
ZeroWaste Dataset: Towards Automated Waste Recycling

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Less than 35% of recyclable waste is being actually recycled in the US [1], which
2 leads to increased soil and sea pollution and is one of the major concerns of
3 environmental researchers as well as the common public. At the heart of the
4 problem is the inefficiencies of the waste sorting process (separating paper, plastic,
5 metal, glass, etc.) due to the extremely complex and cluttered nature of the
6 waste stream. Automated waste detection strategies have a great potential to
7 enable more efficient, reliable and safer waste sorting practices, but the literature
8 lacks comprehensive datasets and methodology for the industrial waste sorting
9 solutions. In this paper, we take a step towards computer-aided waste detection
10 and present the first in-the-wild industrial-grade waste detection and segmentation
11 dataset, ZeroWaste. This dataset contains over 1800 fully segmented video
12 frames collected from a real waste sorting plant along with waste material labels
13 for training and evaluation of the segmentation methods, as well as over 6000
14 unlabeled frames that can be further used for semi-supervised and self-supervised
15 learning techniques. ZeroWaste also provides frames of the conveyor belt
16 before and after the sorting process, comprising a novel setup that can be used
17 for weakly-supervised segmentation. We present baselines for fully-, semi- and
18 weakly-supervised segmentation methods. Our experimental results demonstrate
19 that state-of-the-art segmentation methods struggle to correctly detect and classify
20 target objects which suggests the challenging nature of our proposed in-the-wild
21 dataset. We believe that ZeroWaste will catalyze research in object detection and
22 semantic segmentation in extreme clutter as well as applications in the recycling
23 domain. Our project page can be found at <http://ai.bu.edu/zerowaste/>.

24 1 Introduction

25 As the world population grows and gets increasingly urbanized, waste production is estimated to
26 reach 2.6 billion tonnes a year in 2030, an increase from its current level of around 2.1 billion tonnes
27 [5]. Efficient recycling strategies are critical to reduce the devastating environmental effects of rising
28 waste production. Materials Recovery Facilities (MRFs) are at the center of the recycling process.
29 These facilities are where the collected recyclable waste is sorted into separate bales of plastic, paper,
30 metal and glass. The accuracy of the sorting directly determines the quality of the recycled material;
31 for high-quality, commercially viable recycling, the contamination levels (anything but the desired
32 material) need to be less than a few percent of the bale. Even though the MRFs utilize a large
33 number of machinery alongside manual labor [6], the extremely cluttered nature of the waste stream
34 makes automated waste detection (*i.e.* detection of waste objects that should be removed from the
35 conveyor belt) very challenging to achieve, and the recycling rates as well as the profit margins stay at
36 undesirably low levels (e.g. less than 35% of the recyclable waste actually got recycled in the United
37 States in 2018 [1]). Another crucial aspect of manual waste sorting is the safety of the workers that
38 risk their lives daily picking up unsanitary objects (*e.g.* medical needles).

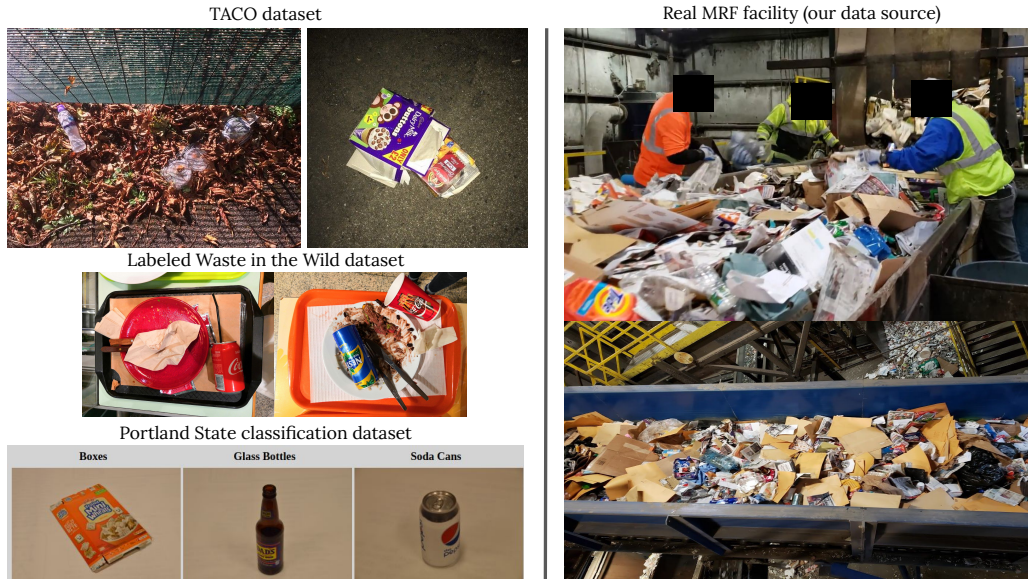


Figure 1: **Left:** examples of the existing waste detection and classification datasets (top to bottom): Trash Annotation in Context (TACO) [2], Labeled Waste in the Wild [3], Portland State University Recycling [4] datasets. **Right:** footage of the waste sorting process at a real Materials Recovery Facilities (MRF). The domain shift between the simplified datasets with solid background and little to no clutter and the real images of the conveyor belt from the MRF makes it impossible to use models trained on these datasets for automated detection on real waste processing plants. In this paper, we propose a new `ZeroWaste` dataset collected from a real waste sorting plant. Our dataset includes a set of densely annotated frames for training and evaluation of the detection and segmentation models, as well as a large number of unlabeled frames for semi- and self-supervised learning methods. We also include frames of the conveyor belt before and after manual collection of foreground objects to facilitate research on weakly supervised detection and segmentation. Please see Figure 2 for the illustration of our `ZeroWaste` dataset.

39 Recent advances in object classification and segmentation provide a great potential to make the
 40 recycling process more efficient, more profitable and safer for the workers. Accurate waste clas-
 41 sification and detection algorithms have a potential to enable new sorting machinery (e.g. waste
 42 sorting robots), improve the performance of existing machinery (e.g. optical sorters [6]), and allow
 43 automatic quality control of the MRFs’ output. Unfortunately, the research community is lacking the
 44 gold-standard in-the-wild datasets to train and evaluate the classification and segmentation algorithms
 45 for industrial waste sorting. While several companies do development on this subject (e.g. [7, 8, 9]),
 46 they keep their dataset private, and the few existing open-source datasets [10, 4, 3, 2] are very limited
 47 in data amount and/or generated in uncluttered environments, not representing the complexity of the
 48 domain (see Figure 1). In this paper, we propose a first large-scale in-the-wild waste detection dataset
 49 `ZeroWaste` that is specifically designed for the industrial waste detection. `ZeroWaste` is a dataset
 50 that is fundamentally different from the popular detection and segmentation benchmarks: high level
 51 of clutter, visual diversity of the foreground and background objects that are often severely deformed,
 52 as well as a fine-grained difference between the object classes (e.g. brown paper vs. cardboard, soft
 53 vs. rigid plastic) – all these aspects pose a unique challenge for the automated vision. We envision
 54 that our open-access dataset will enable computer vision and robotics communities to develop more
 55 robust and data-efficient algorithms for object detection, robotic grasping and other related problems.
 56 Our contributions can be summarized as follows:

- 57 1. We propose the first fully-annotated `ZeroWaste-f` dataset specifically designed for industrial
 58 waste object detection. The proposed `ZeroWaste-f` dataset contains video frames from a real
 59 MRF conveyor belt densely annotated with instance segmentation and proposes a challenging
 60 real-life computer vision problem of detecting highly deformable objects in severely cluttered
 61 scenes. In addition to the fully annotated frames from `ZeroWaste-f` set, we include the
 62 unlabeled `ZeroWaste-s` set for semi-supervised learning.
- 63 2. We introduce a novel before-after data collection setup and propose the `ZeroWaste-w` dataset
 64 for binary classification of frames before and after the collection of target objects. This binary



Figure 2: Examples of images (left) and the corresponding polygon annotation (right) of the proposed ZeroWaste dataset. At the end of this conveyor belt, only paper objects must remain. Therefore, we annotated the objects of four material types that should be removed from the conveyor belt as foreground: soft plastic, rigid plastic, cardboard and metal. The background includes the conveyor belt and paper objects. Severe clutter and occlusions, high variability of the foreground object shapes and textures, as well as severe deformations of objects usually not present in other segmentation datasets, make this domain very challenging for object detection. More examples of our annotated data can be found in Section B.3 of the Appendix (best viewed in color).

65 classification setup allows much cheaper data annotation and allows further development of
 66 weakly supervised segmentation and detection methods.
 67 3. We implement the fully-supervised detection and segmentation baselines for the ZeroWaste-
 68 f dataset and semi- and weakly-supervised baselines for ZeroWaste- s and ZeroWaste-
 69 w datasets. Our experimental results show that popular detection and segmentation methods
 70 struggle to generalize to our proposed data, which indicates a challenging nature of our in-the-wild
 71 dataset and suggests that more robust and data-efficient methods must be developed to solve the
 72 waste detection problem.

73 2 Related Work

74 **Detection and Segmentation Datasets** Many datasets for image segmentation have been proposed
 75 with the goal of densely recognizing general objects and “stuff” in image scenes like street view [11,
 76 12, 13], natural scenes [14, 15, 16, 17, 18], and indoor spaces [19, 20, 21]. Yet, few of them have
 77 been designed for the more challenging vision task required in automated waste recycling, aiming to
 78 densely identify and segment deformable recyclable materials, many of which look very similar to
 79 each other, from a highly cluttered background [6]. Several related datasets have been proposed that
 80 contain only image-level labels. For example, *Portland State University Recycling* [4] consists of
 81 11500 labeled images of five common recyclable types: box-board, glass bottles, soda cans, crushed
 82 soda cans and plastic bottles. Similarly, *Stanford TrashNet* [10] presents 400 images containing a
 83 single waste object from six predefined classes. Though beneficial for image-level classification
 84 in well-defined conditions, images of in these two datasets have very simple background and do
 85 not apply to waste object localization. To enable localization tasks, *Labeled Waste in the Wild* [3]
 86 annotated bounding boxes for objects of 20 classes in 1002 food tray photos. *Annotation in Context*
 87 (*TACO*) [2] went one step further by densely annotating 60 litter objects from 1500 images. Yet
 88 TACO contains deliberately collected outdoor scenes with one or a few foreground objects that are
 89 rarely occluded, which makes it less practical for materials recovery scenarios. In contrast, our
 90 ZeroWaste was collected from the front lines of a waste sorting plant where the collected objects

91 are frequently severely deformed and occluded, which makes both detection and segmentation a
92 significantly more challenging and practical task.

93 **Detection and Segmentation Methods** Image segmentation is an essential component in robotic
94 systems like automated waste sorters [6], as it partitions images into multiple regions or objects
95 suitable for grasping. Image segmentation can be formulated as a task that classifies each pixel into
96 a set of labels [22]. Recent semantic segmentation models [23, 24, 25, 26] have achieved state-of-
97 the-art performance for recognizing general object/stuff classes from natural scene images. Instance
98 segmentation [27, 28, 29, 30] works by further consider instance identity for objects. Representative
99 frameworks like MaskRCNN [31] effectively detect objects in images and simultaneously generate
100 high-quality masks, which enables efficient interaction between robots and target objects. Yet due to
101 their data-hungry nature, these methods rely on large volumes of annotated data for training, which
102 can be challenging and expensive, especially in specialized application scenarios [32]. Recycling
103 annotation in particular requires expert labelers and is thus even more costly. Semi-supervised
104 segmentation methods have been proposed to address such limitations by jointly learning from
105 both annotated and unannotated images [33, 34, 35, 36, 37, 38]. Weakly-supervised segmentation
106 methods exploit annotations that are even easier to obtain, e.g. image-level tags [39, 40, 41]. These
107 methods typically utilize class activation maps (CAM) [42] to select the most discriminative regions,
108 which are later used as pixel-level supervision for segmentation networks [43, 44, 45]. All these
109 advanced segmentation models are trained on general-purpose data, and applying them to waste
110 sorting scenarios presents challenges like domain shift. To study the effectiveness of existing models
111 and enable further improvement for the waste sorting task, we test our proposed ZeroWaste with
112 previous state-of-the-art methods and report their performance as baselines.

113 3 ZeroWaste Dataset

114 In this section, we describe our ZeroWaste- f dataset for fully supervised detection and evaluation,
115 unlabeled ZeroWaste- s data for semi-supervised learning and ZeroWaste- w dataset of images
116 before and after the removal of target objects for weakly supervised detection. The datasets are
117 licensed under the Creative Commons Attribution-NonCommercial 4.0 International License [46].
118 The MRF at which the data was collected agreed to release the data for any non-commercial purposes
119 and decided to remain unacknowledged.

120 **Data Collecion and Pre-processing** The data was collected from a high-quality paper conveyor
121 of a single stream recycling facility in Massachusetts. The sorting operation on this conveyor aims to
122 keep high quality paper and consider anything else as contaminants including non-paper items (*e.g.*
123 metal, plastic, brown paper, cardboard, boxboard). We collected data during the regular operation
124 of the MRF using two compact recording installations at the start and end of the conveyor belt (see
125 Fig. 3, right), that is, footage is captured simultaneously both at the unsorted and sorted sections of
126 the same conveyor. The recording apparatus is designed to fit the constraints of the facility: In order
127 not to disrupt the MRF operation and be able to work in confined spaces available near the conveyor
128 the recording platform needs to be compact, non-intrusive (to the workers), and portable (easy to
129 move, battery-powered). Note that the cameras are not directly mounted on the conveyor but to a
130 stand-alone platform, to reduce vibrations transmitted to the cameras. Additional considerations are
131 made (see Figure 3, center): (1) Damping pads are installed to counter the ground vibrations of the
132 heavy machinery and reduce vibrations on the camera even further; (2) Weighted bases lower the
133 center of mass to keep the apparatus stable.

134 We used the GoPro Hero 7 for RGB footage, and we additionally collected the the near-infrared
135 (NIR) footage simultaneously with the RGB footage using the MAPIR Survey3W NIR camera for the
136 future work (specifically, it captures at a wavelength of 850 nm). The cameras in their encasings meet
137 both the portability and ruggedness requirements. To maintain consistent lighting, two LitraTorch 2.0
138 portable lamps are installed with a light diffuser. This softens the light and spreads it more evenly in
139 the scene. Both cameras were installed at around 100 cm above the conveyor, and the light sources at
140 around 80 cm. Sequences of twelve videos of total length of 95 minutes and 14 seconds with FPS
141 120 and size 1920×1080 were collected and processed. The preprocessing of the collected data
142 involved the following steps:

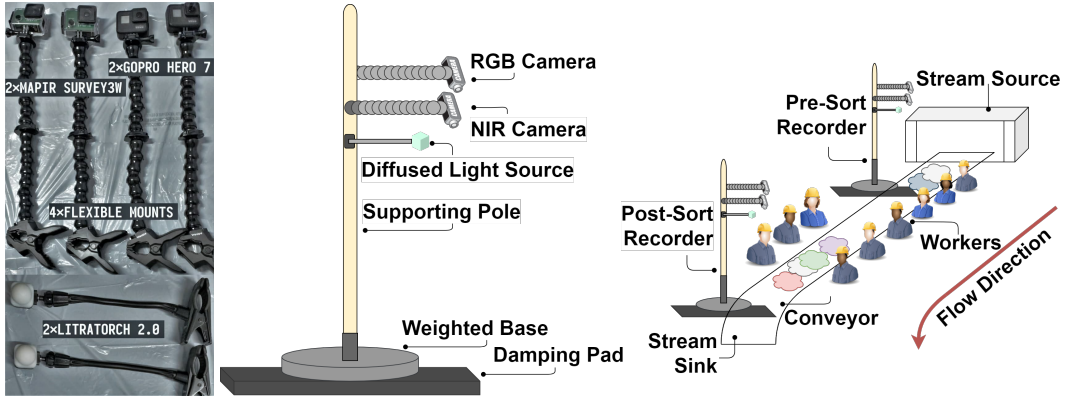


Figure 3: The footage recording setup is designed to fit the constraints of the facility environment. **Left:** The specific cameras and lamps used. **Center:** Assembly of each recording apparatus. **Right:** Layout of the recording setup in the recycling environment.



Figure 4: **Left:** example of an image from *ZeroWaste-f* dataset. **Right:** the corresponding ground truth instance segmentation. Expert training and common sense knowledge are required to distinguish between the cardboard object on the left (red circle) and the brown paper on the right (blue circle), as they are visually very similar but differ in thickness and rigidity (*best viewed in color*).

- 143 1. Rotation and cropping. The frames were rotated so that the conveyor belt is parallel to the
- 144 frame borders and cropped to remove the regions outside the conveyor belt. We ensured that any
- 145 personal information or identifiable footage of the workers at the conveyor belt was excluded
- 146 from our data.
- 147 2. Optical distortion. We removed the distortion [47] using the OpenCV [48] library to compensate
- 148 for the fish-eye effect caused by the proximity of the cameras to the conveyor belt.
- 149 3. Deblurring. We used the SRN-Deblur [49] method to remove motion blur resulting from the
- 150 fast-moving conveyor belt. According to our visual inspection, SRN-Deblur achieves satisfactory
- 151 deblurring and does not introduce the undesired artifacts that usually appear when classical
- 152 deconvolution-based methods are used.
- 153 4. Subsampling. We sampled every tenth frame from the video to avoid redundancy.

154 The illustration of the original frames shot at the beginning of the conveyor belt and the corresponding
 155 preprocessing results can be found on Figure 8 in Section B.3 of the Appendix.

156 **Densely Annotated ZeroWaste-f and Unlabeled ZeroWaste-s Datasets** The fully anno-
 157 tated *ZeroWaste-f* dataset consists of 1874 frames sampled from the processed videos and the
 158 corresponding ground truth polygon segmentation. We used the open-source CVAT [50] annotation
 159 toolkit to manually collect the polygon annotations of objects of four material types: cardboard, soft
 160 plastic, rigid plastic and metal. We chose this set of class labels following the MRF’s guidelines
 161 for the workers to collect cardboard, plastic and metal into separate bins, as well as the fact that
 162 grasping of rigid and non-rigid objects might require the use of fundamentally different kinds of
 163 robotic systems. The polygon annotation was performed according to the following set of rules:

Split	#Images	Carboard	Soft Plastic	Rigid Plastic	Metal	#Objects
Train	1245	4038	1550	460	114	6162
Validation	312	795	310	195	24	1324
Test	317	1216	466	242	53	1977
Unlabeled	6212	-	-	-	-	-
Total	8086	6049	2326	897	191	9463

Table 1: Statistics of the training, validation and test splits of our ZeroWaste-*f* dataset *w.r.t.* the number of labeled objects, and the additional unlabeled ZeroWaste-*s* set of images for semi-supervised learning.

- 164 1. Objects of four material types were annotated as foreground: cardboard (including parcel pack-
165 ages, boxboard such as cereal boxes and other carton food packaging), soft plastic (*e.g.* plastic
166 bags, wraps), rigid plastic (*e.g.* food containers, plastic bottles) and metal (*e.g.* metal cans). Paper
167 objects were treated as background.
- 168 2. The entire object must be within the corresponding polygon.
- 169 3. If an object is partially occluded and separate parts are visible, we annotated them as separate
170 objects.

171 Each annotated video frame was validated by an independent reviewer to pass the standards above
172 (see Figure 2). Both the annotation and the review process were performed by the students and
173 researchers with a computer science background specifically trained to perform the annotation.
174 We did not delegate the annotation to the crowd-sourcing platforms, such as Amazon Mechanical
175 Turk [51], due to the complexity of the domain that requires expert knowledge to be able to detect
176 and correctly classify the foreground objects (see the illustration on Figure 4). The estimated average
177 cost of the annotation and review is about 12.5 minutes per frame. The dataset was split into training,
178 validation and test splits and stored in the widely used MS COCO [18] format for object detection and
179 segmentation using the open-source Voxel51 toolkit [52]. Please refer to Table 1 for more details about
180 the class-wise statistics of all splits. In addition to the fully annotated ZeroWaste-*f* examples,
181 we provide 6212 unlabeled images that can be used to refine the detection using semi-supervised or
182 self-supervised learning methods. We refer to this unlabeled set of images as ZeroWaste-*s* data
183 later on in this paper.

184 **ZeroWaste-*w* Dataset for Binary Classification** We leverage the videos taken of the conveyer
185 belt before and after the removal of the foreground objects to create a weakly-supervised
186 ZeroWaste-*w* dataset. This dataset contains 1202 frames with the foreground objects (*before*
187 class) and 1208 frames without the foreground objects (*after* class). One advantage of such a setup is
188 that it is relatively cheap to acquire the ground truth labels (only an image-level inspection is required
189 to ensure there are no false negatives in the *after* class subset). The ZeroWaste-*w* dataset is
190 specifically collected to be used in the weakly-supervised setup and is meant to provide an alternative
191 and more data-efficient solution to the problem. The ground truth instance segmentation is available
192 for all images of the *before* class as it overlaps with the ZeroWaste-*f* dataset. Please see Figure 5
193 for an illustration of the ZeroWaste-*w* examples.

194 4 Experiments

195 In this section, we provide baseline results for our proposed ZeroWaste dataset. We perform
196 fully supervised instance and semantic segmentation on ZeroWaste-*f* using the most widely
197 used Mask R-CNN [31] and DeepLabV3+ [53] respectively. We also perform fully- and semi-
198 supervised semantic segmentation on ZeroWaste-*s* using the CCT [33] method, and report the
199 initial segmentation quality of CAMs produced by a classifier trained on ZeroWaste-*w* dataset as
200 a weakly-supervised baseline. The implementation of our experiments and the detailed description of
201 the experimental setup are available at <https://github.com/dbash/zerowaste>.

202 4.1 Object Detection

203 **Experiments with COCO-pretrained Networks** It has been shown that pretraining the model
204 on a large-scale dataset, such as MS COCO [18], improves generalization and helps to prevent
205 severe overfitting in case when the target dataset is relatively small [54, 55, 56]. Therefore, in our



Figure 5: We installed two stationary cameras above the conveyor belt: one at the beginning of the line and another one at the end. At this particular conveyor belt, workers are asked to remove objects of any material other than paper, such as cardboard, plastic and metal. Therefore, the footage collected from the beginning of the line contains the “foreground” objects that need to be removed, and the frames from the end of the conveyor belt are supposed to only contain the “background” paper objects. We used this setup as a foundation of our `ZeroWaste-w` dataset.

206 first experiments, we used the initialized the model with weights learned on COCO and further
 207 finetuned it with our `ZeroWaste-f` dataset. We used a standard implementation of the popular
 208 Mask R-CNN with ResNet-50 [57] backbone provided in the popular Detectron2 [58] library in
 209 all of the experiments. The model was finetuned for 40000 iterations on the training set of our
 210 `ZeroWaste-f` dataset on a single Geforce GTX 1080 GPU with batch size 8. To compensate for
 211 a relatively small number of examples in the training set and to avoid overfitting, we leveraged heavy
 212 data augmentation, including random rotation and cropping, adjustment of brightness and hue, *etc.*
 213 We report the experimental results in Table 2 (COCO \rightarrow `ZeroWaste` section). A more detailed
 214 description of the results can be found in Section B.1 of Appendix.

215 **Experiments with TACO-pretrained Mask RCNN** In the next set of experiments, we utilize the
 216 TACO dataset for waste detection in the outdoor scenes distributed under Attribution 4.0 International
 217 (CC BY 4.0) license. We trained Mask R-CNN for 40000 epochs on the modified TACO dataset
 218 with the material-based labels (cardboard, soft plastic, rigid plastic, metal and other) initialized with
 219 weights from MS COCO. We then finetuned the model on the training set of `ZeroWaste-f` data
 220 and report the results in Table B.1 (TACO \rightarrow `ZeroWaste` section).

221 **Results** The experimental results with Mask RCNN indicate severe overfitting to the training data,
 222 hence the model fails to generalize to the unseen examples. The model pretrained on the TACO
 223 dataset performs poorly on both TACO and `ZeroWaste-f` datasets, which shows that, despite its
 224 remarkable efficiency on the large-scale datasets with natural scenes, such as MS COCO or Pascal
 225 VOC [59], Mask RCNN cannot generalize to our relatively small, extremely cluttered data with very
 226 diverse deformable objects. Recalling the history of success with other complex segmentation and
 227 detection datasets (*e.g.* from mIoU 57% in 2015 [60] to 84% in 2020 [61] on CityScapes [11], or
 228 from 51.6% in 2014 [62] to 90% in 2020 [63] on PASCAL VOC 2012 [59]), and knowing that the
 229 task *can* be solved by humans with a little training, we believe that the computer vision community
 230 will eventually come up with efficient methods for this challenging task.

231 4.2 Semantic Segmentation

232 **Fully supervised experiments** We used the state-of-the-art DeeplabV3+ model as a fully-supervised
 233 semantic segmentation baseline for our dataset. DeeplabV3+ is an efficient segmentation model
 234 that combines the atrous convolutions to extract the features in multiple scales, and an encoder-

	TACO \rightarrow ZeroWaste			COCO \rightarrow ZeroWaste		
	AP	AP50	AP75	AP	AP50	AP75
Train	39.11	54.77	44.58	62.55	81.59	71.59
Validation	15.86	28.83	16.37	14.99	23.62	16.09
Test	14.55	25.9	14.81	14.79	25.94	14.82

Table 2: Instance segmentation results of Mask R-CNN pretrained on TACO dataset (**left**) and MS COCO dataset (**right**). The model pretrained on MS COCO overfits to the training split, while pretraining on TACO dataset significantly reduces overfitting but does not yield a significant improvement in detection accuracy on the validation and test sets. Please refer to Tables 6 in the Appendix for class-wise results.

235 decoder paradigm to gradually sharpen the object boundary using the intermediate features. As in the
 236 detection experiments, we used a standard implementation of DeeplabV3+ from Detectron2 library.
 237 We used the model with ResNet-101 backbone with three 3×3 convolutions instead of the first
 238 7×7 convolution that was pretrained on Cityscapes dataset [11]. We froze the first three stages
 239 of the backbone (convolution and two first residual block groups) and finetuned the model on the
 240 training set of *ZeroWaste-f* for 10000 iterations with starting learning rate 0.01 and batch size 40
 241 on a single GPU RTX A6000 which took approximately 14 hours. As in the previous experiments,
 242 we augmented the data extensively to prevent overfitting. The results of our experiments on all
 243 *ZeroWaste-f* splits can be found on the Table 3.

244 **Semi-supervised experiments** For a semi-supervised segmentation baseline, we used an official
 245 implementation of Cross-Consistency Training CCT [33] method. CCT uses a shared encoder and
 246 several auxiliary decoders each of which performs various augmentations, such as spatial dropout,
 247 random noise, cutout of object regions *etc.*, and a cross-entropy-based loss to force the unlabeled
 248 predictions to be consistent across all decoders. Since CCT uses a different backbone architecture
 249 from DeeplabV3+, we first trained CCT on the labeled *ZeroWaste-f* data only for comparison
 250 with the semi-supervised setting. We used the same default hyperparameters reported in the paper
 251 for both supervised and semi-supervised experiments (the exact configuration can be found in our
 252 project). We report the mean Intersection over Union (mIoU) as well as mean pixel accuracy for both
 253 setups in Table 3, and more details can be found in Section B.2 of the Appendix.

254 **Weakly-supervised baseline** As a baseline for weakly-supervised segmentation, we trained a binary
 255 classifier on the before and after collection frames of the *ZeroWaste-w* dataset. We used a standard
 256 Pytorch [64] implementation of ResNet50 [57] pretrained on ImageNet [65] for our classifier, and
 257 trained it for 5 epochs with learning rate 5×10^{-4} using the binary cross-entropy loss. The resulting
 258 classifier obtained over 98% accuracy on the test set. We then used RISE [66], a black-box saliency
 259 generating technique, to extract the class activation maps (CAMs). RISE masks the input image with
 260 a set of random binary masks and returns the linear combination of the resulting CAMs weighted with
 261 the corresponding masks. The maps generated by RISE are then normalized and thresholded with
 262 0.621 that results in highest mIoU on the training set. For comparison, we computed the mean pixel
 263 accuracy and mIoU on randomly generated masks with the probability of each pixel belonging to the
 264 foreground class equal to the average fraction of the foreground pixels in the *ZeroWaste-w* dataset
 265 14.9% and report these results in Table 3. The visualization of the resulting CAMs can be found in
 266 Figure 9 in Section B.3 of the Appendix.

267 **Results** Experimental results in Table 3 indicate that our *ZeroWaste* dataset proposes a challenging
 268 semantic segmentation task with an unusual for the standard segmentation datasets level of clutter,
 269 diversity of the foreground objects and, at the same time, their visual similarity with the background
 270 objects (all methods often tend to mistake the paper objects for cardboard and vice versa, and have
 271 a hard time distinguishing between soft and rigid plastic objects). The semi-supervised learning
 272 results indicate that the unlabeled examples from the *ZeroWaste-s* subset do not significantly help
 273 CCT improve the overall segmentation quality. As seen from the class-wise segmentation results
 274 on Table 8 in Section B.2 of Appendix, additional training of CCT with unlabeled data results in
 275 higher segmentation accuracy of the most frequent classes (*e.g.* cardboard and background), but
 276 degrades the performance on the objects of the rare classes (*e.g.* metal). Additionally, the binary
 277 classification results show that a simple CAM-based approach with cheap *ZeroWaste-w* data
 278 provides meaningful localization cues that can be further used for weakly- and semi-supervised
 279 segmentation.

	Supervision	Train		Validation		Test	
		mIoU	Pixel Acc.	mIoU	Pixel Acc.	mIoU	Pixel Acc.
<i>Random</i>	none	7.2	74.7	7.2	75.3	8.4	71.8
<i>CAM</i>	weak	15.7	43.9	16.3	47.5	18.6	43.2
<i>CCT semi</i>	semi	61.2	97.4	29.40	83.3	30.0	83.6
<i>CCT</i>	full	65.38	97.9	29.80	83.4	29.20	81.2
<i>DeeplabV3+</i>	full	88.5	98.19	40.16	91.23	39.06	88.47

Table 3: Results of CAMs produced by RISE [66] with a binary classifier trained on `ZeroWaste-w` before and after frames, CCT [33] trained only using the `ZeroWaste-f`, CCT trained with `ZeroWaste-f` and `ZeroWaste-s`, and DeepLabV3+ [53] on our `ZeroWaste-f` dataset. Results indicate that 1) severe overfitting occurs in the supervised scenario; 2) unlabeled `ZeroWaste-s` images do not significantly improve the segmentation quality of CCT and 3) the binary classifier trained on `ZeroWaste-w` provides plausible localization guidance that can serve as cues for weakly-supervised segmentation. Please refer to Tables 7 and 8 for class-wise segmentation results and Figure 7 in the Appendix for confusion matrices on all splits.

280 5 Impact and Limitations of ZeroWaste

281 **Machine Learning Research** `ZeroWaste` provides a gold standard for the evaluation of different
282 waste sorting methods. It will catalyze research in the areas of fully, semi, and weakly supervised
283 segmentation, data-efficient learning and domain adaptation. Our dataset provides a real-world
284 application that is significantly more challenging than the previously used benchmarks for these tasks.

285 **Robotics Research** This dataset will enable the development of robotic manipulation algorithms for
286 waste sorting. It will facilitate research in object picking algorithms that can work with extremely
287 cluttered scenes using realistic segmentation polygons. Integrating high-level reasoning about object
288 classes and properties (e.g. hard/soft materials) to the picking algorithm will provide novel research
289 avenues and can significantly boost the picking accuracy.

290 **Limitations and Future Directions** Despite the fact that `ZeroWaste` is to the date the largest
291 public dataset for waste detection and segmentation, it is still smaller than the standard large-scale
292 benchmarks due to the fact that the annotation process for this domain is very expensive. For this
293 reason, state-of-the-art detection and segmentation methods tend to overfit to the training data and
294 therefore do not generalize well to the unseen examples. As future work, we plan to increase the
295 diversity of our dataset by using synthetic-to-real domain adaptation and other data augmentation
296 techniques. Another important future direction is to utilize visual signals of other modalities, e.g.
297 near infrared footage that can be especially useful for distinguishing different material types.

298 **Societal Impact** This paper is a part of a collaboration project that investigates the implications of
299 deploying new AI and Robotics algorithms to MRFs [67]. We believe that human-robot collaboration
300 is essential for more efficient computer-aided recycling, quality control of the sorting process, as well
301 as in establishing safer work conditions for the MRF workers (e.g. by detecting dangerous waste items,
302 such as sharp or explosive objects). This dataset can potentially be used to develop fully-automated
303 MRFs with waste sorting robots, which may compromise the financial security of the MRF workers.
304 However, after consulting with experts, we found that such fully-automated solutions would be far
305 from sufficient to meet the contamination levels required in recycling, especially considering the
306 complex, cluttered and varying nature of the waste stream. Given that only a small portion of the
307 recyclable waste is currently getting recycled, achieving an efficient human-robot collaboration has a
308 potential to solve the pressing problem of water and soil pollution.

309 6 Conclusion

310 This work introduces the largest public dataset for waste detection. `ZeroWaste` is designed
311 as a benchmark for training and evaluation of fully, weakly, and semi-supervised detection and
312 segmentation methods, and can be directly used for a broader category of tasks including transfer
313 learning, domain adaptation and label-efficient learning. We provide baseline results for the most
314 popular fully, weakly, semi-supervised, and transfer learning techniques. Our results show that
315 current state-of-the-art detection and segmentation methods cannot efficiently handle this complex
316 in-the-wild domain. We anticipate that our dataset will motivate the computer vision community to
317 develop more data-efficient methods applicable to a wider range of real-world problems.

References

- 318
- 319 [1] National overview: Facts and figures on materials, wastes and recycling. [EB/OL].
- 320 [2] Pedro F Proença and Pedro Simões. Taco: Trash annotations in context for litter detection.
321 *arXiv preprint arXiv:2003.06975*, 2020.
- 322 [3] Joao Sousa, Ana Rebelo, and Jaime S Cardoso. Automation of waste sorting with deep learning.
323 In *2019 XV Workshop de Visão Computacional (WVC)*, pages 43–48. IEEE, 2019.
- 324 [4] Anthony Martin. Recycling image classification. [EB/OL]. [http://web.cecs.pdx.edu/
325 ~singh/rcyc-web/index.html/](http://web.cecs.pdx.edu/~singh/rcyc-web/index.html/) Accessed May 22, 2021.
- 326 [5] Silpa Kaza, Lisa C. Yao, Perinaz Bhada-Tata, and title = Van Woerden, Frank.
- 327 [6] Sathish Paulraj Gundupalli, Subrata Hait, and Atul Thakur. A review on automated sorting of
328 source-separated municipal solid waste for recycling. *Waste management*, 60:56–74, 2017.
- 329 [7] AMP Robotics. <https://www.amprobotics.com/>. Accessed: 2020-05-30.
- 330 [8] Waste-Robotics. <https://wasterobotic.com/>. Accessed: 2020-05-30.
- 331 [9] Zen Robotics. <https://zenrobotics.com/>. Accessed: 2020-05-30.
- 332 [10] Mindy Yang and Gary Thung. Classification of trash for recyclability status. *CS229 Project
333 Report*, 2016, 2016.
- 334 [11] Marius Cordts, Mohamed Omran, Sebastian Ramos, Timo Rehfeld, Markus Enzweiler, Rodrigo
335 Benenson, Uwe Franke, Stefan Roth, and Bernt Schiele. The cityscapes dataset for semantic
336 urban scene understanding. In *CVPR*, 2016.
- 337 [12] Fisher Yu, Haofeng Chen, Xin Wang, Wenqi Xian, Yingying Chen, Fangchen Liu, Vashisht
338 Madhavan, and Trevor Darrell. Bdd100k: A diverse driving dataset for heterogeneous multitask
339 learning. In *CVPR*, 2020.
- 340 [13] Gabriel J Brostow, Jamie Shotton, Julien Fauqueur, and Roberto Cipolla. Segmentation and
341 recognition using structure from motion point clouds. In *ECCV*, 2008.
- 342 [14] M. Everingham, L. Van Gool, C. K. I. Williams, J. Winn, and A. Zisserman. The pascal visual
343 object classes (voc) challenge. *IJCV*, 88:303–338, 2010.
- 344 [15] Holger Caesar, Jasper Uijlings, and Vittorio Ferrari. Coco-stuff: Thing and stuff classes in
345 context. In *CVPR*, 2018.
- 346 [16] Roozbeh Mottaghi, Xianjie Chen, Xiaobai Liu, Nam-Gyu Cho, Seong-Whan Lee, Sanja Fidler,
347 Raquel Urtasun, and Alan Yuille. The role of context for object detection and semantic
348 segmentation in the wild. In *CVPR*, 2014.
- 349 [17] Bolei Zhou, Hang Zhao, Xavier Puig, Sanja Fidler, Adela Barriuso, and Antonio Torralba.
350 Scene parsing through ade20k dataset. In *CVPR*, 2017.
- 351 [18] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan,
352 Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *ECCV*,
353 2014.
- 354 [19] Nathan Silberman, Derek Hoiem, Pushmeet Kohli, and Rob Fergus. Indoor segmentation and
355 support inference from rgb-d images. In *ECCV*, 2012.
- 356 [20] Angela Dai, Angel X Chang, Manolis Savva, Maciej Halber, Thomas Funkhouser, and Matthias
357 Nießner. Scannet: Richly-annotated 3d reconstructions of indoor scenes. In *CVPR*, 2017.
- 358 [21] Liang-Chieh Chen, Yukun Zhu, George Papandreou, Florian Schroff, and Hartwig Adam.
359 Encoder-decoder with atrous separable convolution for semantic image segmentation. In *ECCV*,
360 2018.

- 361 [22] Shervin Minaee, Yuri Y Boykov, Fatih Porikli, Antonio J Plaza, Nasser Kehtarnavaz, and
362 Demetri Terzopoulos. Image segmentation using deep learning: A survey. *IEEE Trans. on*
363 *PAMI*, 2021.
- 364 [23] Hengshuang Zhao, Jianping Shi, Xiaojuan Qi, Xiaogang Wang, and Jiaya Jia. Pyramid scene
365 parsing network. In *CVPR*, 2017.
- 366 [24] Liang-Chieh Chen, George Papandreou, Florian Schroff, and Hartwig Adam. Rethinking atrous
367 convolution for semantic image segmentation. *arXiv preprint arXiv:1706.05587*, 2017.
- 368 [25] Yuhui Yuan, Xilin Chen, and Jingdong Wang. Object-contextual representations for semantic
369 segmentation. *ECCV*, 2020.
- 370 [26] Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining
371 Guo. Swin transformer: Hierarchical vision transformer using shifted windows. *CVPR*, 2021.
- 372 [27] Shu Liu, Lu Qi, Haifang Qin, Jianping Shi, and Jiaya Jia. Path aggregation network for instance
373 segmentation. In *CVPR*, 2018.
- 374 [28] Jifeng Dai, Yi Li, Kaiming He, and Jian Sun. R-fcn: Object detection via region-based fully
375 convolutional networks. In *NIPS*, 2016.
- 376 [29] Enze Xie, Peize Sun, Xiaoge Song, Wenhai Wang, Xuebo Liu, Ding Liang, Chunhua Shen, and
377 Ping Luo. Polarmask: Single shot instance segmentation with polar representation. In *CVPR*,
378 2020.
- 379 [30] Youngwan Lee and Jongyoul Park. Centermask: Real-time anchor-free instance segmentation.
380 In *CVPR*, 2020.
- 381 [31] Francisco Massa and Ross Girshick. maskrcnn-benchmark: Fast, modular reference imple-
382 mentation of Instance Segmentation and Object Detection algorithms in PyTorch. <https://github.com/facebookresearch/maskrcnn-benchmark>, 2018. Accessed: [Insert date
383 here].
- 384 [32] Adela Barriuso and Antonio Torralba. Notes on image annotation. *arXiv preprint*
385 *arXiv:1210.3448*, 2012.
- 386 [33] Yassine Ouali, Celine Hudelot, and Myriam Tami. Semi-supervised semantic segmentation
387 with cross-consistency training. In *The IEEE/CVF Conference on Computer Vision and Pattern*
388 *Recognition (CVPR)*, June 2020.
- 389 [34] Sudhanshu Mittal, Maxim Tatarchenko, and Thomas Brox. Semi-supervised semantic segmen-
390 tation with high-and low-level consistency. *IEEE Trans. on PAMI*, 2019.
- 391 [35] Robert Mendel, Luis Antonio de Souza, David Rauber, João Paulo Papa, and Christoph Palm. Semi-supervised segmentation based on error-correcting supervision. In *ECCV*, 2020.
- 392 [36] Jongmok Kim, Jooyoung Jang, and Hyunwoo Park. Structured consistency loss for semi-
393 supervised semantic segmentation. *arXiv preprint arXiv:2001.04647*, 2020.
- 394 [37] Geoff French, Timo Aila, Samuli Laine, Michal Mackiewicz, and Graham Finlayson. Semi-
395 supervised semantic segmentation needs strong, high-dimensional perturbations. *arXiv preprint*
396 *arXiv:1906.01916*, 2019.
- 397 [38] Liang-Chieh Chen, Raphael Gontijo Lopes, Bowen Cheng, Maxwell D Collins, Ekin D Cubuk,
398 Barret Zoph, Hartwig Adam, and Jonathon Shlens. Semi-supervised learning in video sequences
399 for urban scene segmentation. *ECCV*, 2020.
- 400 [39] Jiwoon Ahn and Suha Kwak. Learning pixel-level semantic affinity with image-level supervision
401 for weakly supervised semantic segmentation. In *CVPR*, 2018.
- 402 [40] Alexander Kolesnikov and Christoph H Lampert. Seed, expand and constrain: Three principles
403 for weakly-supervised image segmentation. In *ECCV*, 2016.
- 404
405

- 406 [41] Deepak Pathak, Philipp Krahenbuhl, and Trevor Darrell. Constrained convolutional neural
407 networks for weakly supervised segmentation. In *ICCV*, 2015.
- 408 [42] Bolei Zhou, Aditya Khosla, Agata Lapedriza, Aude Oliva, and Antonio Torralba. Learning deep
409 features for discriminative localization. In *CVPR*, 2016.
- 410 [43] Yu-Ting Chang, Qiaosong Wang, Wei-Chih Hung, Robinson Piramuthu, Yi-Hsuan Tsai, and
411 Ming-Hsuan Yang. Weakly-supervised semantic segmentation via sub-category exploration. In
412 *CVPR*, 2020.
- 413 [44] Junsong Fan, Zhaoxiang Zhang, Chunfeng Song, and Tieniu Tan. Learning integral objects with
414 intra-class discriminator for weakly-supervised semantic segmentation. In *CVPR*, 2020.
- 415 [45] Yude Wang, Jie Zhang, Meina Kan, Shiguang Shan, and Xilin Chen. Self-supervised equivariant
416 attention mechanism for weakly supervised semantic segmentation. In *CVPR*, 2020.
- 417 [46] Creative commons attribution-noncommercial 4.0 international license. [EB/OL]. <http://creativecommons.org/licenses/by-nc/4.0/>
418 Accessed May 22, 2021.
- 419 [47] Duane C Brown. Decentering distortion of lenses. *Photogrammetric Engineering*, 1966.
- 420 [48] G. Bradski. The OpenCV Library. *Dr. Dobb's Journal of Software Tools*, 2000.
- 421 [49] Xin Tao, Hongyun Gao, Xiaoyong Shen, Jue Wang, and Jiaya Jia. Scale-recurrent network
422 for deep image deblurring. In *IEEE Conference on Computer Vision and Pattern Recognition*
423 *(CVPR)*, 2018.
- 424 [50] Boris Sekachev, Nikita Manovich, Maxim Zhiltsov, Andrey Zhavoronkov, Dmitry Kalinin, Ben
425 Hoff, TOsmanov, Dmitry Kruchinin, Artyom Zankevich, DmitriySidnev, Maksim Markelov,
426 Johannes222, Mathis Chenuet, a andre, telenachos, Aleksandr Melnikov, Jijoong Kim, Liron
427 Ilouz, Nikita Glazov, Priya4607, Rush Tehrani, Seungwon Jeong, Vladimir Skubriev, Sebastian
428 Yonekura, vugia truong, zliang7, lizhming, and Tritin Truong. *opencv/cvat: v1.1.0*, August
429 2020.
- 430 [51] Kevin Crowston. Amazon mechanical turk: A research tool for organizations and information
431 systems scholars. In *Shaping the future of ict research. methods and approaches*, pages 210–221.
432 Springer, 2012.
- 433 [52] B. E. Moore and J. J. Corso. Fiftyone. *GitHub. Note: https://github.com/voxel51/fiftyone*, 2020.
- 434 [53] Liang-Chieh Chen, Yukun Zhu, George Papandreou, Florian Schroff, and Hartwig Adam.
435 Encoder-decoder with atrous separable convolution for semantic image segmentation. In *ECCV*,
436 2018.
- 437 [54] Minyoung Huh, Pulkit Agrawal, and Alexei A Efros. What makes imagenet good for transfer
438 learning? *arXiv preprint arXiv:1608.08614*, 2016.
- 439 [55] Hoo-Chang Shin, Holger R Roth, Mingchen Gao, Le Lu, Ziyue Xu, Isabella Nogues, Jian-
440 hua Yao, Daniel Mollura, and Ronald M Summers. Deep convolutional neural networks for
441 computer-aided detection: Cnn architectures, dataset characteristics and transfer learning. *IEEE*
442 *transactions on medical imaging*, 35(5):1285–1298, 2016.
- 443 [56] Ross Girshick, Jeff Donahue, Trevor Darrell, and Jitendra Malik. Rich feature hierarchies for
444 accurate object detection and semantic segmentation. In *Proceedings of the IEEE conference on*
445 *computer vision and pattern recognition*, pages 580–587, 2014.
- 446 [57] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image
447 recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*,
448 pages 770–778, 2016.
- 449 [58] Yuxin Wu, Alexander Kirillov, Francisco Massa, Wan-Yen Lo, and Ross Girshick. Detectron2.
450 <https://github.com/facebookresearch/detectron2>, 2019.

- 451 [59] M. Everingham, S. M. A. Eslami, L. Van Gool, C. K. I. Williams, J. Winn, and A. Zisserman.
452 The pascal visual object classes challenge: A retrospective. *International Journal of Computer*
453 *Vision*, 111(1):98–136, January 2015.
- 454 [60] Vijay Badrinarayanan, Alex Kendall, and Roberto Cipolla. Segnet: A deep convolutional
455 encoder-decoder architecture for image segmentation. *IEEE transactions on pattern analysis*
456 *and machine intelligence*, 39(12):2481–2495, 2017.
- 457 [61] Bowen Cheng, Maxwell D Collins, Yukun Zhu, Ting Liu, Thomas S Huang, Hartwig Adam,
458 and Liang-Chieh Chen. Panoptic-deeplab: A simple, strong, and fast baseline for bottom-up
459 panoptic segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and*
460 *Pattern Recognition*, pages 12475–12485, 2020.
- 461 [62] Bharath Hariharan, Pablo Arbeláez, Ross Girshick, and Jitendra Malik. Simultaneous detection
462 and segmentation. In *European Conference on Computer Vision*, pages 297–312. Springer,
463 2014.
- 464 [63] Barret Zoph, Golnaz Ghiasi, Tsung-Yi Lin, Yin Cui, Hanxiao Liu, Ekin Dogus Cubuk, and
465 Quoc Le. Rethinking pre-training and self-training. *Advances in Neural Information Processing*
466 *Systems*, 33, 2020.
- 467 [64] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan,
468 Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. Pytorch: An imperative
469 style, high-performance deep learning library. *arXiv preprint arXiv:1912.01703*, 2019.
- 470 [65] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep
471 convolutional neural networks. *Advances in neural information processing systems*, 25:1097–
472 1105, 2012.
- 473 [66] Vitali Petsiuk, Abir Das, and Kate Saenko. Rise: Randomized input sampling for explanation
474 of black-box models. *arXiv preprint arXiv:1806.07421*, 2018.
- 475 [67] NSF Award Abstract # 1928506. Fw-htf-rl: Collaborative research: Shared autonomy for the
476 dull, dirty, and dangerous: Exploring division of labor for humans and robots to transform
477 the recycling sorting industry. [https://www.nsf.gov/awardsearch/showAward?AWD_ID=](https://www.nsf.gov/awardsearch/showAward?AWD_ID=1928506&HistoricalAwards=false)
478 [1928506&HistoricalAwards=false](https://www.nsf.gov/awardsearch/showAward?AWD_ID=1928506&HistoricalAwards=false). Accessed: 2020-05-30.