

# Appendix for SODA10M: A Large-Scale 2D Self/Semi-Supervised Object Detection Dataset for Autonomous Driving

## A SODA10M dataset

We publish the SODA10M dataset, benchmark, data format and annotation instructions at our website <https://soda-2d.github.io>. It is our priority to protect the privacy of third parties. We bear all responsibility in case of violation of rights, etc., and confirmation of the data license.

**License, Terms of use, Terms of privacy.** The SODA10M dataset is published under CC BY-NC-SA 4.0 license, which means everyone can use this dataset for non-commercial research purpose. Find more details in Appendix F.

**Participants' instructions:** Full participants' instructions can be found in google drive: <https://drive.google.com/file/d/1RVxXltsxkoo-XuzVvAplbZHQs5BGDpEp/view?usp=sharing>.

**Dataset documentation.** <https://soda-2d.github.io/documentation.html> shows the dataset documentation and intended uses.

**Data maintenance.** <https://soda-2d.github.io/download.html> provides data download links for users (in Google Drive & Baidu YunPan). We will maintain the data for a long time and check the data accessibility regularly. We also plan to extend the scale of unlabeled data to a 100-million level to help develop robust models.

**Benchmark and code.** The codebases used in our benchmark are open-source. More details of the reproduction code and experiment settings are illustrated in Appendix B.

**Data format & Evaluation metrics.** Annotation files in SODA10M are stored in standard COCO format, which can be easily accessible to most object detection codebases. We follow the popular COCO API [14] to utilize the Average Precision metric for evaluating detection performance.

**Limitations.** The major limitation of SODA10M is that some domains exist in the unlabeled set but not in labeled sets, which raises the problem that we can not verify the domain adaptation abilities in these domains. To overcome the limitation, we plan to provide more labeled data in these domains in the future.

**Discussion of round1:** The main issues raised by round1 can be summarized and addressed as follows:

Q1: Need to show some interesting or distinct findings of self-supervised learning methods (refer to Sec. 4.3 and Sec. 4.5).

R1: Some interesting observations can be found under autonomous driving scope: Firstly, dense contrastive method (*i.e.*, DenseCL [19]) preforms better than Image-wise methods (*i.e.*, MoCo-v1 [10]) when pre-training on ImageNet (39.9% vs. 39.0%). However, Dense contrastive method preforms worse when pre-training on SODA10M (38.1% vs. 38.9%) due to the reason that pixel-wise contrastive loss may not suitable for complex driving scenario. Secondly, pre-training on ImageNet brings almost equal improvement in both day and night domain (+1.1% vs. +0.9%), and we assume that is because ImageNet [8] contains common images which achieve no special helps to the performance of night domain. On the other hand, pre-training on SDOA10M, which contains massive images collected at night, brings double improvement in night domain (+1.6%) compared to day domain (+0.7%). More observations can be found in Sec. 4.3 and Sec. 4.5.

Q2: Detailed comparison of current autonomous driving datasets (refer to Sec. 3.3).

R2: We compare the SODA10M with other large-scale autonomous driving datasets (including BDD100K [23], nuScenes [1] and Waymo [18]) in the field of scale, driving time, collecting cities, driving conditions and experiment results as upstream pre-training dataset.

Firstly, observed from Table 1 in Introduction, the number of images, driving time and collecting cities of SODA10M is 10M, 27833 hrs and 32 respectively, which is much larger than the current datasets BDD100K [23] (0.1M, 1111 hrs, 4 cities), nuScenes [1] (1.4M, 5.5 hrs, 2 cities) and Waymo [18] (1M, 6.4 hrs, 3 cities). Secondly, as shown the Table 1 in Appendix A, our SODA10M is more diverse in driving conditions compared with nuScenes [1] and Waymo [18], and achieves competitive results with BDD100K [23]. Finally, with above characteristics, SODA10M achieves better generalization

ability and best performance compared with the other three datasets in almost all (9/10) downstream detection and segmentation tasks when regarded as the upstream pre-training dataset.

## B Implement Details

In order to make the experiment results reproducible, we further provide detailed experimental settings for each method in this paper. The primary differences compared with default setting is that the class number is set to 6, syncBN is on and multi-scale training (in supervised and semi-supervised methods) is utilized with scale  $1920 \times (864, 907.2, 950.4, 993.6, 1036.8, 1080)$ . Without specifying, all self-supervised and semi-supervised methods adopt ResNet50 [11] backbone. Our models are trained on servers with 8 Nvidia V100 GPU(32GB) cards with Intel Xeon Platinum 8168 CPU(2.70GHz).

### B.1 Open-Source Codebase

Table 1: The codebases used in our benchmark.

Codebase	link
Detectron2 [20]	<a href="https://github.com/facebookresearch/detectron2">https://github.com/facebookresearch/detectron2</a>
OpenSelfSup	<a href="https://github.com/open-mmlab/OpenSelfSup">https://github.com/open-mmlab/OpenSelfSup</a>
VINCE [9]	<a href="https://github.com/danielgordon10/vince">https://github.com/danielgordon10/vince</a>
STAC [17]	<a href="https://github.com/google-research/ssl_detection">https://github.com/google-research/ssl_detection</a>
Unbiased Teacher [15]	<a href="https://github.com/facebookresearch/unbiased-teacher">https://github.com/facebookresearch/unbiased-teacher</a>

### B.2 Supervised Methods

For the 1x schedule, the learning rate is set to 0.02, decreased by a factor of 10 at 8th, 11th epoch of total 12 epochs. Random crop is used as the only data augmentation method and SGD optimizer is adopted with momentum set as 0.9.

Table 2: Implement details for supervised learning benchmark on SODA10M with 8 Tesla V100.

Model	Train split	Default Setting	Difference	GPU hours
RetinaNet [13] 1x	train set	Detectron2	1. backbone no freeze. 2. turn on precise_bn.	$0.78 \times 8$
Faster RCNN [16] 1x	train set	Detectron2	same as above	$0.83 \times 8$
Cascaded RCNN [2] 1x	train set	Detectron2	same as above	$0.92 \times 8$

### B.3 Self-supervised Methods

We follow the default settings in OpenSelfSup<sup>1</sup> to train six state-of-the-art standard self-supervised learning methods, including MoCo-v1 [10], MoCo-v2 [5], SimCLR [4], SwAV [3], DetCo [21], DenseCL [19], and evaluate their performance by fine-tuning the pre-trained models on the SODA10M labeled data and other self-driving datasets like BDD100K [23] and Cityscapes [6] to verify the generalization ability. Due to the limit of hardware resources, we only use a 5-million unlabeled subset in each experiment by default, while we also make full use of the other 5-million subset in a sequential training manner (mentioned in †), following Hu et al. [12]. Specifically, the model pre-trained on the first subset will be used as initialization to continue pre-training on the second one. We follow the standard data augmentation pipeline adopted by MoCov2 [5], which consists of random resized crop, color jitter, gaussian blur and random horizontal flip, for all considered models except MoCov1 [10], which implements all augmentations except gaussian blur. Auto Augmentation [7] is further used in DetCo [21] following the original paper. Cosine learning rate decay is adopted for all models except MoCov1, which uses step-wise learning rate decay to decrease learning rate by 10x at 40 and 50 epochs. The base learning rates are 0.03, 0.03, 0.3, 4.8, 0.06, 0.03 sequentially for the models in Table 3. LARS optimizer [22] is used in SimCLR [4], while all other models implement SGD optimizer with momentum set as 0.9.

<sup>1</sup><https://github.com/open-mmlab/OpenSelfSup>

For video-based self-supervised learning, MoCo-v1 [10], MoCo-v2 [5] and VINCE [9] are adopted. To ensure fairness, we apply the same data augmentation with VINCE to MoCo-v1 and MoCo-v2 to exploit temporal information and extra jigsaw augmentation to VINCE for better results.

We adopt 3700-epoch, 220-epoch, 325-epoch and 60-epoch pre-training on BDD100K [23], nuScenes [1] and Waymo [18] and SODA unlabeled set for image-based methods respectively, to maintain similar GPU hours with pre-training 200 epochs on ImageNet for fair comparison. Video-based approaches are trained for 800 epochs by considering time limit.

Table 3: Implement details for semi-supervised learning benchmark on SODA10M with 8 Tesla V100.

Model	Train split	Default Setting	Difference	GPU days
MoCov1 [10], MoCov2 [5], SimCLR [4], SwAV [3], DenseCL [19]	5-million unlabeled	OpenSelfSup	60 epochs	$8.40 \times 8$
DetCo [21]	5-million unlabeled	OpenSelfSup	same as above	$14.1 \times 8$
MoCov1 [10]†, MoCov2 [5]†, SimCLR [4]†	10-million unlabeled	OpenSelfSup	same as above	$16.8 \times 8$
MoCov1 [10], MoCov2 [5], VINCE [9], VINCE+Jigsaw [9] on VIDEO	5-million unlabeled to 90k videos	VINCE	1. 800 epochs 2. VINCE augmentations	$2.80 \times 8$

When finetuning on SODA10M labeled set and other datasets (Cityscape [6] & BDD100K [23]), the setting keeps consistent with the supervised methods and MoCo [10] respectively.

#### B.4 Semi-supervised Methods

For semi-supervised methods, considering the time limit, only 1-million unlabeled images (split 0) of SODA10M are used. Compared with self-supervised methods, semi-supervised methods are much more efficient. The learning rate of pseudo labeling, STAC [17] and unbiased teacher [15] is set to 0.02, 0.01 and 0.02 respectively. Color jitter, gaussian blur, affine transformation, cutout augmentation and random crop are adopted in STAC, while color jitter, gray scale, gaussian blur, random erasing and random crop are adopted in unbiased teacher. SGD optimizer is adopted in both methods with momentum set as 0.9.

Table 4: Implement details for semi-supervised learning benchmark on SODA10M with 8 Tesla V100.

Model	Train split	Default Setting	Difference	GPU days
Pseudo Labeling(50K)	50K unlabeled	Detectron2	1. backbone no freeze. 2. turn on precise_bn.	$0.21 \times 8$
Pseudo Labeling(100K)	100K unlabeled	Detectron2	same as above	$0.39 \times 8$
Pseudo Labeling(500K)	500K unlabeled	Detectron2	same as above	$2.00 \times 8$
STAC [17]	1-million unlabeled	STAC	same as default	$2.50 \times 8$
Unbiased Teacher [15]	1-million unlabeled	Unbiased Teacher	same as default	$2.80 \times 8$

## C More Experiments

To show that how the number of labeled images affects the final results in self/semi-supervised object detection task, we further conduct the experiments to show the self-supervised (with MoCov1 [10]) object detection performance under various downstream SODA10M dataset sizes (20%, 50%, 100%) and upstream pre-training datasets (Waymo [18], SODA10M). The results are shown in following:

Table 5: Self-supervised (with MoCov1 [10]) object detection performance mAP(%) under various downstream SODA10M dataset sizes (20%, 50%, 100%) and upstream pre-training datasets (Waymo [18], SODA10M). 20%, 50%, 100% denote for using 20%, 50%, 100% of SODA10M labeled set for downstream fine-tuning.

	20%	50%	100%
pre-train on Waymo [18]	26.6	33.2	37.1
pre-train on SODA10M	$29.2^{+2.6}$	$35.3^{+2.1}$	$38.9^{+1.8}$

Observation can be made that more downstream labeled data bridge the gap between different self-supervised models. On the other hand, too few labeled images tend to have a large standard

deviation (SD) of the final performance, e.g. pre-training on SODA10M receives an 0.75 SD on 20% SODA10M labeled set, compared with the 0.11 SD on 100% SODA10M labeled set.

Table 6 shows the performance of existing self-supervised methods evaluated on SODA10M labeled set with a longer schedule and instance segmentation result on Cityscape [6] dataset. We observe that dense contrastive methods (Detco [21], DenseCL [19]) show excellent results when pre-trained on ImageNet [8], but relatively poor pre-trained on SODA10M unlabeled set. For semantic segmentation performance on Cityscapes with MoCo-v1 [10], the model pre-trained on SODA10M even surpasses the one pre-trained in ImageNet by 1.6%, further verifying the generalization ability of pre-training on SODA10M.

Table 6: Detection results(%) of self-supervised models evaluated on SODA10M labeled dataset (with object detection task) and Cityscapes (CS) dataset (with instance segmentation task).

Pre-train Dataset	Method	Faster-RCNN 2x (SODA)			RetinaNet 2x (SODA)			Mask-RCNN 1x	
		mAP	AP50	AP75	mAP	AP50	AP75	mAP (CS)	AP50 (CS)
ImageNet [8]	random init	29.6	49.8	31.2	20.9	35.4	21.4	25.4	51.1
	super. IN	38.7	61.0	41.5	35.0	57.0	36.0	32.9	59.6
	MoCo-v1 [10]	39.3	60.9	42.5	35.9	57.4	37.3	32.3	59.3
	MoCo-v2 [5]	40.4	62.7	43.6	37.4	59.1	39.3	33.9	60.8
	SimCLR [4]	37.9	61.0	40.4	32.7	53.3	33.8	32.8	59.4
	SwAV [3]	38.2	61.9	40.9	32.6	53.4	33.8	33.9	62.4
	DetCo [21]	39.8	62.1	43.3	35.8	57.5	37.5	34.7	63.2
	DenseCL [19]	40.6	62.9	43.8	37.5	59.4	39.2	34.3	62.5
SODA10M	MoCo-v1 [10]	38.7	60.9	41.1	33.4	56.2	34.3	33.9	60.6
	MoCo-v2 [5]	39.1	60.8	42.6	33.6	56.2	34.8	33.7	61.0
	SimCLR [4]	36.7	59.6	39.1	31.6	53.8	32.3	30.2	57.0
	SwAV [3]	36.0	59.8	37.9	29.7	50.0	30.4	29.4	57.7
	DetCo [21]	37.2	58.9	39.8	31.2	53.5	31.3	32.5	59.8
	DenseCL [19]	38.9	61.0	41.9	33.2	55.4	33.7	33.1	60.7

## D Domain Illustration & Driving Conditions Comparison

The distribution of each fine-grained domain in the validation set, testing set, unlabeled set and labeled set is shown in the Table 7, Table 8, Table 9 and Fig. 1.

Table 7: The number of images in each domain in validation set.

	Daytime			Night		
	City street	Highway	Country road	City street	Highway	Country road
Clear	383	961	63	137	167	312
Overcast	597	517	240	288	627	59
Rainy	177	406	0	0	66	0

Table 8: The number of images in each domain in testing set.

	Daytime			Night		
	City street	Highway	Country road	City street	Highway	Country road
Clear	490	2024	250	1591	498	146
Overcast	1216	1103	481	361	237	133
Rainy	917	520	24	9	0	0

The driving conditions comparison is shown in Table 10.

Table 9: The number of images in each domain in unlabeled set.

	Daytime				Night				Dawn/Dusk			
	Clear	Overcast	Rainy	Snowy	Clear	Overcast	Rainy	Snowy	Clear	Overcast	Rainy	Snowy
City street	2247K	1483K	458K	140K	1274K	582K	157K	71K	325K	215K	62K	22K
Highway	506K	311K	114K	4K	186K	58K	24K	0.70K	37K	27K	9K	0.17K
Country road	499K	333K	69K	7K	170K	40K	12K	1K	38K	23K	5K	0.50K
Residential	146K	154K	29K	20K	61K	40K	7K	10K	9K	1K	2K	2K

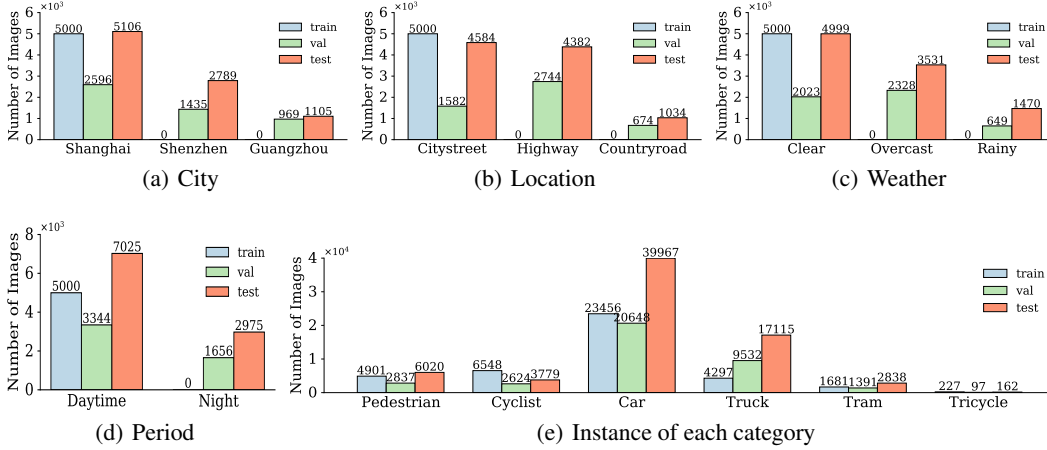


Figure 1: Statistics of the labeled set. (a) Number of images in each city. (b) Number of images in each location. (c) Number of images in each weather condition. (d) Number of images in each period. (e) Number of instances in each category.

Table 10: Driving conditions comparison between SODA10M and other autonomous driving datasets (i.e., nuScenes [1], Waymo [18] and BDD100K [23]), where '-' denotes for not having annotations in this field.

Dataset	Period	Weather	Location
nuScenes [1]	Day: 1045K(88.3%), Night: 138K(11.7%)	Sunny: 952K(80.4%), Rain: 231K(19.6%)	-
Waymo [18]	Day: 798K(80.7%), Night: 96K(9.8%), Dawn/Dusk: 93K(9.5%)	Sunny: 983K(99.4%), Rain: 5K(0.6%)	-
BDD100K [23]	Daytime: 41K(52.6%), Night: 31K(40.1%), Dawn/Dusk: 5K(7.3%)	Clear: 42K(60.6%), Overcast: 10K(14.2%), Rainy: 5K(8.1%), Snowy: 6K(8.9%), Partly cloudy: 5K(8.0%), Foggy: 143(0.2%)	City street: 49K(62.3%), Highway: 19K(25.1%), Residential: 9K(11.8%), Parking lot: 426(0.5%), Gas stations: 34(0.1%), Tunnel: 156(0.2%)
SODA10M	Daytime: 6536K(65.4%), Night: 2697K(26.9%), Dawn/Dusk: 786K(7.7%)	Clear: 5507K(55.7%), Overcast: 3283K(33.6%), Rainy: 949K(8.5%), Snowy: 278K(2.2%)	City street: 7045K(70.7%), Highway: 1283K(12.3%), Country road: 1199K(12.1%), Residential: 490K(4.9%)

More images in each domain are shown in Fig. 2.

## E Acknowledgements

We thank our two data suppliers, named Testin<sup>2</sup> and Speechocean<sup>3</sup> (collected from King-IM-055), for helping us collect and annotate SODA10M dataset.

<sup>2</sup><http://www.testin.cn>

<sup>3</sup><http://en.speechocean.com>





Figure 2: More examples of challenging environments in our dataset.

## F Terms of Use and Licenses

**Description.** Huawei Technologies Co. Ltd (the ‘Organizers’ ,we", "us", and "our",) provides public access to and use of data that it collects and publishes. The data are organized in datasets (the “Datasets”) may be accessed at <https://sslad2021.github.io/index.html>. Any individual or entity (hereinafter “You” or “Your”) with access to the Datasets free of charge subject to the terms of this agreement (hereinafter “Dataset Terms”). By using or downloading the Datasets, you are agreeing to comply with the Dataset Terms and any licensing terms referenced below. Use of any data derived from the Datasets, which may appear in any format such as tables and charts, is also subject to these Dataset Terms.

**Licenses.** Unless specifically labeled otherwise, these Datasets are provided to You under the terms of the Creative Commons Attribution-NonCommercial-ShareAlike 4.0 International Public License

(“CC BY-NC-SA 4.0”), with the additional terms included herein. The CC BY-NC-SA 4.0 may be accessed at <https://creativecommons.org/licenses/by-nc-sa/4.0/legalcode>. When You download or use the Datasets from the Website or elsewhere, You are agreeing to comply with the terms of CC BY-NC-SA 4.0, and also agreeing to the Dataset Terms. Where these Dataset Terms conflict with the terms of CC BY-NC-SA 4.0, these Dataset Terms shall prevail. We reiterate once again that this dataset is used only for non-commercial purposes such as academic research, teaching, or scientific publications. We prohibits You from using the dataset or any derivative works for commercial purposes, such as selling data or using it for commercial gain.

**Sharing.** We prohibits You from distributing this dataset or modified versions. It is permissible to distribute derivative works in as far as they are abstract representations of this dataset (such as models trained on it or additional annotations that do not directly include any of our data).

**Trademark.** All logos and trademarks used on this website are the properties of us or other third parties as stated if applicable. No content provided on the website shall be deemed as granting approval or the right to use any trademark or logo aforesaid by implication, lack of objection, or other means without the prior written consent of us or any third party which may own the mark. No individual shall use the name, trademark, or logo of us by any means without the prior written consent of us.

**Privacy.** We will take reasonable care to remove or scrub personally identifiable information (PII) including, but not limited to, faces of people and license plates of vehicles. Furthermore, We prohibits You from using the Datasets in any manner to identify or invade the privacy of any person even when such use is otherwise legal. If You have any privacy concerns, including to remove your name or other PII from the Dataset, please contact us by sending an e-mail to [xu.hang@huawei.com](mailto:xu.hang@huawei.com).

**Warranties.** The datasets and the website (including, without limitation, all content and modifications of original datasets posted on the website) are provided “as is” and “as available” and without warranty of any kind, express or implied, including, but not limited to, the implied warranties of title, non-infringement, merchantability and fitness for a particular purpose, and any warranties implied by any course of performance or usage of trade, all of which are expressly disclaimed. Without limiting the foregoing, HUAWEI does not warrant that: (a) the content or modifications to the dataset are timely, accurate, complete, reliable or correct in their posted forms at the website; (b) the website will be secure; (c) the website will be available at any particular time or location; (d) any defects or errors will be corrected; (e) the website, content or any modifications are free of viruses or other harmful components; or (f) the results of using the website will meet your requirements. Your use of the website, the datasets, and any content is solely at your own risk. Any entity or individual who suspects that the content on the website (including but not limited to the datasets posted on the website) infringes upon legal rights or interests shall notify our contact [xu.hang@huawei.com](mailto:xu.hang@huawei.com) in written form and provide the identity, ownership certification, associated link (url), and proof of infringement. We will remove the content related to the alleged infringement by law upon receiving the foregoing legal documents.

**Limitation of liability.** In no event shall HUAWEI and its affiliates, or their directors, employees, agents, partners, or suppliers, be liable under contract, tort, strict liability, negligence or any other legal theory with respect to the website, the datasets, or any content or user submissions (i) for any direct damages, or (ii) for any lost profits or special, indirect, incidental, punitive, or consequential damages of any kind whatsoever.

**Applicable Law and Dispute Resolution.** Access and all related activities on or through the website shall be governed by, construed, and interpreted in accordance with the laws of the People’s Republic of China. You agree that any dispute between the parties arising out of or in connection with this legal notice or your access and all related activities on or through this website shall be governed by a court with jurisdiction in Shenzhen, Guangdong Province of the People’s Republic of China.

## References

- [1] H. Caesar, V. Bankiti, A. H. Lang, S. Vora, V. E. Liong, Q. Xu, A. Krishnan, Y. Pan, G. Baldan, and O. Beijbom. nuscenes: A multimodal dataset for autonomous driving. In *CVPR*, 2020.
- [2] Z. Cai and N. Vasconcelos. Cascade R-CNN: delving into high quality object detection. In *CVPR*, 2018.

- [3] M. Caron, I. Misra, J. Mairal, P. Goyal, P. Bojanowski, and A. Joulin. Unsupervised learning of visual features by contrasting cluster assignments. In *NeurIPS*, 2020.
- [4] T. Chen, S. Kornblith, M. Norouzi, and G. Hinton. A simple framework for contrastive learning of visual representations. In *ICLR*, 2020.
- [5] X. Chen, H. Fan, R. Girshick, and K. He. Improved baselines with momentum contrastive learning. *arXiv:2003.04297*, 2020.
- [6] M. Cordts, M. Omran, S. Ramos, T. Rehfeld, M. Enzweiler, R. Benenson, U. Franke, S. Roth, and B. Schiele. The cityscapes dataset for semantic urban scene understanding. In *CVPR*, 2016.
- [7] E. D. Cubuk, B. Zoph, J. Shlens, and Q. V. Le. Randaugment: Practical automated data augmentation with a reduced search space. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops*, pages 702–703, 2020.
- [8] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei. Imagenet: A large-scale hierarchical image database. In *CVPR*, 2009.
- [9] D. Gordon, K. Ehsani, D. Fox, and A. Farhadi. Watching the world go by: Representation learning from unlabeled videos. *arXiv:2003.07990*, 2020.
- [10] K. He, H. Fan, Y. Wu, S. Xie, and R. Girshick. Momentum contrast for unsupervised visual representation learning. In *CVPR*, 2020.
- [11] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *CVPR*, 2016.
- [12] D. Hu, Q. Lu, L. Hong, H. Hu, Y. Zhang, Z. Li, A. Shen, and J. Feng. How well self-supervised pre-training performs with streaming data? *arXiv:2104.12081*, 2021.
- [13] T.-Y. Lin, P. Goyal, R. Girshick, K. He, and P. Dollár. Focal loss for dense object detection. In *ICCV*, 2017.
- [14] T.-Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. L. Zitnick. Microsoft COCO: common objects in context. In *ECCV*, 2014.
- [15] Y.-C. Liu, C.-Y. Ma, Z. He, C.-W. Kuo, K. Chen, P. Zhang, B. Wu, Z. Kira, and P. Vajda. Unbiased teacher for semi-supervised object detection. In *ICLR*, 2021.
- [16] S. Ren, K. He, R. Girshick, and J. Sun. Faster r-cnn: Towards real-time object detection with region proposal networks. In *NeurIPS*, 2015.
- [17] K. Sohn, Z. Zhang, C.-L. Li, H. Zhang, C.-Y. Lee, and T. Pfister. A simple semi-supervised learning framework for object detection. *arXiv:2005.04757*, 2020.
- [18] P. Sun, H. Kretzschmar, X. Dotiwalla, A. Chouard, V. Patnaik, P. Tsui, J. Guo, Y. Zhou, Y. Chai, B. Caine, V. Vasudevan, W. Han, J. Ngiam, H. Zhao, A. Timofeev, S. Ettinger, M. Krivokon, A. Gao, A. Joshi, Y. Zhang, J. Shlens, Z. Chen, and D. Anguelov. Scalability in perception for autonomous driving: Waymo open dataset. In *CVPR*, 2020.
- [19] X. Wang, R. Zhang, C. Shen, T. Kong, and L. Li. Dense contrastive learning for self-supervised visual pre-training. In *CVPR*, 2021.
- [20] Y. Wu, A. Kirillov, F. Massa, W.-Y. Lo, and R. Girshick. Detectron2. <https://github.com/facebookresearch/detectron2>, 2019.
- [21] E. Xie, J. Ding, W. Wang, X. Zhan, H. Xu, Z. Li, and P. Luo. DetCo: Unsupervised contrastive learning for object detection. *arXiv:2102.04803*, 2021.
- [22] Y. You, I. Gitman, and B. Ginsburg. Large batch training of convolutional networks. *arXiv preprint arXiv:1708.03888*, 2017.
- [23] F. Yu, H. Chen, X. Wang, W. Xian, Y. Chen, F. Liu, V. Madhavan, and T. Darrell. BDD100K: A diverse driving dataset for heterogeneous multitask learning. In *CVPR*, 2020.