligned images that match wrong colors. Specifically, the test prompts are constructed as phrases, *e.g.*, “a red vehicle,” where the color term is dynamically substituted from a predefined color list, such as “red,” “yellow,” “green,” “white,” and “black.” The cosine similarity between image and test text in the feature level is:

$$\text{sim}_k = \frac{\mathbf{f}_I \cdot \mathbf{f}_{T_k}}{\|\mathbf{f}_I\| \|\mathbf{f}_{T_k}\|}. \quad (1)$$

The predicted color \hat{k} is identified as: $\hat{k} = \arg \max_k(\text{sim}_k)$. We then compare the predicted color to the expected color, which is specified within the prompt used to generate the image. If the predicted color matches the expected color, the image is preserved; otherwise, it is discarded.

D. ReID Learning

In this paper, we do not pursue the network structure, but focus on the data aspect. Our generated data is compatible with different networks, and we are free to the reID model selection. Here, we take the typical transformer, Swin-V2 [44], as an example (please see the bottom of Fig. 2). We follow the GoogleNet [45] and existing works [46] to add an auxiliary classifier to facilitate the backward gradients, especially for the large-scale dataset. To optimize the network, we adopt the classification loss [34], [35] and the circle loss [36] as $\mathcal{L}_{total} = \mathcal{L}_{ce} + \mathcal{L}_{circle}$, where \mathcal{L}_{ce} is the cross-entropy loss to classify different vehicles, and the \mathcal{L}_{circle} is to optimize the representation space by pulling closer positive images, while pushing away the negative samples. We apply the same loss terms to both the primary and auxiliary classifiers. It is worth noting that our synthetic data can be combined with real-world data to improve performance even further. In practice, we find that a balance sampler can be useful for reID learning. Experiment contains more details.

IV. EXPERIMENT

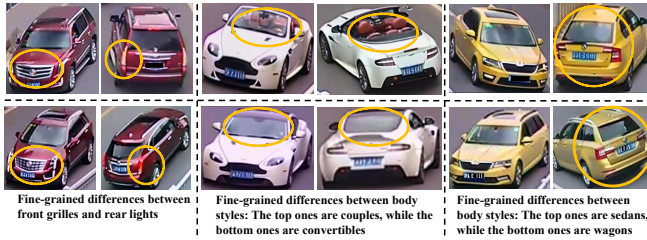
A. Implementation Details

Synthetic data generation. The Diffusion Fine-tuning process uses the Adam optimizer [47], with a learning rate of 0.0001 at the start and a polynomial scheduler for scheduling. We train the diffusion model for 100 epochs, with the first 20 serving as a warm-up. During inference, we set the guidance scale to 8, and the diffusion step to 50. The output size is set to 1024×1024 . The vehicle detection threshold is set to 0.65. Our generation pipeline, Vehicle-Diff, yields 149,472 images of 4,940 vehicles on VeRi-776.

ReID baseline training. We deploy three widely-used baselines to assess the efficacy of our pipeline. Following the setting of existing works [8], [9], we mainly study two CNN-based models, *i.e.*, IDE [35], DenseNet121 [48], and one transformer model, *i.e.*, Swin Transformer V2 [44].

B. Comparison with the State-of-the-art

In Table I, we show the statistics of dataset generated by our Vehicle-Diff and other existing vehicle re-ID datasets. We observe that our pipeline could synthesize more high-fidelity images with more identities, *i.e.*, 4 times larger number of images and IDs compared with VehicleX [9]. It is worth noting that our proposed Vehicle-Diff could further generate more images, if more text prompts are provided. In Table II and Table III, we compare our proposed Vehicle-Diff with existing vehicle re-ID methods on two real-world datasets, *i.e.*, VeRi-776 [55] and CityFlowV2 [57], respectively. For a fair comparison, we follow the setting in the existing work [9] and utilize the same number synthetic image during the reID model training. As shown in Table II,



(a) Inter-class Discrepancy.



(b) Intra-class Variance.

Fig. 3: Our pipeline could reflect the fine-grained discrepancy between two appearance-similar vehicles, *e.g.*, front grilles, rear lights, and body types, while we also depict reasonable intra-class variations of the same vehicle, such as vehicle pose. Please zoom in to get the best view.

	Dataset	#IDs	#Images	#Cameras	Attr
Real	StanfordCars [49]	196	16,185	N/A	✓
	PKU-Vehicle [50]	N/A	10,000,000	N/A	✗
	CompCar [51]	4,701	136,726	N/A	✗
	PKU-VD1 [52]	1,232	1,097,649	1	✓
	PKU-VD2 [52]	1,112	807,260	1	✓
	VehicleID [53]	26,328	222,629	2	✗
	VehicleReID [54]	N/A	47,123	2	✗
	VeRi-776 [55]	776	49,357	20	✓
	CityFlow [56]	666	56,277	40	✗
	CityFlowV2 [57]	440	52,717	46	✗
VRIC [58]	5,622	60,430	60	✗	
Synthetic	PAMTRI [8]	402	41,000	Varied	✓
	VehicleX [9]	1,362	75,516 [†]	Varied	✓
	Vehicle-Diff	4,896	149,472 [‡]	Varied	✓

TABLE I: Comparisons with public real-world and synthetic vehicle re-ID datasets in terms of the number of vehicle IDs, images, and viewpoints, and the availability of attributes. [†]: Number of images used in their code. [‡]: If more text prompts are given, we could generate more images as other synthetic methods.

Query ID	Model	Query	Rank-1	Rank-2	Rank-3	Rank-4	Rank-5
ID 1130	Baseline						
	VehicleX						
	Vehicle-Diff						

Fig. 4: Qualitative retrieval results. Here we compare our method with both our baseline and VehicleX. The ranking list is presented in descending order from left to right based on the similarity score. The images in **red** boxes are false-matched, whereas the **green** ones are true-matched.

Vehicle-Diff enables to achieve competitive vehicle re-ID accuracy on VeRi-776. This indicates that our proposed coarse-to-fine generation pipeline adapts well to vehicle re-ID, and enables to generate high-fidelity training images, even through our generative diffusion model is fine-tuned only with 1% of the unlabeled training data. Specifically, when the reID model is trained solely on synthetic data, our approach improves mAP by 0.92% compared with VehicleX on VeRi-776. When the reID backbone is switched to SwinV2-Base, we observe a consistent mAP improvement,

Method	Backbone	Data	Mix	Rank-1	Rank-5	mAP
VehicleX [9]	Res50	S	-	51.25	67.70	21.29
Vehicle-Diff	Res50	S	-	57.87	74.97	22.21
VehicleX [9]	SwinV2-B	S	-	66.87	79.80	28.33
Vehicle-Diff	SwinV2-B	S	-	74.14	84.45	34.73
VANet [59]	Res50	R	-	89.78	95.99	66.34
AAVER [60]	Res101	R	-	90.17	94.34	66.35
baseline (IDE [35])	Res50	R	-	92.73	96.78	66.54
VehicleX [9]	Res50	R+S	D	93.44	97.26	70.62
Vehicle-Diff	Res50	R+S	D	94.52	97.97	71.50
PAMTRI [8]	DenseNet121	R+S	D	92.86	96.97	71.88
SAN [61]	Res50	R	-	93.30	-	72.50
VehicleGAN [10]	Res50	R+S	D	93.60	97.30	74.20
CAL [62]	Res50	R	-	95.40	97.90	74.30
MSDeep [15]	Res50	R	-	95.10	-	74.50
VehicleX (PCB) [9]	Res50	R+S	D	94.34	97.91	74.51
Vehicle-Diff (PCB)	Res50	R+S	D	94.40	97.56	75.45
Vanilla Diffusion [29]	SwinV2-B	R+S	B	95.53	98.03	75.95
baseline	SwinV2-B	R	-	96.72	98.57	77.99
CLIP-ReID [63]	ViT-B/16	R	-	95.70	-	79.30
DCAL [64]	ViT-B/16	R	-	96.90	-	80.20
GiT [65]	GiT	R	-	96.86	-	80.34
TransReID [14]	ViT-B/16	R	-	96.90	-	80.60
PCL-CLIP [66]	ViT-B/16	R	-	97.10	98.60	82.50
CLIP-ReID [63]	ViT-B/16	R	-	97.40	-	83.30
VehicleX [9]	SwinV2-B	R+S	D	97.32	98.69	80.36
Vehicle-Diff	SwinV2-B	R+S	D	97.38	98.51	80.98
VehicleX [9]	SwinV2-B	R+S	B	97.08	98.81	81.39
Vehicle-Diff	SwinV2-B	R+S	B	97.68	98.93	83.79

TABLE II: Comparisons with the state-of-the-art methods on VeRi-776 [55]. “S” and “R” denote synthetic and real data, respectively. “B” indicates that each training batch selects equal amounts of synthetic and real data (as introduced in § III-D), whereas “D” indicates that synthetic and real data are combined randomly. Results on two backbones, *i.e.*, Res50 and SwinV2-B, are both reported.

i.e., +6.40%. Furthermore, combined with the original real-world training set, our generated dataset can further improve the reID performance. In particular, our approach achieves 0.94% and 3.57% improvements in mAP compared with VehicleX and PAMTRI, respectively, when jointly trained with the original VeRi-776 training set in Res50 backbone [67]. For SwinV2-Base reID backbone, our method shows a consistent improvement. In VeRi-776 dataset, Vehicle-Diff outperforms VehicleX by 0.62% on mAP when using random combination strategy (“D” in Table II) and 2.4% on mAP when using balanced sampling combination strategy (“B” in Table II). Besides, compared with other state-of-the-art methods, Vehicle-Diff also shows competitive performances. Our Vehicle-Diff method achieves 97.68% Rank-1 and 83.79% mAP, which surpasses CLIP-ReID [63] of 97.40% Rank-1

Method	Data	Rank-1	Rank-5	Rank-10
CityFlowV2				
VehicleX [9]	S	22.21	28.83	35.09
Vehicle-Diff	S	26.38	33.09	36.54
CityFlowV2→VeRi-776				
VehicleX [9]	S	62.04	76.16	81.59
Vehicle-Diff	S	66.81	77.18	83.61

TABLE III: Comparisons with the state-of-the-art method on CityFlowV2 [57]. Generative model is fine-tuned on CityFlowV2, and we do not use any labels in CityFlowV2.

Method	FID↓	
	VeRi-776	CityFlowV2
VehicleGAN [10]	233.0	-
PTGAN [11]	231.1	-
VehicleX	88.20	77.87
Vehicle-Diff	44.84	54.84

TABLE IV: Quantitative comparisons with the state-of-the-art methods on data generation. For a fair comparison, Vehicle-Diff is trained on 1% unlabeled images while VehicleX is trained on 1% labeled images on VeRi-776.

and 83.30% mAP. In CityFlowV2, Vehicle-Diff outperforms VehicleX by 4.17% on Rank-1 and 4.26% on Rank-5 (see the upper part of Table III). We also conduct a series of experiments to verify the generalization ability of Vehicle-Diff. As shown in the bottom of Table III, we apply the reID model trained on the source-domain synthesized data to evaluate performance on the target domain. It should be noted that we only utilize the images in CityFlowV2 to fine-tune the generative model, not the labels. Nonetheless, our Vehicle-Diff consistently outperforms VehicleX.

We further evaluate the quality of the generated data through both quantitative and qualitative evaluation. For the quantitative assessment, we utilize the Frechet Inception Distance (FID) [68], a widely recognized evaluation metric. Unfortunately, since the PAMTRI dataset is not publicly available, we are unable to calculate its FID score. To ensure a fair comparison, we randomly selected 1% of the training datasets to train VehicleX and generate sample images. As shown in Table IV, Vehicle-Diff achieves a lower FID score compared to all other generative methods. For qualitative comparison, we visualize the sample outputs of competitive generative methods in Fig. 1. The images in the first row are from the real-world dataset, while the images in the remaining five rows are from different synthetic data pipeline based on both 3D engines and GAN. We could observe that Vehicle-Diff produces images that are visually closer to the real-world dataset while keeping the fine-grained texture.

C. Ablation Studies and Further Discussion

Effectiveness of the coarse-to-fine strategy. In Vehicle-Diff, we adopt a coarse-to-fine generation strategy. Here we study the effectiveness of each component in our pipeline. Although the filtering process has little effect on the visual gap and the FID change after fine-tuning is negligible, the reID model performance steadily improves (see Table V). Table V validates that quality matters more than quantity,

Components		#IDs	#Imgs	Rank-1	mAP	FID
DFT	SF					
		5,305	191,720	33.19	8.26	126.24
✓		4,940	160,758	58.34	22.00	44.35
✓	✓	4,896	149,472	58.76	22.33	44.78

TABLE V: Ablation study on components, *i.e.*, diffusion fine-tuning (DFT) and semantic filtering (SF).

Baseline	#IDs	#imgs	Rank-1	Rank-5	mAP
IDE	4,894	45,338	57.87	74.97	22.21
	4,896	149,472	58.76	74.43	22.33

TABLE VI: Ablation study on the number of synthetic images for training the reID model on the IDE baseline.

and Table VI shows that more high-quality data leads to better results.

Effectiveness of the balanced sampling strategy. Previous methods, such as VehicleX and PAMTRI, typically conduct random sampling on mixed real and synthetic data to train the model. As a by-product of our pipeline, we introduce a balanced sampling strategy. We merge two mini-batch samples from real and synthetic datasets as a new mini-batch for training. We find that our balanced sampling strategy improves model learning on both VehicleX and Vehicle-Diff data. As shown in the last four rows of Table II, compared to the vanilla sampling strategy, our balanced sampling strategy yields a +1.06% boost in mAP for VehicleX and +2.81% boost in mAP for Vehicle-Diff.

Retrieval visualization. As shown in Fig. 4, we conduct the qualitative image retrieval comparison on VeRi-776. Our method has successfully recalled the target vehicle in the top-5 of the ranking list, surpassing the same model trained on real data or VehicleX. It is because that our Vehicle-Diff contains a large number of vehicle images with fine-grained attributes and intra-class variances such as camera angle, facilitating the discriminative feature learning (see Fig. 3). Therefore, the model trained on our Vehicle-Diff is able to handle challenging matches with fine-grained differences and significant camera angle variations.

V. CONCLUSION

In this paper, we explore the efficacy of state-of-the-art synthetic data generated by a text-to-image model for vehicle re-identification (reID). We introduce Vehicle-Diff, a new coarse-to-fine cross-modality generation pipeline that crafts a synthetic reID dataset tailored to specific downstream scenarios using only 1% of unlabeled images from the original dataset. Our extensive experiments show that this pipeline significantly narrows the gap between synthetic and real-world data, thereby enhancing subsequent reID model performance. Notably, our method achieves a competitive 83.79% mAP on the VeRi-776 dataset. Additionally, we analyze the strengths and limitations of synthetic data across various settings and identify optimal strategies for its utilization. We anticipate that our work will contribute to applications such as privacy protection and intelligent traffic systems.

REFERENCES

- [1] F. Rollo, A. Zunino, N. Tsagarakis, E. M. Hoffman, and A. Ajoudani, "Continuous adaptation in person re-identification for robotic assistance," in *2024 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2024, pp. 425–431.
- [2] X. Liu, W. Liu, H. Ma, and H. Fu, "Large-scale vehicle re-identification in urban surveillance videos," in *2016 IEEE International Conference on Multimedia and Expo (ICME)*. IEEE, 2016, pp. 1–6.
- [3] P. Shyam, K.-J. Yoon, and K.-S. Kim, "Adversarially-trained hierarchical feature extractor for vehicle re-identification," in *2021 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2021, pp. 13 400–13 407.
- [4] Z. Lu, R. Lin, Q. He, and H. Hu, "Mask-aware pseudo label denoising for unsupervised vehicle re-identification," *IEEE Transactions on Intelligent Transportation Systems*, vol. 24, no. 4, pp. 4333–4347, 2023.
- [5] A. Khurshudov, "The smart city conundrum: technology, privacy, and the quest for convenience," *Smart and Sustainable Built Environment*, 2024.
- [6] L. Zheng, L. Shen, L. Tian, S. Wang, J. Wang, and Q. Tian, "Scalable person re-identification: A benchmark," in *Proceedings of the IEEE International Conference on Computer Vision*, 2015, pp. 1116–1124.
- [7] A. Wang, D. Sato, Y. Corzo, S. Simkin, A. Biswas, and A. Steinfeld, "Tbd pedestrian data collection: Towards rich, portable, and large-scale natural pedestrian data," in *2024 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2024, pp. 637–644.
- [8] Z. Tang, M. Naphade, S. Birchfield, J. Tremblay, W. Hodge, R. Kumar, S. Wang, and X. Yang, "Pamtri: Pose-aware multi-task learning for vehicle re-identification using highly randomized synthetic data," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2019, pp. 211–220.
- [9] Y. Yao, L. Zheng, X. Yang, M. Naphade, and T. Gedeon, "Simulating content consistent vehicle datasets with attribute descent," in *Computer Vision—ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part VI 16*. Springer, 2020, pp. 775–791.
- [10] B. Li, P. Liu, L. Fu, J. Li, J. Fang, Z. Xu, and H. Yu, "Vehiclegan: Pair-flexible pose guided image synthesis for vehicle re-identification," *arXiv:2311.16278*, 2023.
- [11] C.-S. Hu, S.-W. Tseng, X.-Y. Fan, and C.-K. Chiang, "Vehicle view synthesis by generative adversarial network," in *ICASSP 2023-2023 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2023, pp. 1–5.
- [12] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio, "Generative adversarial networks," *Communications of the ACM*, vol. 63, no. 11, pp. 139–144, 2020.
- [13] J. Gu, K. Wang, H. Luo, C. Chen, W. Jiang, Y. Fang, S. Zhang, Y. You, and J. Zhao, "Msinet: Twins contrastive search of multi-scale interaction for object reid," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023, pp. 19 243–19 253.
- [14] S. He, H. Luo, P. Wang, F. Wang, H. Li, and W. Jiang, "Transreid: Transformer-based object re-identification," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021, pp. 15 013–15 022.
- [15] Y. Cheng, C. Zhang, K. Gu, L. Qi, Z. Gan, and W. Zhang, "Multi-scale deep feature fusion for vehicle re-identification," in *ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2020, pp. 1928–1932.
- [16] Z. Zheng, T. Ruan, Y. Wei, Y. Yang, and T. Mei, "VehicleNet: Learning robust visual representation for vehicle re-identification," *IEEE Transactions on Multimedia*, vol. 23, pp. 2683–2693, 2020.
- [17] Y. Xu, N. Jiang, L. Zhang, Z. Zhou, and W. Wu, "Multi-scale vehicle re-identification using self-adapting label smoothing regularization," in *ICASSP 2019-2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2019, pp. 2117–2121.
- [18] J. Yu, J. Kim, M. Kim, and H. Oh, "Camera-tracklet-aware contrastive learning for unsupervised vehicle re-identification," in *2022 International Conference on Robotics and Automation (ICRA)*. IEEE, 2022, pp. 905–911.
- [19] B. Jiao, L. Yang, L. Gao, P. Wang, S. Zhang, and Y. Zhang, "Vehicle re-identification in aerial images and videos: Dataset and approach," *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 34, no. 3, pp. 1586–1603, 2024.
- [20] R. Kishore, N. Aslam, and M. H. Kolekar, "Patriid: Pose apprise transformer network for vehicle re-identification," *IEEE Transactions on Emerging Topics in Computational Intelligence (Early Access)*, pp. 1–12, 2024.
- [21] Z. Zheng, L. Zheng, and Y. Yang, "Unlabeled samples generated by gan improve the person re-identification baseline in vitro," in *Proceedings of the IEEE International Conference on Computer Vision*, 2017, pp. 3754–3762.
- [22] S. Yang, Y. Zhou, Z. Zheng, Y. Wang, L. Zhu, and Y. Wu, "Towards unified text-based person retrieval: A large-scale multi-attribute and language search benchmark," in *Proceedings of the 31st ACM International Conference on Multimedia*, 2023, pp. 4492–4501.
- [23] Z. Zheng, X. Wang, N. Zheng, and Y. Yang, "Parameter-efficient person re-identification in the 3d space," *IEEE Transactions on Neural Networks and Learning Systems*, vol. 35, no. 6, pp. 7534–7547, 2024.
- [24] K. Chen, W. Chen, T. He, R. Du, F. Wang, X. Sun, Y. Guo, and G. Ding, "Tagperson: A target-aware generation pipeline for person re-identification," in *Proceedings of the 30th ACM International Conference on Multimedia*, 2022, pp. 560–571.
- [25] S. Xiang, D. Qian, M. Guan, B. Yan, T. Liu, Y. Fu, and G. You, "Less is more: Learning from synthetic data with fine-grained attributes for person re-identification," *ACM Transactions on Multimedia Computing, Communications and Applications*, vol. 19, no. 5s, pp. 1–20, 2023.
- [26] J. Sohl-Dickstein, E. Weiss, N. Maheswaranathan, and S. Ganguli, "Deep unsupervised learning using nonequilibrium thermodynamics," in *International Conference on Machine Learning*. PMLR, 2015, pp. 2256–2265.
- [27] J. Ho, A. Jain, and P. Abbeel, "Denoising diffusion probabilistic models," *Advances in Neural Information Processing Systems*, vol. 33, pp. 6840–6851, 2020.
- [28] R. Rombach, A. Blattmann, D. Lorenz, P. Esser, and B. Ommer, "High-resolution image synthesis with latent diffusion models," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 10 684–10 695.
- [29] D. Podell, Z. English, K. Lacey, A. Blattmann, T. Dockhorn, J. Müller, J. Penna, and R. Rombach, "Sdxl: Improving latent diffusion models for high-resolution image synthesis," *arXiv:2307.01952*, 2023.
- [30] Midjourney, "Midjourney - official website," <https://www.midjourney.com/>, accessed: 2024-09-15.
- [31] R. He, S. Sun, X. Yu, C. Xue, W. Zhang, P. Torr, S. Bai, and X. Qi, "Is synthetic data from generative models ready for image recognition?" in *The Eleventh International Conference on Learning Representations*, 2022.
- [32] S. Azizi, S. Kornblith, C. Saharia, M. Norouzi, and D. J. Fleet, "Synthetic data from diffusion models improves imagenet classification," *arXiv:2304.08466*, 2023.
- [33] A. Nichol, P. Dhariwal, A. Ramesh, P. Shyam, P. Mishkin, B. McGrew, I. Sutskever, and M. Chen, "Glide: Towards photorealistic image generation and editing with text-guided diffusion models," *arXiv:2112.10741*, 2021.
- [34] Z. Zheng, L. Zheng, and Y. Yang, "A discriminatively learned cnn embedding for person reidentification," *ACM Transactions on Multimedia Computing, Communications and Applications*, vol. 14, no. 1, pp. 1–20, 2017.
- [35] L. Zheng, Z. Bie, Y. Sun, J. Wang, C. Su, S. Wang, and Q. Tian, "Mars: A video benchmark for large-scale person re-identification," in *Computer Vision—ECCV 2016: 14th European Conference, Amsterdam, The Netherlands, October 11–14, 2016, Proceedings, Part VI 14*. Springer, 2016, pp. 868–884.
- [36] Y. Sun, C. Cheng, Y. Zhang, C. Zhang, L. Zheng, Z. Wang, and Y. Wei, "Circle loss: A unified perspective of pair similarity optimization," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 6398–6407.
- [37] J. C. Bezdek and R. J. Hathaway, "Some notes on alternating optimization," in *Advances in Soft Computing—AFSS 2002: 2002 AFSS International Conference on Fuzzy Systems Calcutta, India, February 3–6, 2002 Proceedings*. Springer, 2002, pp. 288–300.
- [38] Y. Liang, J. He, G. Li, P. Li, A. Klimovskiy, N. Carolan, J. Sun, J. Pont-Tuset, S. Young, F. Yang, et al., "Rich human feedback for text-to-image generation," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2024, pp. 19 401–19 411.
- [39] X. Lin, Z. Dai, A. Verma, S.-K. Ng, P. Jaillet, and B. K. H.

- Low, "Prompt optimization with human feedback," *arXiv preprint arXiv:2405.17346*, 2024.
- [40] J. Li, D. Li, S. Savarese, and S. Hoi, "Blip-2: Bootstrapping language-image pre-training with frozen image encoders and large language models," in *International Conference on Machine Learning*. PMLR, 2023, pp. 19 730–19 742.
- [41] E. J. Hu, Y. Shen, P. Wallis, Z. Allen-Zhu, Y. Li, S. Wang, L. Wang, and W. Chen, "Lora: Low-rank adaptation of large language models," *arXiv:2106.09685*, 2021.
- [42] R. Sapkota, R. Qureshi, M. F. Calero, M. Hussain, C. Badjugar, U. Nepal, A. Poulouse, P. Zeno, U. B. P. Vaddevolu, H. Yan, *et al.*, "Yolov10 to its genesis: A decadal and comprehensive review of the you only look once series," *arXiv:2406.19407*, 2024.
- [43] A. Radford, J. W. Kim, C. Hallacy, A. Ramesh, G. Goh, S. Agarwal, G. Sastry, A. Askell, P. Mishkin, J. Clark, *et al.*, "Learning transferable visual models from natural language supervision," in *International Conference on Machine Learning*. PMLR, 2021, pp. 8748–8763.
- [44] Z. Liu, H. Hu, Y. Lin, Z. Yao, Z. Xie, Y. Wei, J. Ning, Y. Cao, Z. Zhang, L. Dong, *et al.*, "Swin transformer v2: Scaling up capacity and resolution," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 12 009–12 019.
- [45] C. Szegedy, W. Liu, Y. Jia, P. Sermanet, S. Reed, D. Anguelov, D. Erhan, V. Vanhoucke, and A. Rabinovich, "Going deeper with convolutions," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2015, pp. 1–9.
- [46] Z. Zheng and Y. Yang, "Unsupervised scene adaptation with memory regularization in vivo," *IJCAI*, 2020.
- [47] D. Kinga, J. B. Adam, *et al.*, "A method for stochastic optimization," in *International Conference on Learning Representations (ICLR)*, vol. 5. San Diego, California, 2015, p. 6.
- [48] G. Huang, Z. Liu, L. Van Der Maaten, and K. Q. Weinberger, "Densely connected convolutional networks," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 4700–4708.
- [49] J. Krause, M. Stark, J. Deng, and L. Fei-Fei, "3d object representations for fine-grained categorization," in *Proceedings of the IEEE International Conference on Computer Vision Workshops*, 2013, pp. 554–561.
- [50] Y. Bai, Y. Lou, F. Gao, S. Wang, Y. Wu, and L.-Y. Duan, "Group-sensitive triplet embedding for vehicle reidentification," *IEEE Transactions on Multimedia*, vol. 20, no. 9, pp. 2385–2399, 2018.
- [51] L. Yang, P. Luo, C. Change Loy, and X. Tang, "A large-scale car dataset for fine-grained categorization and verification," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2015, pp. 3973–3981.
- [52] K. Yan, Y. Tian, Y. Wang, W. Zeng, and T. Huang, "Exploiting multi-grain ranking constraints for precisely searching visually-similar vehicles," in *Proceedings of the IEEE International Conference on Computer Vision*, 2017, pp. 562–570.
- [53] H. Liu, Y. Tian, Y. Yang, L. Pang, and T. Huang, "Deep relative distance learning: Tell the difference between similar vehicles," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2016, pp. 2167–2175.
- [54] D. Zapletal and A. Herout, "Vehicle re-identification for automatic video traffic surveillance," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*, 2016, pp. 25–31.
- [55] X. Liu, W. Liu, T. Mei, and H. Ma, "A deep learning-based approach to progressive vehicle re-identification for urban surveillance," in *Computer Vision—ECCV 2016: 14th European Conference, Amsterdam, The Netherlands, October 11–14, 2016, Proceedings, Part II 14*. Springer, 2016, pp. 869–884.
- [56] Z. Tang, M. Naphade, M.-Y. Liu, X. Yang, S. Birchfield, S. Wang, R. Kumar, D. Anastasiu, and J.-N. Hwang, "Cityflow: A city-scale benchmark for multi-target multi-camera vehicle tracking and re-identification," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019, pp. 8797–8806.
- [57] M. Naphade, S. Wang, D. C. Anastasiu, Z. Tang, M.-C. Chang, X. Yang, Y. Yao, L. Zheng, P. Chakraborty, C. E. Lopez, A. Sharma, Q. Feng, V. Ablavsky, and S. Sclaroff, "The 5th ai city challenge," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) Workshops*, June 2021, pp. 4263–4273.
- [58] A. Kanaci, X. Zhu, and S. Gong, "Vehicle re-identification in context," in *Pattern Recognition - 40th German Conference, GCPR 2018, Stuttgart, Germany, September 10–12, 2018, Proceedings*, 2018.
- [59] R. Chu, Y. Sun, Y. Li, Z. Liu, C. Zhang, and Y. Wei, "Vehicle re-identification with viewpoint-aware metric learning," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2019, pp. 8282–8291.
- [60] P. Khorramshahi, A. Kumar, N. Peri, S. S. Rambhatla, J.-C. Chen, and R. Chellappa, "A dual-path model with adaptive attention for vehicle re-identification," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2019, pp. 6132–6141.
- [61] X. Jin, C. Lan, W. Zeng, G. Wei, and Z. Chen, "Semantics-aligned representation learning for person re-identification," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 34, no. 07, 2020, pp. 11 173–11 180.
- [62] Y. Rao, G. Chen, J. Lu, and J. Zhou, "Counterfactual attention learning for fine-grained visual categorization and re-identification," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021, pp. 1025–1034.
- [63] S. Li, L. Sun, and Q. Li, "Clip-reid: exploiting vision-language model for image re-identification without concrete text labels," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 37, no. 1, 2023, pp. 1405–1413.
- [64] H. Zhu, W. Ke, D. Li, J. Liu, L. Tian, and Y. Shan, "Dual cross-attention learning for fine-grained visual categorization and object re-identification," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 4692–4702.
- [65] F. Shen, Y. Xie, J. Zhu, X. Zhu, and H. Zeng, "Git: Graph interactive transformer for vehicle re-identification," *IEEE Transactions on Image Processing*, vol. 32, pp. 1039–1051, 2023.
- [66] J. Li and X. Gong, "Prototypical contrastive learning-based clip fine-tuning for object re-identification," *arXiv preprint arXiv:2310.17218*, 2023.
- [67] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2016, pp. 770–778.
- [68] M. Heusel, H. Ramsauer, T. Unterthiner, B. Nessler, and S. Hochreiter, "Gans trained by a two time-scale update rule converge to a local nash equilibrium," *Advances in Neural Information Processing Systems*, vol. 30, 2017.