# DOMAIN ADAPTATION FOR UNKNOWN IMAGE DISTORTIONS IN INSTANCE SEGMENTATION

*Maximiliane Gruber, Fabian Brand, Alina Mosebach, Jürgen Seiler, and André Kaup*

Multimedia Communications and Signal Processing
Friedrich-Alexander-Universität Erlangen-Nürnberg (FAU)

## ABSTRACT

Data-driven techniques for machine vision heavily depend on the training data to sufficiently resemble the data occurring during test and application. However, in practice unknown distortion can lead to a domain gap between training and test data, impeding the performance of a machine vision system. With our proposed approach this domain gap can be closed by unpaired learning of the pristine-to-distortion mapping function of the unknown distortion. This learned mapping function may then be used to emulate the unknown distortion in the training data. Employing a fixed setup, our approach is independent from prior knowledge of the distortion. Within this work, we show that we can effectively learn unknown distortions at arbitrary strengths. When applying our approach to instance segmentation in an autonomous driving scenario, we achieve results comparable to an oracle with knowledge of the distortion. An average gain in mean Average Precision (mAP) of up to 0.19 can be achieved.

***Index Terms***— Image Distortions, Unpaired Image-to-Image Translation, Unsupervised Domain Adaptation, Instance Segmentation, Autonomous Driving
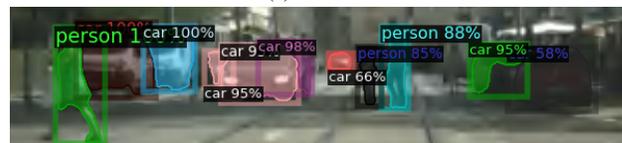
## 1. INTRODUCTION

Depending on the employed image acquisition setup and environmental conditions, images and videos contain different image distortions like blur, noise, low contrast, low resolution, coding artifacts, or combinations thereof. In previous work, it has been shown that the performance of Deep Neural Network (DNN)-based techniques for machine vision decreases if the input images or videos are subject to such distortions. This negative impact has been shown for image classification [1–5], semantic segmentation [6], object detection [7, 8], instance segmentation [8] and license plate recognition [9]. A common approach to encounter this decrease in performance is to enlarge data sets by data augmentation, i.e., extending them with modified versions of the original images by applying expected distortions synthetically [10].

However, not all image and video distortions may be easily modeled and applied to pristine images. For example, translating images between different camera setups is a more complex problem. Therefore, first attempts at learning pristine-to-distorted mapping functions have been made [11]. A pristine-to-distorted mapping function translates images from the source domain $X$ containing pristine, undistorted images to the target domain $Y$ containing distorted images. Learning mappings between images belonging to different domains is also termed image-to-image translation.

(a) Baseline



(b) Our approach

**Fig. 1**. Visual comparison of instance segmentation results on Cityscapes image `frankfurt_0000000_001236` for test data subject to an unknown distortion (here: JPEG2000 at a PSNR of 32 dB) with and without our proposed approach for unpaired learning of unknown distortions.

For paired training of image-to-image translation, images from domain $X$ and $Y$ related by the true pristine-to-distorted mapping function $\mathcal{C}$ are required. In [12], Isola et al. have proposed the `pix2pix` framework for paired image-to-image translation. This method employs adversarial learning, and is shown to be a flexible framework able to learn a large variety of mappings between domains. Paired image-to-image translation is impossible when regarding *unknown distortions*, since access to the true distortion function $\mathcal{C}$ is required to obtain the training data. In [11], Chen et al. show that their `pix2pix`-based approach is capable of learning distortions from paired data. The visual impression is verified by an evaluation in terms of natural scene statistics. However, the effectiveness of the approach is not investigated with regards to machine vision tasks. In [13], Zhu et al. extend `pix2pix` to `CycleGAN` by introducing a cycle-consistency loss. This enables the unpaired learning of mappings between domains in a similar flexible framework. For unpaired image-to-image translation unrelated images may be taken from domains $X$ and $Y$ enabling the learning of *unknown distortions*.

In contrast to utilizing image-to-image translation to align data of differing domains on a pixel-level, a domain gap may also be overcome by aligning the training and testing data on feature-level. However, the advantage of pixel-level approaches is that the alignment of data and the method to solve the machine vision task may be regarded independently. For this reason, we only regard alignment on pixel-level, i.e., pixel-level domain adaptation, by means of image-to-image translation in this work.

In this work, we show that with our approach for unpaired learning of unknown distortions, blur, white noise, JPEG coding, JPEG2000 coding, and HEIF coding may be emulated at several

levels of distortion. We compare the performance of instance segmentation in an autonomous driving scenario when adapting the model by means of corresponding true and learned pristine-to-distorted mapping functions. In this context, we also perform an extensive benchmark of the impact of image distortions on instance segmentation employing Mask Region-based Convolutional Neural Network (R-CNN).

A visual impression of the advantage of employing our proposed unpaired learning of unknown distortions in instance segmentation is given in Figure 1. In Figure 1, multiple instances in a distorted image are not recognized by the baseline system trained on pristine data. By adapting the instance segmentation to the unknown distortion (here: JPEG 2000 at a PSNR of $32\,\text{dB}$) with our proposed approach, more instances are recognized.

## 2. UNSUPERVISED DOMAIN ADAPTATION FOR UNKNOWN DISTORTIONS

DNNs for instance segmentation in autonomous driving are commonly trained on pristine data. However, in practical applications, the test data is often subject to unknown distortions impeding the performance of the DNN. We propose to overcome this domain shift between the source domain $X$ containing pristine, undistorted data and the target domain $Y$ subject to an unknown distortion by means of unpaired image-to-image translation. The pristine-to-distorted mapping function $\mathcal{C} : X \rightarrow Y$ represents the unknown mapping from the source to the target domain. Learning this mapping function from unpaired data enables the emulation of the unknown distortion on the labeled training data. With this unsupervised domain adaptation, the performance in the target domain may be improved without access to labeled training data.

### 2.1. Unpaired Learning of Unknown Distortions

For unpaired learning of the unknown pristine-to-distorted mapping function $\mathcal{C}$, the image-to-image translation technique `CycleGAN` is employed [13]. This unpaired approach consists of a generator $G$ to translate images from the undistorted to the distorted domain $G : X \rightarrow Y$ and a generator $F$ to translate images from the distorted to the undistorted domain $F : Y \rightarrow X$. Employing a cycle-consistency loss, the difference between input images $\boldsymbol{x} \in X$ and $\boldsymbol{y} \in Y$ and their respective translations employing both generators $\hat{\boldsymbol{x}} = F(G(\boldsymbol{x}))$ and $\hat{\boldsymbol{y}} = G(F(\boldsymbol{y}))$ is minimized. Furthermore, discriminators $D_X$ and $D_Y$ exist to distinguish between true and generated samples of domains $X$ and $Y$, respectively. In our approach, the learned mapping function $\tilde{\mathcal{C}}$ for the unknown distortion is given by generator $G$. The optimal generator $G^*$ is obtained by solving

$$G^* = \text{argmin}_{G,F,D_X,D_Y} \, \mathcal{L}(G, F, D_X, D_Y). \quad (1)$$

For further details on `CycleGAN` and the full objective function $\mathcal{L}(G, F, D_X, D_Y)$, we refer the reader to the original publication and the reference implementation provided by the authors [13].

With the goal of obtaining one set of parameters to learn different unknown distortions at arbitrary strengths, multiple preliminary experiments were conducted. We empirically found the setup presented in [13] for the translation between paintings and photos to perform best over a broad range of distortions at different strengths. These training parameters entail a cycle-consistency loss weighing factor of 10, an identity mapping loss weighing factor of 0.5 and a batch size of 1. The models are trained from scratch, employing a constant learning rate of 0.0002 during the first 100 epochs,
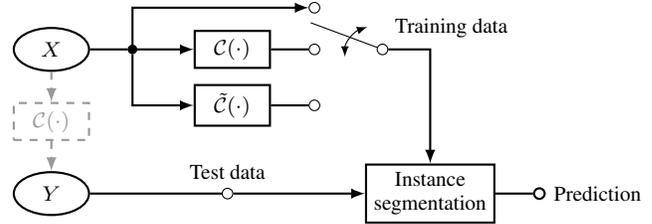


**Fig. 2**. Experimental setup of unsupervised domain adaptation for unknown distortions. Mapping function $\mathcal{C}$ represents the unknown mapping from undistorted source domain $X$ to distorted target domain $Y$. Mapping function $\tilde{\mathcal{C}}$ is learned by means of unpaired image-to-image translation.

and linearly decaying it to zero over the following 100 epochs. The generator networks $G$ and $F$ consist of 9 residual blocks. Discriminators $D_X$ and $D_Y$ are $70 \times 70$ PatchGANs [13]. The training data is randomly cropped to a size of $256 \times 256$.

An advantage of our proposed approach is the independence from prior knowledge of the unknown distortions by employing a single set of parameters for a wide range of distortions at arbitrary levels. The flexibility to adapt to a wide range of distortions at various levels is of particular importance, since commonly images and videos are not only subject to a single distortion, but combinations thereof.

Another benefit of our approach is the unnecessity of additional data sets. We employ the instance segmentation test images from the target domain $Y$ and the instance segmentation training images from the undistorted source domain $X$ to train the image-to-image translation in an unpaired manner. The learned mapping $\tilde{\mathcal{C}}$ is then employed to emulate the unknown distortion on the instance segmentation training data.

### 2.2. Adapting Instance Segmentation to Unknown Distortions

In [8], Fischer et al. show that robustness against distortion may be obtained either by including the degraded images into the training data, or by fine-tuning a network trained on pristine data on the respective distortion. Therefore, we choose the same pre-trained Mask R-CNN [14] provided in [15] and perform a fine-tuning on training data distorted by means of the learned pristine-to-distortion mapping function $\tilde{\mathcal{C}}$.

The `Detectron2` framework [15] is employed to fine-tune Mask R-CNN [14] for instance segmentation. The employed pre-trained model has a ResNet50 [16] backbone, and is pre-trained on COCO [17] and pristine Cityscapes [18]. With a batch size of 4, the fine-tuning is performed for a maximum of $24\,000$ iterations, to improve the DNN's performance on the target domain. Starting with a learning rate of 0.01, the learning rate is decreased after $18\,000$ iterations to 0.001.

## 3. EXPERIMENTAL SETUP

The experimental setup to evaluate our proposed approach is depicted in Figure 2. We employ data from the undistorted source domain $X$ as training data. We regard three different scenarios illustrated by the three different branches:

**Baseline** In the top branch, we directly employ data from the undistorted source domain $X$ for training. Hence, the instance segmentation system is not adapted to the distorted domain $Y$ and the distortion remains unknown.
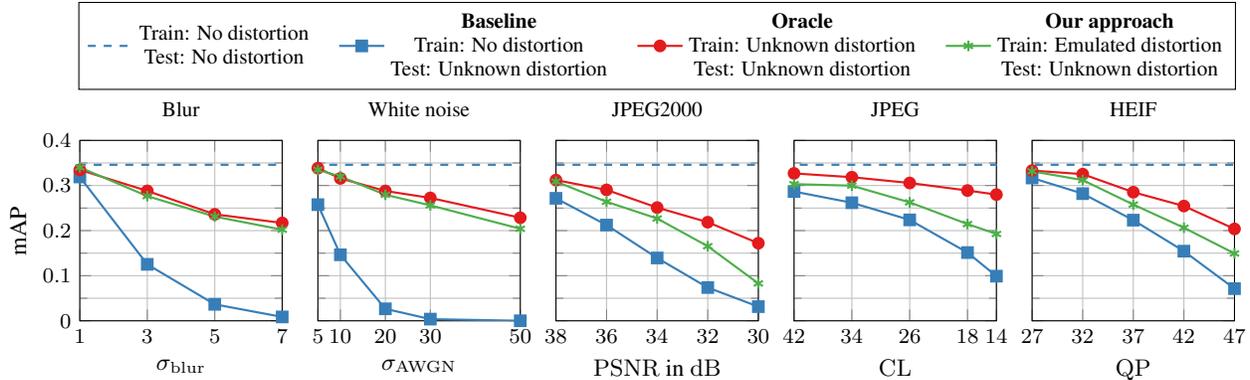
**Fig. 3**. Results of instance segmentation measured as mean Average Precision (mAP) over distortion level for various distortion types.

**Oracle-based approach** In the middle branch, the true pristine-to-distorted mapping function $\mathcal{C}$ is employed to adapt the instance segmentation. This represents the best-case scenario. However, since our proposed approach is meant to adapt instance segmentation to unknown distortions of arbitrary strength, this *oracle-based* approach is unattainable in practice.

**Our proposed approach** As detailed in Section 2.1, we learn a mapping function $\tilde{\mathcal{C}}$ from pristine images to an unknown distortion by means of unpaired image-to-image translation. We emulate the unknown distortion on the training data employing $\tilde{\mathcal{C}}$ to adapt instance segmentation to the unknown distortion.

The evaluation is carried out employing test data from the distorted domain $Y$. This distorted test data is obtained by applying the true distortion mapping function $\mathcal{C}$ to the Cityscapes validation split. The performance of Mask R-CNN is evaluated in terms of mAP and reported over the level of distortion. The mAP is calculated as described in [18], by first calculating the Average Precision (AP) for each class for various overlaps, and then taking the mean over all classes.

We employ the Cityscapes dataset introduced in [18]. We choose this data set, since it contains relatively little distortions in comparison to other data sets for autonomous driving. For each distortion employed in this work, we degrade all splits of the Cityscapes data set in order to train the unpaired image-to-image translation as well as the instance segmentation.

To show the flexibility of our proposed approach, we choose a diverse set of distortions at different levels of distortion, namely blur, additive white Gaussian noise (AWGN) and different image compression techniques. We show the applicability of our approach for compression with a fixed block size (JPEG), with adaptive block-size (HEIF), and wavelet-based coding (JPEG2000).

The distortion level of *blur* is denoted by standard deviation $\sigma_{\text{blur}}$ of the two-dimensional symmetric Gaussian kernel, with which the pristine image is convolved. For *white noise*, random values are drawn independently for each pixel from a zero-mean Gaussian distribution with standard deviation $\sigma_{\text{AWGN}}$ and added to the undistorted image. For *JPEG* and *JPEG2000* en- and decoding methods provided by `ImageMagick` [19] are employed. The JPEG quality is varied by means of Compression Level (CL), with zero leading to the strongest and 100 to the lowest level of distortion. The quality of JPEG2000 is controlled by the Peak Signal-to-Noise Ratio (PSNR) in dB between pristine and coded image. *HEIF* en- and decoding is performed employing `libheif` [20]. The distortion level is controlled by the Quantization Parameter (QP) ranging from zero to 51. With a QP of 51 the strongest distortions are introduced.

| Distortion | Oracle | Our approach | Difference |
|---|---|---|---|
| Blur | 0.15 | 0.14 | −0.01 |
| White noise | 0.20 | 0.19 | −0.01 |
| JPEG2000 | 0.10 | 0.06 | −0.04 |
| JPEG | 0.10 | 0.05 | −0.05 |
| HEIF | 0.07 | 0.04 | −0.03 |

**Table 1**. Average mean Average Precision (mAP) gains over baseline model without adaption to distortions, when employing an oracle and our proposed approach. The last column shows the difference between our proposed approach and the oracle-based approach.

## 4. RESULTS AND DISCUSSION

The experimental results in terms of mAP over the level of distortion are shown in Figure 3. The distortion levels are sorted so that the leftmost value corresponds to the lowest and the rightmost value to the highest distortion. The blue dashed line shows the mAP obtained when testing the pre-trained baseline model on pristine data. With this setup the best results are obtained reaching an mAP of 0.35. For all solid plots, the test data is distorted with the respective true distortion function $\mathcal{C}$. The blue solid line depicts the results for the baseline scenario without fine-tuning of the model. For all types of distortion the mAP decreases with the increasing level of distortion. The red solid line depicts the results of the oracle with knowledge of the true distortion function $\mathcal{C}$. Here each distortion at each strength is regarded as its own target domain, i.e., for each red point a specific instance segmentation model was adapted. For very low distortion like a blur of $\sigma_{\text{blur}} = 1$, the difference between the pre-trained baseline approach and the adapted oracle-based approach is very small. For stronger degradations, the mAP is increased by the specialization on the regarded distortion.

The green solid line depicts the results for our proposed unpaired learning of unknown distortions. For all distortions at all distortion levels, the mAP is increased in comparison to the baseline, i.e., the case without adaptation to the unknown distortion. For blur and white noise, mAPs very close to the oracle-based approach are reached. For the different image coding techniques, the obtained mAP is also quite close to the mAP achieved by the oracle. For stronger distortions, the gap between mAPs obtained by the true and the learned distortion function becomes larger.

In Table 1, the gains in mAP over the baseline are averaged over the level of distortion. It can be seen that we achieve results comparable to the oracle-based approach. The largest average gain in mAP of 0.19 is observed for white noise. The average gain in mAP for

**Fig. 4**. Visual examples of applying true (top row) and learned (bottom row) distortion mapping function to Cityscapes image `bremen_000117_000019` for various distortion types. *(Best to be viewed enlarged on a monitor.)*

(a) Pristine original image

(b) Blur $\sigma_{\mathrm{blur}} = 3$

(c) White noise $\sigma_{\mathrm{AWGN}} = 20$

(d) JPEG2000 PSNR $= 32\,\mathrm{dB}$

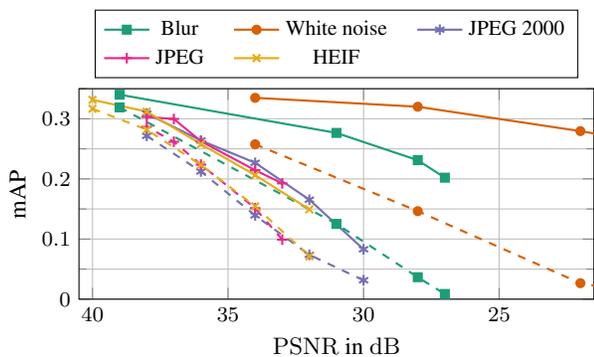(e) JPEG CL $= 18$

(f) HEIF QP $= 42$



**Fig. 5**. Results of instance segmentation measured as mean Average Precision (mAP) over Peak Signal-to-Noise Ratio (PSNR) for various distortion types. Dashed plots denote results obtained with the baseline, solid lines denote results obtained with our proposed approach.

white noise and blur achieved with our proposed approach is only 0.01 lower than the oracle-based approach. For image compression with our proposed approach average gains ranging from 0.04 to 0.06 are achieved. The gains obtained with the oracle are between 0.03 to 0.05 higher.

Visual examples for one distortion level of each distortion type are shown in Figure 4. In Figure 4(a), a pristine example image is depicted. In Figure 4(b)-(f), this image is distorted by the true pristine-to-distorted mapping $\mathcal{C}$ in the top row, and the the learned mapping $\tilde{\mathcal{C}}$ in the bottom row. For blur and white noise it can be seen that the learned distortion is visually very similar to the true distortion. In the case of the different image compression techniques, the learned and true distortions are also visually close. However, not all types of occurring artifacts can be reproduced yet.

For a better comparison in terms of objective image quality, the results of the instance segmentation measured as mAP, are shown as a function of PSNR in Figure 5. Therefore, the mean PSNR was computed for each distortion level of each distortion type. The dashed plots denote the results obtained with the baseline, the solid plots represent the results of our proposed approach. In the baseline approach, for the same distortion level in terms of PSNR, the highest mAPs are observed for white noise, followed by blur. All regarded image compression techniques have a similar influence, resulting in a lower mAP. Our proposed unpaired learning of unknown distortions improves the instance segmentation results for all PSNR

levels. The sensitivity towards the different distortion types remains the same.

It can be seen, that the employed instance segmentation model exhibits a strong sensitivity towards various image distortion. The instance segmentation seems to be very sensitive towards image compression, while exhibiting more resiliency toward blur and white noise. As expected, for the *unknown distortion*, i.e., model trained on undistorted source domain $X$ and tested on distorted target domain $Y$, the lowest mAPs are achieved. With the regarded *oracle-based approach*, i.e., model specialized and tested on distorted target domain $Y$, the highest mAPs are obtained. With our proposed *unpaired learning of unknown distortions* results comparable to the *oracle-based approach* are obtained. Moreover, we can achieve mAPs for blur and white noise closely matching the oracle-based approach. For the different image coding techniques, especially at a higher level of distortion there is still room for improvement. A potential reason for the larger gap in image coding techniques is the inability of the employed generator network to reproduce block artifacts.

## 5. CONCLUSION

In this work we propose to learn the mapping from pristine data to data subject to an unknown distortion. We show that this mapping may be learned with a fixed setup for a wide range of distortions and for a diverse set of distortion levels. This learned mapping is then employed to emulate unknown distortions on pristine, labeled training data for instance segmentation. By fine-tuning instance segmentation with this training data, we achieve results that are comparable to an oracle with knowledge of the true pristine-to-distorted mapping. The largest average gain in mAP of 0.19 is obtained for white noise. For image coding and blur we achieve average gains in mAP between 0.04 and 0.14. Our approach has the advantage that neither prior knowledge of the distortion nor additional data sets are required. Furthermore, we performed an extensive benchmark on the influence of image distortions on instance segmentation.

In future work, the difference between specialization of DNNs to certain degradations of certain strengths could be compared to a model trained on various degradations. Furthermore, an evaluation in terms of visual closeness of the emulated and true distortions may be conducted. Applying our proposed approach to more complex combinations of distortions could be performed to investigate the generalization capabilities of our approach.

## 6. REFERENCES

[1] S. Dodge and L. Karam, "Understanding How Image Quality Affects Deep Neural Networks," in *Proc. International Conference on Quality of Multimedia Experience (QoMEX)*, June 2016, pp. 1–6.

[2] S. Ghosh, R. Shet, P. Amon, A. Hutter, and A. Kaup, "Robustness of Deep Convolutional Neural Networks for Image Degradations," in *Proc. IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, Apr. 2018, pp. 2916–2920.

[3] D. Hendrycks and T. Dietterich, "Benchmarking Neural Network Robustness to Common Corruptions and Perturbations," in *Proc. International Conference on Learning Representations (ICLR)*, May 2019, pp. 1–16.

[4] Y. Pei, Y. Huang, Q. Zou, X. Zhang, and S. Wang, "Effects of Image Degradation and Degradation Removal to CNN-Based Image Classification," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 43, no. 4, pp. 1239–1253, Apr. 2021.

[5] K. Endo, M. Tanaka, and M. Okutomi, "CNN-Based Classification of Degraded Images With Awareness of Degradation Levels," *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 31, no. 10, pp. 4046–4057, Oct. 2021.

[6] C. Kamann and C. Rother, "Benchmarking the Robustness of Semantic Segmentation Models with Respect to Common Corruptions," *International Journal of Computer Vision*, vol. 129, no. 2, pp. 462–483, Feb. 2021.

[7] R. Geirhos, C. R. M. Temme, J. Rauber, H. H. Schütt, M. Bethge, and F. A. Wichmann, "Generalisation in Humans and Deep Neural Networks," in *Proc. Advances in Neural Information Processing Systems (NeurIPS)*, Dec. 2018, pp. 7538–7550.

[8] K. Fischer, C. Blum, C. Herglotz, and A. Kaup, "Robust Deep Neural Object Detection and Segmentation for Automotive Driving Scenario with Compressed Image Data," in *Proc. IEEE International Symposium on Circuits and Systems (ISCAS)*, May 2021, pp. 1–5.

[9] P. Kaiser, F. Schirrmacher, B. Lorch, and C. Riess, "Learning to Decipher License Plates in Severely Degraded Images," in *Proc. Pattern Recognition. ICPR International Workshops and Challenges*, Jan. 2021, pp. 544–559.

[10] C. Shorten and T. M. Khoshgoftaar, "A Survey on Image Data Augmentation for Deep Learning," *Journal of Big Data*, vol. 6, pp. 1–48, July 2019.

[11] L.-H. Chen, C. G. Bampis, Z. Li, and A. C. Bovik, "Learning to Distort Images Using Generative Adversarial Networks," *IEEE Signal Processing Letters*, vol. 27, pp. 2144–2148, Nov. 2020.

[12] P. Isola, J.-Y. Zhu, T. Zhou, and A. A. Efros, "Image-to-Image Translation with Conditional Adversarial Networks," in *Proc. IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, July 2017, pp. 5967–5976.

[13] J. Zhu, T. Park, P. Isola, and A. A. Efros, "Unpaired Image-to-Image Translation Using Cycle-Consistent Adversarial Networks," in *Proc. IEEE International Conference on Computer Vision (ICCV)*, Oct. 2017, pp. 2242–2251.

[14] K. He, G. Gkioxari, P. Dollár, and R. Girshick, "Mask R-CNN," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 42, no. 2, pp. 386–397, Feb. 2020.

[15] Y. Wu, A. Kirillov, F. Massa, W.-Y. Lo, and R. Girshick, "Detectron2," https://github.com/facebookresearch/detectron2, Oct. 2019.

[16] K. He, X. Zhang, S. Ren, and J. Sun, "Deep Residual Learning for Image Recognition," in *Proc. IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2016, pp. 770–778.

[17] 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 *Proc. European Conference on Computer Vision (ECCV)*, Sept. 2014, pp. 740–755.

[18] 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 *Proc. IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2016, pp. 3213–3223.

[19] The ImageMagick Development Team, "ImageMagick," https://imagemagick.org.

[20] Struktur AG, "Libheif," https://github.com/strukturag/libheif.