

# DeePaste - Inpainting for Pasting

Levi Kassel

The Hebrew University of Jerusalem  
Jerusalem, Israel

levi.kassel@mail.huji.ac.il

Michael Werman

The Hebrew University of Jerusalem  
Jerusalem, Israel

michael.werman@mail.huji.ac.il

## Abstract

*One of the challenges of supervised learning training is the need to procure an substantial amount of tagged data. A well-known method of solving this problem is to use synthetic data in a copy-paste fashion, so that we cut objects and paste them onto relevant backgrounds. Pasting the objects naively results in artifacts that cause models to give poor results on real data. We present a new method for cleanly pasting objects on different backgrounds so that the dataset created gives competitive performance on real data. The main emphasis is on the treatment of the border of the pasted object using inpainting. We show state-of-the-art results both on instance detection and foreground segmentation.*

## 1. Introduction

In recent years, with the dramatic development of detection and segmentation in computer vision with the help of AI, the demand has increased to put these capabilities into existing systems such as mobile apps, autonomous vehicles, robots, etc. One of the biggest problems with state of the art systems is the amount of tagged data needed for each task or scene to train them.

For each new task, as well as for any change in the environment, such as a change in the objects you wish to detect or detecting objects in a different scene, one is required to create thousands or more tagged images based on the same task or scene. Since data tagging is very labor intensive, deploying those tasks is not carried out in most cases.

In recent years more and more researchers have been able to overcome this challenge by creating synthetic data from synthetically rendered scenes and objects [15] [34] [16] [31], so that they render objects and scenes for training detection and segmentation systems. Although these systems bring promising results, they require a high level of graphics know how and a lot of effort to make objects and scenes look real .

Moreover, models that train on such data find it diffi-

cult to give satisfactory results on real data because of the change in the statistics of the image [5] [29]. To overcome this difficulty and to still generate tagged data faster, more and more researchers are working on composing real images [9] [48] [11]. The general idea is to cut real objects and paste them on real scenes, thus getting free tagged images that fit the environment in which we want to work.

However, if we use this paradigm naively - we will get poor results on real data. This is because existing systems are more sensitive to local region based features than to global scene layout. The implication is that when we paste an object on an image we are creating subtle pixel artifacts on the background images. This phenomenon prevents the system from generalizing to real images. Dwibedi et al. [9] create a framework that combines all kinds of simple blending and blurring types on the objects and achieves surprisingly impressive results.

In this article, we present a new method of blending to overcome the difficulty of local-level realism. We use the power of inpainting trained on relevant scenes and thus succeed in filling in the gaps in required areas - in our case - the area between the object and the background, the *blending gap*.

Because we know the environment in which we want to work, we train an inpainting model that knows how to fill the blending gap so that the detection and segmentation systems can be generalized to real data.

In recent years there have also been works on improving the blending and improving the harmonization of objects that are pasted on a background image such as Deep Harmonization [42] [7] or Deep Blending [47] [52]. These systems were not compared because the purpose of these systems is graphical, and perform the learning on each image and object individually [52] - which makes the purpose of creating a large dataset non-feasible. What we do is different. The images we produce are not graphically perfect - but are good enough to cheat the learning systems so that they can generalize well on real data.

In this article, we evaluate this method and its effectiveness and show that it achieves state of the art results on in-

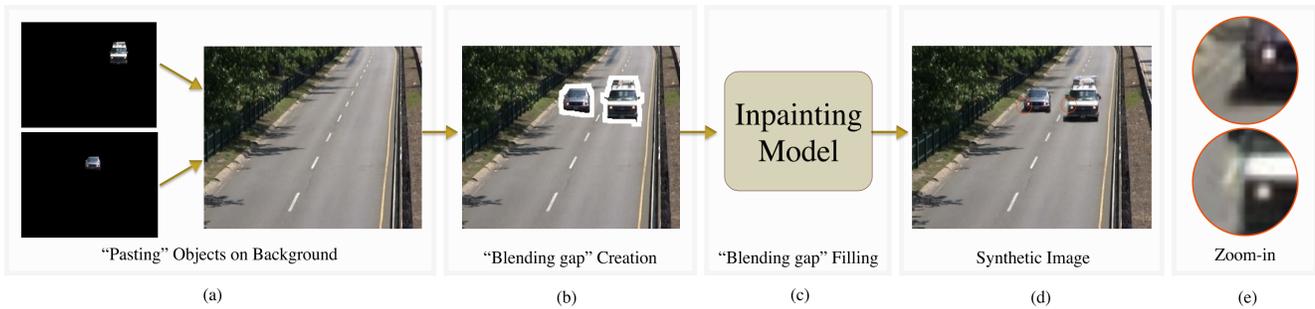


Figure 1. This is the main pipeline of our system. First we paste objects on the background image (a). Second, we create the blending gap mask (b). We use the trained inpainting model to fill the blending gap pixels (c). In (d) you can see the final image. The pixels are filled with the right context (e).

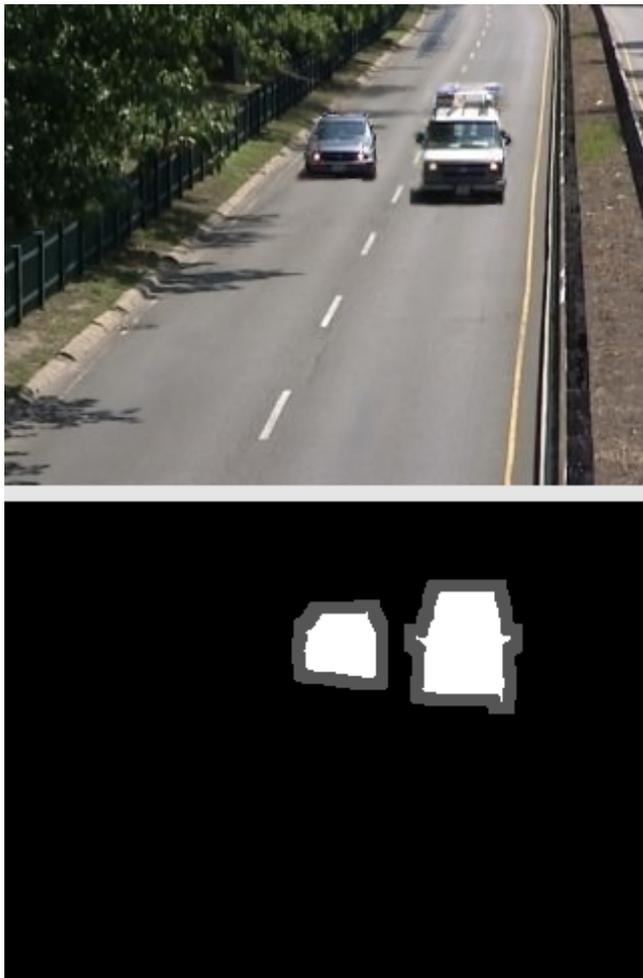


Figure 2. The pasting in the top image is over the white masks. The inpainting is on the gray area.

stance detection and foreground segmentation.

## 2. Related Work

Instance detection is where we try to locate a particular object within a target image. Early approaches, such as [6] tried to solve this task by using local features of the desired object image such as SIFT [26], FAST [8], or SURF [4], and try to match them with local features extracted from the target image. These methods did not work well when the objects we wanted to detect were partially occluded or did not have enough features [14] [17].

Another well-known computer vision task is foreground segmentation (also called background subtraction), determining which pixels belong to the background and which to the foreground, mainly from videos. Early approaches, such as [39] [38] [43] [3], were mainly based on analyzing the displacement of certain parts of the image and trying to decide if it is a displacement of a foreground object or a natural change that occurs in the background such as shading, change in lighting, or natural displacement of dynamic objects belonging to the background.

Modern methods in these two areas are mainly deep learning systems based on convolution network architectures so that the extraction of the features is done more semantically even in difficult cases [33] [32] [24] [45]. Most of the best performing systems are supervised - which means that they learn directly from tagged images and are trained end to end.

In recent years, with the availability of powerful hardware, these methods are more suitable for real-time applications and thus the demand is increasing to use these capabilities in more and more areas like robots, navigation, and surveillance [46] [28].

These methods require large amounts of annotated data. There is a bottleneck as each environment and scene requires a collection of annotated data which significantly slows the deployment of these systems.

One method that solves this problem is to create a synthetic dataset. One setting where synthetic data can be generated is by rendering the environment and thus automati-

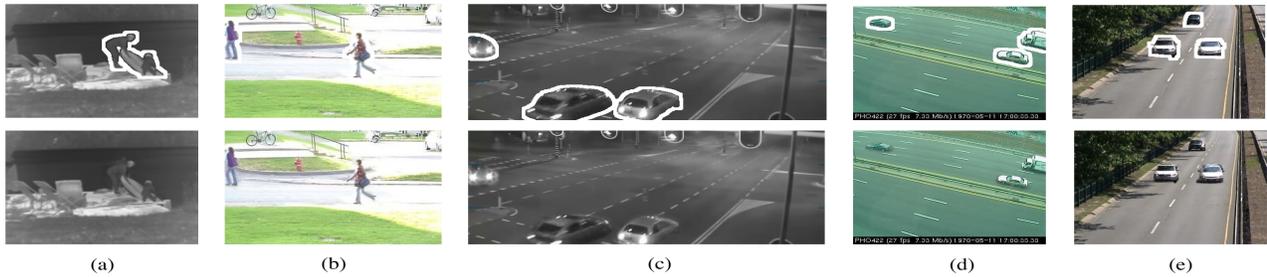


Figure 3. Sample pairs of blending gap filling from the CDnet 2014 dataset [44] in thermal (a), baseline (b,e), night videos (c), and low frame-rate (d) scenes.

cally creating a dataset tailored to the specific problem we want to solve [15] [34]. A big advantage of this method is that the space of possibilities is large - you can create any environment you want and you can change the objects and background simply as needed [16] [31]. Unfortunately, when trying to generalize those systems from synthetic data to real world images - we encounter poor performance due to the difference in the statistics between the images space [5] [29].

Another method to produce synthetic data is by composing real images. In general, the method works by taking existing objects, which you want to detect or to segment, cut them, with a corresponding image mask, select a background image, and paste the object in a certain position in the background image. There are certain algorithms where you can automatically choose where to paste an object to maintain the realism of the created image [41] [20]. This process is very scalable and you can easily create a variety of data for the required purpose.

Dwibedi et al. [9] took this process and examined it on the task of instance segmentation. In their paper, they notice that modern deep learning based systems care more about local region-based features for detection than the global scene layout. This is why it is very important to pay attention to how the pasting process takes place. When we naively place objects in scenes subtle pixel artifacts are introduced into the images. As these minor imperfections in the pixel space feed forward deeper into the layers of a ConvNet [22], they lead to noticeably different features and the training algorithm focuses on these discrepancies to detect objects, often ignoring to model their complex visual appearance. To overcome this problem they used blending methods such as Gaussian Blurring and Poisson editing [30] and found that the use of a blending method in the pasting process can alter the detection performance dramatically. In addition, they found that if they synthesize the same scene with the same object placement only varying the type of blending used makes the training algorithm invariant to these blending factors and improves performance

by a big margin.

Kassel et al [20] used this framework in foreground segmentation. They used a weak foreground segmenter to extract objects from the training images and insert them in their original position into a background image. It is especially pertinent to static cameras as the objects found are automatically in the right location, being of the right size, color, and shape, and in the right lighting conditions.

Another task in computer vision that is related to image synthesis is image inpainting. The main idea of image inpainting is to fill missing pixels in the image. The main difficulty in this task is to synthesize visually realistic and semantically plausible pixels for the missing regions that are coherent with existing ones.

Early work on this subject [2] [13] attempted to solve this task by searching for similar areas in the image and attempting to paste patches into the missing holes while maintaining global consistency. These works work well especially in images that are characterized by stationary textures but are limited to non-stationary data such as natural images [35]. In recent years, deep learning has also entered this field and GAN-based [12] approaches have yielded promising results [18] [23]. Yang et al. [49] and [50] devised feature shift and contextual attention operations, respectively, to allow the model to borrow feature patches from distant areas of the image. Another work by Yu et al. [51] tries to handle irregular holes by filling using gated convolutions.

### 3. Approach Overview

As in *Cut, Paste, and Learn* [9], we focus only on the process of pasting objects on the background images, as part of the creation of synthetic data. From their paper, it can be seen that there is a high correlation between the pasting process and how the system focuses on the appearance of the pasted object and so manages to generalize well on real data. As mentioned, if you paste the object somehow on the background image, there are pixel artifacts in the gap between the object and the background image which greatly affects how the system treats real data. In other words, the

object is fine, the background is fine, but the border between them causes problems. It can be seen that if we perform even just Gaussian smoothing at the edges of the object - the generalization significantly increases [9]. In this paper, we used this insight, and the power of inpainting, to make this connection as natural as possible in terms of convolution networks.

Here we list the stages we did to enhance the pasting process.

- **Train an inpainting model** We train a model to learn how to paste objects on background images in a natural way. We train an inpainting model that learns the statistics of the scene and can fill designated holes naturally. We tune our model in filling the blending gap as well as in filling irregular holes in this environment.
- **Paste objects on background** When pasting the object on the background image we automatically create a blending gap mask that separates the object from the background image, Figure 1 (b).
- **Fill the blending gap** We use the inpainting model trained in Stage 1 to fill in the blending gap, Figures 1 (c,d).

## 4. Approach Details

We now present our approach in detail.

### 4.1. Train an inpainting model

To fill the blending gap correctly we need to use a model that can fill this area so that it will be as natural as possible. This model needs to have the capability to harmoniously fill many *irregular* holes in the image. For that we chose to train the DeepFill-V2 inpainting model [51], which is designed especially to deal with filling *irregular* structures. This model uses a simple algorithm to automatically generate random free-form masks on-the-fly during training. So that the model will succeed to fill the blending gap correctly we train it on the background images. In addition to enhancing the inpainting model’s capabilities specifically on the pasting work, we add to the training set some fixed mask around real objects, with a weak segmenter, such that filling the blending gap will be fine-tuned.

### 4.2. Paste objects on background

To paste the object on the background image we need to have an image of the object together with its corresponding mask. Our method is agnostic to how we find the segmentation mask of the object we need to paste. If you do not have the ground truth of the mask, you can use pre-trained foreground segmenters such as [19] depending on the task required. Appropriate augmentations can now be performed

on the object such as rotation, scaling, illumination, etc. Then you have to choose where to paste the object on the image. There are settings where a random placement [9] is sufficient and there are settings where the location is important for improving the performance of the model [20].

The blending gap around object mask,  $\mathcal{M}$ , is produced in the instance segmentation case as  $Dilate(\mathcal{M}) \setminus Erode(\mathcal{M})$  and in the foreground segmentation as  $Dilate(\mathcal{M}) \setminus \mathcal{M}$  as the background of the object is correct for the scene and we do not have to worry about an incorrect background leaking into the object, Figure 2.

One of the advantages of our method is that even if the object segmentation is not that accurate at the border of the object, where mistakes are prone to happen - our method can handle it.

### 4.3. Fill the blending gap

Once we have the background image along with the objects pasted on it with the blending gap along with the corresponding inpainting mask (Figure 1 b), they can be inserted as input into the inpainting model that we trained in the first step. The inference process is fast and is performed for all images in the environment in the same way. The results may not be perfect in terms of graphics, but they do exactly the job we are interested in - cause the learning models to focus on appearance and not recognize that it is a pasted object.

## 5. Experiments

Our method is general and can be applied in creating synthetic data for a wide variety of tasks. This paper explored the capabilities of our method in two main computer vision tasks: instance detection and foreground segmentation.

### 5.1. Training and Evaluation on Instance Detection

To test our method in instance detection task we used the Dwibedi et al. [9] environment as in the cut, paste, and learn paper. We use a total of 33 object instances from the Big-BIRD Dataset [37] overlapping with the 11 instances from GMU Kitchen Dataset [10]. For the foreground/background masks of the objects, we use a pretrained foreground segmentation network [19]. We use backgrounds from the UW Scenes Dataset [21]. We trained a inpainting model for this task on 1000 images, 600 of them were from the various background images from the UW Scenes Dataset, and another 400 randomly picked images from the BigBIRD Dataset from which all the objects that we pasted are. All the images were trained on inpainting irregular holes. We generated a synthetic dataset with approximately 6000 images using all the augmentations they suggested e.g. scale, rotation, etc. where position and the background were chosen randomly. Each background appears roughly few times in the generated dataset with different objects.

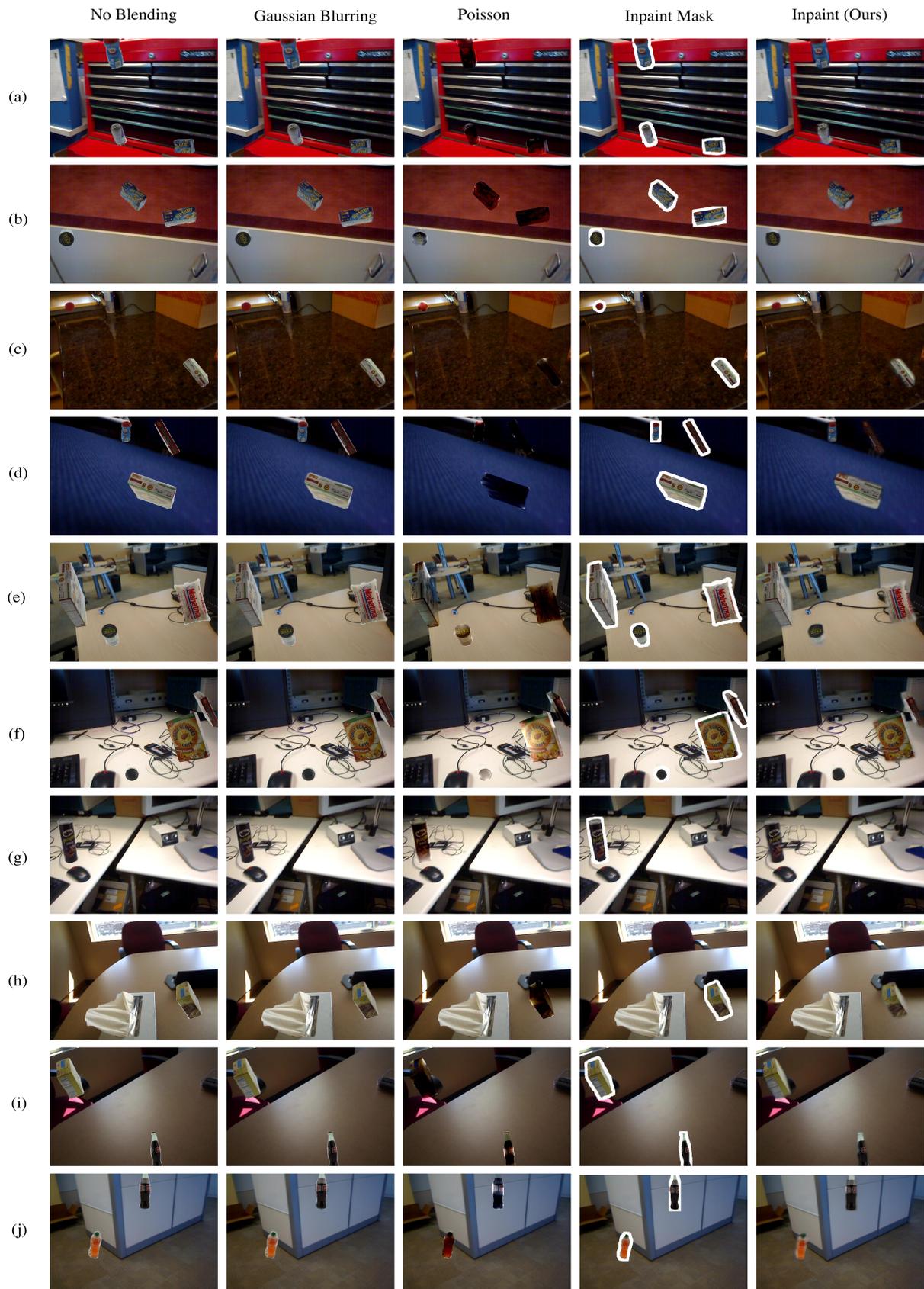


Figure 4. Sample synthetically generated images from the instance detection setting. Each column refers to a different blending technique. We add the inpainting mask column to show which pixels were filled in with the inpainting model.

Category	Coca Cola	Coffee Mate	Honey Bunches	Hunt's Sauce	Mahatama Rice	Nature V1	Nature V2	Palmolive Orange	Pop Secret	Pringles BBQ	Red Bull	mAP
No blending	64.3	87.4	83.5	57.9	61.4	92.3	79.4	59.2	54.8	44.4	32.2	65.2
Gaussian Blurring	65.4	86.4	80.4	69.3	64.3	90.3	83.4	59.3	59.4	66.4	42.5	69.8
Poisson	63.6	83.2	67.3	55.3	29.2	85.5	66.4	60.3	70.8	50.3	19.0	59.2
Inpaint(Ours)	70.3	90.4	84.0	<b>73.7</b>	62.4	<b>94.3</b>	83.5	73.4	65.4	68.3	44.5	73.7
All Blend+Same Image(No Inpaint)	77.3	91.0	79.3	69.3	67.2	92.9	81.2	65.2	77.1	<b>71.8</b>	40.2	73.9
All Blend+Same Image(With Inpaint)	<b>80.4</b>	<b>93.0</b>	<b>84.6</b>	73.3	<b>72.9</b>	94.1	<b>88.4</b>	<b>75.4</b>	<b>78.9</b>	66.8	<b>47.9</b>	<b>77.8</b>

Table 1. Evaluation results on the GMU Dataset [10] from models that were trained with different pasting techniques.

Unsupervised Foreground Segmentation methods Comparison				
	Baseline	Night Videos	Low Frame Rate	Thermal
<b>BSUV-Net 2.0</b>	0.962	0.585	0.790	0.893
<b>SemanticBGS</b>	0.960	0.501	0.788	0.821
<b>IUTIS-5</b>	0.956	0.529	0.774	0.830
<b>Ours</b>	<b>0.991</b>	<b>0.752</b>	<b>0.893</b>	<b>0.902</b>

Table 2. F-measure comparison of state of the art unsupervised methods from CDnet dataset.

We created 4 different synthetic datasets that differ in terms of the pasting process e.g. No blending, Gaussian blurring, Poisson editing, and Inpainting (ours). In addition, we generated another dataset called All Blend + same image e.g. synthesize the same scene with the same object placement, and only vary the type of blending used to make the training algorithm further ignore the effects of blending as they suggested. Samples of the generated images compared to other blending methods can be seen in Figure 4.

We can see that the inpainting model fills the missing pixels with existing shapes of the background and combines them with the texture of the object. The border of the object is slightly blurred and the transition is smooth. This causes the detection system not to focus on the border of the object but on its appearance. In contrast, the transition between the background and the objects in the no blending and gaussian blur methods are more sharp and thus more noticeable. In addition, we can see that the Poisson blending method often causes large parts of the object to be occluded in order to create a good blending.

We use a Faster R-CNN model [33] based on VGG-16 [36] pre-trained weights on the MSCOCO [25] detection task to train the detector on the synthetic dataset we created. For evaluation, we use the GMU Kitchen Dataset [10] which contains 9 kitchen scenes with 6,728 images. We evaluate on the 11 objects present in the dataset overlapping with the BigBIRD [37] objects. We report Average Precision (AP) at IOU of 0.5 in all our experiments for the task of instance localization. Table 1 shows the evaluation results. From the results, we can see that our approach beats almost all compared blending methods by a big margin. In addition, we also made an ablation study in the All Blend + same image setting. The model that was trained with data that

was generated with our blending method performed better than without it.

## 5.2. Training and Evaluation on Foreground Segmentation

To test our method in foreground segmentation we used the setup in Kassel et al. [20]. The idea is to use an unsupervised segmenter to extract the foreground objects from the training frames and when pasted it will be pasted in its original location. This is well suited for the static camera setting where the objects are placed in the right location, being of the right size, color and shape, and in the right lighting conditions.

For this purpose we use four scene categories from the Change Detection 2014 (CDnet) dataset [44] baseline, night videos, low frame-rate, and thermal. The background images are pixel-wise medians of a sequence of 50 frames. The frames used to extract objects were the same 200 used by FgSegNet V2 [24]. For the proposal on extracting the foreground objects we use the state of the art unsupervised method on the CDnet dataset [44], BSUV-Net v2 [40]. For every scene, we trained a unique inpainting model that was trained on 200 training images in random masking holes and also a fixed set of blending gap filling that we created with the BSUV-Net v2 results. For each scene, we generated 500 training frames. Samples from the generated frames can be seen in Figure 3.

We can see that inpainting models fill the missing pixels with the right context. e.g. shadow/light of the car, lines on the road, etc. In general, the object is inserted in a good and consistent way to blend with the background (Figure 1e). To test our generated data we use the FgSegNet V2 by Lim et al. [24], which is a state-of-the-art method in the CDnet 2014 Challenge [44], SBI2015 [27], and UCSD Background Subtraction [1]. We ignored blending gap pixels during the learning process since we didn't know if the inpainting model will generate foreground or background pixels. Since this approach is totally unsupervised we compared it to the state of the art unsupervised methods. Results are shown in Table 2. We can see that the model that trained on data that was generated with our method performs better than all the state of the art methods by a big margin, improving the existing segmentation in the unsupervised setting.

## 6. Summary

We presented a novel approach to paste objects on background images when synthesizing annotated training images for tasks like detection and segmentation. Inspired by the notion that local realism is sufficient for training convolutions based models we focused our efforts on the connection between the pasted objects and the background image - we called this connection the blending gap. Our approach uses the inpainting to learn the task environment and fill the blending gap missing pixels successfully and improves the generalization by a big margin.

## References

- [1] Ucsd anomaly detection dataset. In [//www.svcl.ucsd.edu/projects/anomaly/dataset.htm](http://www.svcl.ucsd.edu/projects/anomaly/dataset.htm), 2010. 6
- [2] Connelly Barnes, Eli Shechtman, Adam Finkelstein, and Dan B Goldman. Patchmatch: A randomized correspondence algorithm for structural image editing. *ACM Trans. Graph.*, 28(3):24, 2009. 3
- [3] Olivier Barnich and Marc Van Droogenbroeck. Vibe: A universal background subtraction algorithm for video sequences. *IEEE Transactions on Image processing*, 20(6):1709–1724, 2010. 2
- [4] Herbert Bay, Tinne Tuytelaars, and Luc Van Gool. Surf: Speeded up robust features. In *European conference on computer vision*, pages 404–417. Springer, 2006. 2
- [5] Wenzheng Chen, Huan Wang, Yangyan Li, Hao Su, Zhenhua Wang, Changhe Tu, Dani Lischinski, Daniel Cohen-Or, and Baoquan Chen. Synthesizing training images for boosting human 3d pose estimation. In *2016 Fourth International Conference on 3D Vision (3DV)*, pages 479–488. IEEE, 2016. 1, 3
- [6] Alvaro Collet, Manuel Martinez, and Siddhartha S Srinivasa. The moped framework: Object recognition and pose estimation for manipulation. *The international journal of robotics research*, 30(10):1284–1306, 2011. 2
- [7] Wenyan Cong, Jianfu Zhang, Li Niu, Liu Liu, Zhixin Ling, Weiyuan Li, and Liqing Zhang. Dovenet: Deep image harmonization via domain verification. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8394–8403, 2020. 1
- [8] E. Rosten; T. Drummond. Machine learning for high-speed corner detection. In *European Conference on Computer Vision*, pages 430–443, 2006. 2
- [9] Debidatta Dwibedi, Ishan Misra, and Martial Hebert. Cut, paste and learn: Surprisingly easy synthesis for instance detection. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 1301–1310, 2017. 1, 3, 4
- [10] Georgios Georgakis, Md Alimoor Reza, Arsalan Mousavian, Phi-Hung Le, and Jana Košecká. Multiview rgb-d dataset for object instance detection. In *2016 Fourth International Conference on 3D Vision (3DV)*, pages 426–434. IEEE, 2016. 4, 6
- [11] Golnaz Ghiasi, Yin Cui, Aravind Srinivas, Rui Qian, Tsung-Yi Lin, Ekin D Cubuk, Quoc V Le, and Barret Zoph. Simple copy-paste is a strong data augmentation method for instance segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2918–2928, 2021. 1
- [12] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. *Advances in neural information processing systems*, 27, 2014. 3
- [13] James Hays and Alexei A Efros. Scene completion using millions of photographs. *ACM Transactions on Graphics (ToG)*, 26(3):4-es, 2007. 3
- [14] Stefan Hinterstoisser, Cedric Cagniart, Slobodan Ilic, Peter Sturm, Nassir Navab, Pascal Fua, and Vincent Lepetit. Gradient response maps for real-time detection of textureless objects. *IEEE transactions on pattern analysis and machine intelligence*, 34(5):876–888, 2011. 2
- [15] Stefan Hinterstoisser, Vincent Lepetit, Paul Wohlhart, and Kurt Konolige. On pre-trained image features and synthetic images for deep learning. In *Proceedings of the European Conference on Computer Vision (ECCV) Workshops*, pages 0–0, 2018. 1, 3
- [16] Stefan Hinterstoisser, Olivier Pauly, Hauke Heibel, Marek Martina, and Martin Bokeloh. An annotation saved is an annotation earned: Using fully synthetic training for object detection. In *Proceedings of the IEEE/CVF International Conference on Computer Vision Workshops*, pages 0–0, 2019. 1, 3
- [17] Edward Hsiao, Alvaro Collet, and Martial Hebert. Making specific features less discriminative to improve point-based 3d object recognition. In *2010 IEEE computer society conference on computer vision and pattern recognition*, pages 2653–2660. IEEE, 2010. 2
- [18] Satoshi Iizuka, Edgar Simo-Serra, and Hiroshi Ishikawa. Globally and locally consistent image completion. *ACM Transactions on Graphics (ToG)*, 36(4):1–14, 2017. 3
- [19] Suyog Jain, Bo Xiong, and Kristen Grauman. Pixel objectness. *arXiv preprint arXiv:1701.05349*, 2017. 4
- [20] Levi Kassel and Michael Werman. Using a supervised method without supervision for foreground segmentation. *arXiv preprint arXiv:2011.07954*, 2020. 3, 4, 6
- [21] Kevin Lai, Liefeng Bo, Xiaofeng Ren, and Dieter Fox. A large-scale hierarchical multi-view rgb-d object dataset. In *2011 IEEE international conference on robotics and automation*, pages 1817–1824. IEEE, 2011. 4
- [22] Yann LeCun, Bernhard Boser, John S Denker, Donnie Henderson, Richard E Howard, Wayne Hubbard, and Lawrence D Jackel. Backpropagation applied to handwritten zip code recognition. *Neural computation*, 1(4):541–551, 1989. 3
- [23] Yijun Li, Sifei Liu, Jimei Yang, and Ming-Hsuan Yang. Generative face completion. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3911–3919, 2017. 3
- [24] Long Ang Lim and Hacer Yalim Keles. Learning multi-scale features for foreground segmentation. *Pattern Analysis and Applications*, 23(3):1369–1380, 2020. 2, 6

- [25] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *European conference on computer vision*, pages 740–755. Springer, 2014. 6
- [26] David G Lowe. Distinctive image features from scale-invariant keypoints. *International journal of computer vision*, 60(2):91–110, 2004. 2
- [27] Lucia Maddalena and Alfredo Petrosino. Towards benchmarking scene background initialization. In *International conference on image analysis and processing*, pages 469–476. Springer, 2015. 6
- [28] Shervin Minaee, Yuri Y Boykov, Fatih Porikli, Antonio J Plaza, Nasser Kehtarnavaz, and Demetri Terzopoulos. Image segmentation using deep learning: A survey. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2021. 2
- [29] Xingchao Peng, Baochen Sun, Karim Ali, and Kate Saenko. Learning deep object detectors from 3d models. In *Proceedings of the IEEE international conference on computer vision*, pages 1278–1286, 2015. 1, 3
- [30] Patrick Pérez, Michel Gangnet, and Andrew Blake. Poisson image editing. In *ACM SIGGRAPH 2003 Papers*, pages 313–318, 2003. 3
- [31] Param S Rajpura, Hristo Bojinov, and Ravi S Hegde. Object detection using deep cnns trained on synthetic images. *arXiv preprint arXiv:1706.06782*, 2017. 1, 3
- [32] Joseph Redmon and Ali Farhadi. Yolov3: An incremental improvement. *arXiv preprint arXiv:1804.02767*, 2018. 2
- [33] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. Faster r-cnn: Towards real-time object detection with region proposal networks. *Advances in neural information processing systems*, 28:91–99, 2015. 2, 6
- [34] Artem Rozantsev, Vincent Lepetit, and Pascal Fua. On rendering synthetic images for training an object detector. *Computer Vision and Image Understanding*, 137:24–37, 2015. 1, 3
- [35] Denis Simakov, Yaron Caspi, Eli Shechtman, and Michal Irani. Summarizing visual data using bidirectional similarity. In *2008 IEEE Conference on Computer Vision and Pattern Recognition*, pages 1–8. IEEE, 2008. 3
- [36] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*, 2014. 6
- [37] Arjun Singh, James Sha, Karthik S Narayan, Tudor Achim, and Pieter Abbeel. Bigbird: A large-scale 3d database of object instances. In *2014 IEEE international conference on robotics and automation (ICRA)*, pages 509–516. IEEE, 2014. 4, 6
- [38] Pierre-Luc St-Charles, Guillaume-Alexandre Bilodeau, and Robert Bergevin. Subsense: A universal change detection method with local adaptive sensitivity. *IEEE Transactions on Image Processing*, 24(1):359–373, 2014. 2
- [39] Chris Stauffer and W Eric L Grimson. Adaptive background mixture models for real-time tracking. In *Proceedings. 1999 IEEE computer society conference on computer vision and pattern recognition (Cat. No PR00149)*, volume 2, pages 246–252. IEEE, 1999. 2
- [40] M Ozan Tezcan, Prakash Ishwar, and Janusz Konrad. Bsvnet 2.0: Spatio-temporal data augmentations for video-agnostic supervised background subtraction. *IEEE Access*, 9:53849–53860, 2021. 6
- [41] Shashank Tripathi, Siddhartha Chandra, Amit Agrawal, Ambrish Tyagi, James M Rehg, and Visesh Chari. Learning to generate synthetic data via compositing. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 461–470, 2019. 3
- [42] Yi-Hsuan Tsai, Xiaohui Shen, Zhe Lin, Kalyan Sunkavalli, Xin Lu, and Ming-Hsuan Yang. Deep image harmonization. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3789–3797, 2017. 1
- [43] Marc Van Droogenbroeck and Olivier Paquot. Background subtraction: Experiments and improvements for vibe. In *2012 IEEE computer society conference on computer vision and pattern recognition workshops*, pages 32–37. IEEE, 2012. 2
- [44] Yi Wang, Pierre-Marc Jodoin, Fatih Porikli, Janusz Konrad, Yannick Benezeth, and Prakash Ishwar. Cdnets 2014: An expanded change detection benchmark dataset. In *Proceedings of the IEEE conference on computer vision and pattern recognition workshops*, pages 387–394, 2014. 3, 6
- [45] Yi Wang, Zhiming Luo, and Pierre-Marc Jodoin. Interactive deep learning method for segmenting moving objects. *Pattern Recognition Letters*, 96:66–75, 2017. 2
- [46] Joe Watson, Josie Hughes, and Fumiya Iida. Real-world, real-time robotic grasping with convolutional neural networks. In *Annual Conference Towards Autonomous Robotic Systems*, pages 617–626. Springer, 2017. 2
- [47] Huikai Wu, Shuai Zheng, Junge Zhang, and Kaiqi Huang. Gp-gan: Towards realistic high-resolution image blending. In *Proceedings of the 27th ACM international conference on multimedia*, pages 2487–2495, 2019. 1
- [48] Yue Wu, Wael Abd-Almageed, and Prem Natarajan. Busternet: Detecting copy-move image forgery with source/target localization. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 168–184, 2018. 1
- [49] Chao Yang, Xin Lu, Zhe Lin, Eli Shechtman, Oliver Wang, and Hao Li. High-resolution image inpainting using multi-scale neural patch synthesis. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 6721–6729, 2017. 3
- [50] Jiahui Yu, Zhe Lin, Jimei Yang, Xiaohui Shen, Xin Lu, and Thomas S Huang. Generative image inpainting with contextual attention. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 5505–5514, 2018. 3
- [51] Jiahui Yu, Zhe Lin, Jimei Yang, Xiaohui Shen, Xin Lu, and Thomas S Huang. Free-form image inpainting with gated convolution. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 4471–4480, 2019. 3, 4
- [52] Lingzhi Zhang, Tarmily Wen, and Jianbo Shi. Deep image blending. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 231–240, 2020. 1