

A Benchmark Image Dataset for Industrial Tools

Cai Luo^{a,**}, Leijian Yu^b, Erfu Yang^b, Huiyu Zhou^c, Peng Ren^{d,**}

^athe Department of Mechanical and Electronic Engineering, China University of Petroleum (East China), Qingdao, China

^bthe Department of Design, Manufacture & Engineering Management, University of Strathclyde, Glasgow, UK

^cthe Department of Informatics, University of Leicester, Leicester, UK

^dthe Department of Information and Control Engineering, China University of Petroleum (East China), Qingdao, China

ABSTRACT

Robots and Artificial Intelligence (AI) play an increasingly important role in manufacture. One of the tasks is to identify tools in the scene so that the tools can be applied to different assembly purposes. In the AI community, many datasets have been generated and deployed to train robots to recognize individual items, however, these datasets are scene-specific and lack generic background. In this paper, we report our dataset contains photos of 8 objects types that would be easily recognized by qualified workers. This is achieved by gathering images of common tools in a typical factory. The ground truth categories of our dataset are manually labeled by experienced workers, which would be worthy evaluation tools for the intelligence industrial systems. The equipment used and the image collection process are discussed, along with the data format. The mean average precisions range from 64.37% to 78.20%, which bring the possibility for future improvement. The dataset is ideal to evaluate and benchmark view-point variant, vision-based control algorithm for industry robots. It is now public available from <https://github.com/tools-dataset/Industrial-Tools-Detection-Dataset>.

© 2019 Elsevier Ltd. All rights reserved.

1. Introduction

Every day, industry workers use variety tools for their daily duties, like cutting steel plate, tightening screw, hammering a nail, or measuring length, as shown in Fig.1. By virtue of training and memory, workers can effortlessly identify a tool and know its function. They are also able to choose suitable tools for different needs. While in the machine world, robots are still struggling to acquire the ability to pick correct instruments for assigned tasks through their visual sensors [28, 41, 20]. As robots like SCHAFT, Atlas, Valkyrie and REEM-C begin to manipulate standard tools and equipments commonly available in industrial environment, ranging from small screw drivers to full-size vehicles [38, 32, 6, 34]. The proliferation of AI embodied in robots increases the needs for these humanoid machines can work with their own hands, so they can take the tasks from repairing satellites to working in a remote factory without human intervention [21, 24]. It appears clear that for dealing with such complex scenarios, robust and



Fig. 1. Industrial tools. Variety industrial tools for workers' daily duties have been chosen for our dataset based on the purpose of evaluating view-point variant, vision-based control algorithms for industry robots.

efficient object detection algorithms are very important. Deep learning related methods has make great success in other field [37, 40, 3, 2, 35, 39, 15]. However, for deep learning methods, training datasets play the vital roles[4, 16]. So, in order to identify different tools successfully, specifically designed datasets are needed.

To advance object recognition research in industry, we in-

^{**}Corresponding author: tsai.lo.95@gmail.com (Cai Luo); pengren@upc.edu.cn (Ren Peng)

roduce a dataset for Industrial Tools Detection (ITD). It appears clear that it would bring great possibilities for robots to use a wide variety of instruments if they could distinguish these tools in factories or construction sites [23, 24, 33]. The dataset detailed in this paper is introduced to identify tools at the level of usages, and provide precise predictions for a robot to interact within the industry scenarios. Furthermore, the dataset is a challenging benchmark to evaluate view-point variant, vision-based control algorithms for industry robots.

The main contribution of this paper are as follows:

- We present a new large-scale object dataset, which consists of 8 object categories, 24 common industrial tools overall and multi distinct views of each tool. The dataset provides hand-labeled ground truth for more than 11,000 RGB images.
- We evaluate state-of-the-art object detection algorithms on ITD and define benchmark as baseline references for developing future new algorithms.
- Dataset and code from this work are available on-line at: <https://github.com/tools-dataset/Industrial-Tools-Detection-Dataset>.

2. Related Works

The purpose of choosing a suitable tool is to fulfill a task goal as quickly as possible [14]. The problem of picking and using tools has been widely studied in robotics, computer vision, artificial perception and psychology for many years and will be hot topics for the next decades as well. Many efforts have been dedicated to detect geometric characteristic of tools and how to handle the items correctly and firmly. They often assume prior information of objects shape and general location. Some of them also need the assistance of affordance labels or predefined markers to accomplish the manipulation tasks.

The creation of ground truth image and video datasets helped stimulated a flood of interest in the related areas. Large datasets like MS-COCO [25] is the de facto standard evaluation instrument for object detection. For the object categories classification, the PASCAL VOC [11, 18] and ImageNET [31, 10] are always in the datasets list of researchers. These datasets have proven to be very good performance test fields for computer vision algorithms in natural scenes.

In the field of tools detection, recognition and manipulation, datasets have been play a critical role as an algorithm assess. However, such successes have been slow to industrial field imagery due to the scarcity of optimal annotated datasets for tools in industrial environments. Unlike common daily objects, the collection and classification of industrial tools are much more difficult. Workers in the factories will need some special trainings in order to know the correct usage of tools [7]. They will need another several years to get the experiences to figure out how to choose the most suitable ones according to the tasks. Furthermore, the detection of industrial items are highly dependent on contextual information, which means the items in the datasets should be in their natural environments. Datasets like TAS [19], HRSC2016 [27] and DOTA [36] only contains large

items like vehicles, planes and ships that are difficult to manipulate by robots. Some pioneering works have grounded the tool handling in a constrained testing samples[7, 13]. Deep learning method has been applied by Ian et al. to solve the grasp problem by using a dataset which containing several daily tools. Kuan et al. proposed an affordance learning approach for tool manipulation through pre-selected objects. When it comes to general industrial tools, such as hammers, wrenches or saws, researches are normally depend on their own testing sets. All these datasets are short in the number of tools varieties, which prevent them from being widely usable.

Table 1. Comparison among MS-COCO, PASCAL VOC, ImageNET, TAS, HRSC2016, DOTA, Cornell grasping and ITD.

Datasets	Instances	Objects of interest
MS-COCO	123,287	Nature objects
PASCAL VOC	21,503	Nature objects
ImageNET	349,319	Nature objects
TAS	1319	Aerial targets
HRSC2016	2976	Aerial targets
DOTA	188,282	Aerial targets
Cornell grasping	1,035	Daily tools
ITD	11,000	Industrial tools

Our target is to simulate all possible situation of intelligent industrial systems. When collecting data, we gather the most common posture of the tools and place them in the location where they may found normally. Next, we analyze the properties of ITD in comparison to several other popular datasets. These include MS-COCO, PASCAL VOC 2012 and Cornell grasping dataset. Each of these datasets varies significantly in numbers of tools categories and quantities of images. MS-COCO was created to detect and segment of items occurring in their natural context. PASCAL VOC focuses on object detection in natural images. They both have at least 20 different categories, such as person, animals, aeroplane, chair and monitor. But none of them include the tools, especially industrial tools. Cornell grasping dataset has the largest number of categories in previous common tools datasets. The comparison results can be seen in Table.1. Note that ITD surpass Cornell grasping dataset not only in tools category numbers, but also in total number of tools, as shown in Fig.2.

In our datasets, we strive to collect images rich in classification, illumination and localization. ITD collected 24 daily industry tools. 8 categories are chosen, including cutting tools, fastener tools, adhesive tools, measuring tools, clamp tools, marker, polish tools and protection tools, as shown in Table.2. Fig.1 shows the examples of these tools. Compared with previous datasets, ITD can aid intelligent industrial systems specifically.

3. Industrial Tools Detection Dataset

This section presents how the ITD Dataset are selected. And what are the hardware and software used for the data collection are also described.

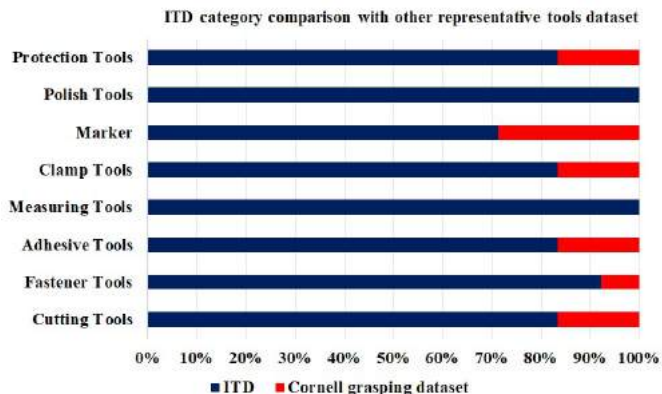


Fig. 2. Category Comparison. We perform an evaluation comparison between ITD and Cornell grasping dataset and responding quantity of items.

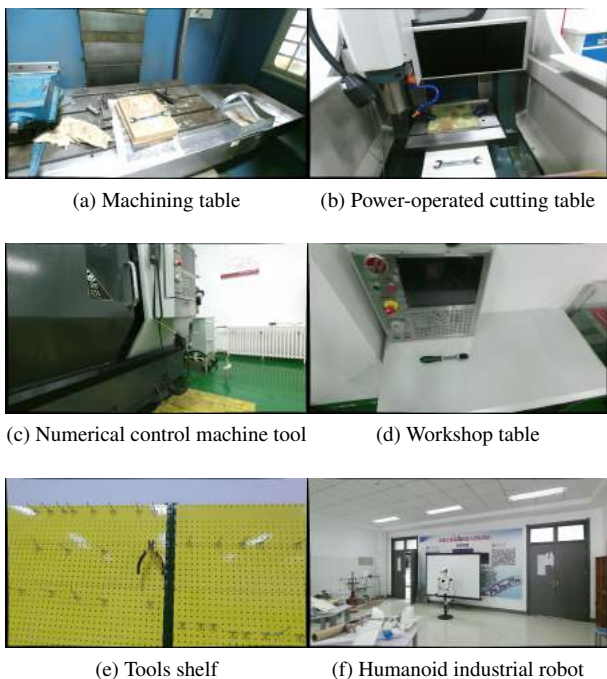


Fig. 3. Tools in different industrial scenes. ITD contains a wide variety of object categories in different industrial environments. We strive to collect images rich in classification, illumination and localization.

3.1. Object categories

As robots begin to manipulate standard tools and equipments available in industry scenarios, they will need to identify the tools and know the usages of them. This is achieved by gathering images of common tools in a typical factory through a computer vision and artificial intelligence study from September 2017 to May 2018. The dataset contains photos of 8 objects types that would be recognized by a qualified worker. The dataset has been collected in five distinct scenarios in factory, workshop, assembly line, and construction site scenarios characterized as shown in Fig.3. When people or industrial robots work in a factory, they are often in a moving state, which can results in view angle change, motion blur, illumination and clutter background. We specially designed dynamic scenes in factory environments to collect data.

3.2. Dataset format

Data was collected using a kinect 2.0 sensor [22] delivering 30 RGB-D frames per second at a resolution 1024×575 pixels + 512×424 depth frames. Since the items are relatively small, we collected data at the distance between 5 meters and 1 meter. Items are placed in their usual posture and environment and the camera point-of-view is that of the worker eyes. The worker was required to walk smoothly around the item while the camera was kept facing the target item consistently.

In order to compute the intrinsic arguments of the camera, we used a calibration checkerboard with known size. The dimension of the checkerboard is 9 squares \times 7 squares, whereas the length of each square is 3 cm. The calibration parameters and OpenCV tools used for calibration are also included in the dataset.

3.3. Ground truth

8 workers with Mechanical Engineer Certificate ranging in experience years from 1 to 10 were hired to label every tools they saw in inside and outside factories. For a given tool, a worker was asked to identify the tool's name, the category it belonged to and the possible usage. This task took a total of ~ 200 worker hours to complete. We assessed the category labeling tasks by comparing to dedicated supervisors. We analyzed precision and recall of five senior workers (managers and supervisors from factories) with the results obtained from the front-line workers. The true positives(TP), false positives(FP) and false negatives(FN) are defined as following [5]:

1. TP means the positive labeling that are categories as the positive class,
2. FP stands for the negative labeling that are categories as the positive class,
3. FN denotes the positive labeling that are categories as the negative class.

The precision and recall rate are computed by:

$$\begin{aligned}
 Precision &= \frac{TP}{TP + FP} \\
 Recall &= \frac{TP}{TP + FN}
 \end{aligned} \tag{1}$$

The results can be seen in Fig.4. It shows that the front-line workers have high recall rate than the senior workers. The labeling results are provided as ground truth in order to evaluate different vision-based target detection algorithms.

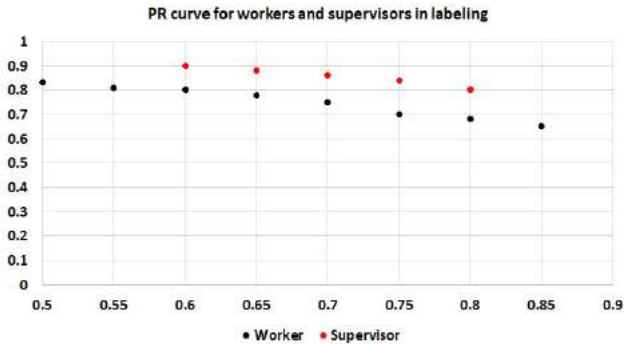


Fig. 4. Precision and recall rate of labeling. 8 workers ranging in experience years from 1 to 10 were hired to label tools in ITD dataset. We assessed the category labeling tasks by comparing to dedicated supervisors. We analyzed precision and recall of five senior workers (managers and supervisors from factories) with the results obtained from the front-line workers.

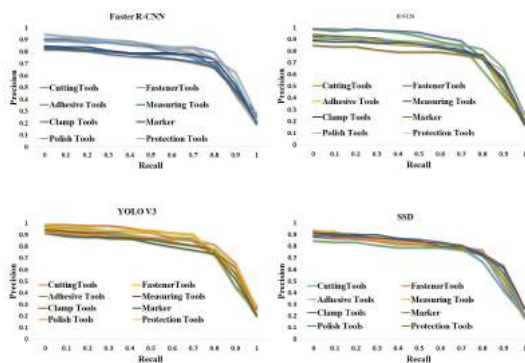


Fig. 5. Precision and recall curves of 4 detection methods. The experiments have been conducted on a PC with a 2.40GHz Intel(R) Xeon(R) CPU E5-2620 CPU, a GTX TITAN X GPU and 128GB memory. As we can see from the results exhibited, performances in clamp tools, marker and measuring tools are suboptimal.

4. Experiments and discussion

4.1. Object recognition evaluation

We evaluate state of the art object detection algorithms on ITD dataset. We carefully choose the Fast Region-based Convolutional Network(Faster R-CNN) [17, 30], Region Fully Convolutional Networks (R-FCN) [8], You Only Look Once (YOLO) V3 [29] and Single Shot MultiBox Detector(SSD) [26] as our benchmark methods for they have been widely used in object detection. We first briefly describe all these representations we have used for assessment.

4.1.1. Faster R-CNN

Faster R-CNN is a hybrid of deep convolutional network and region detector. The deep convolutional network combines a

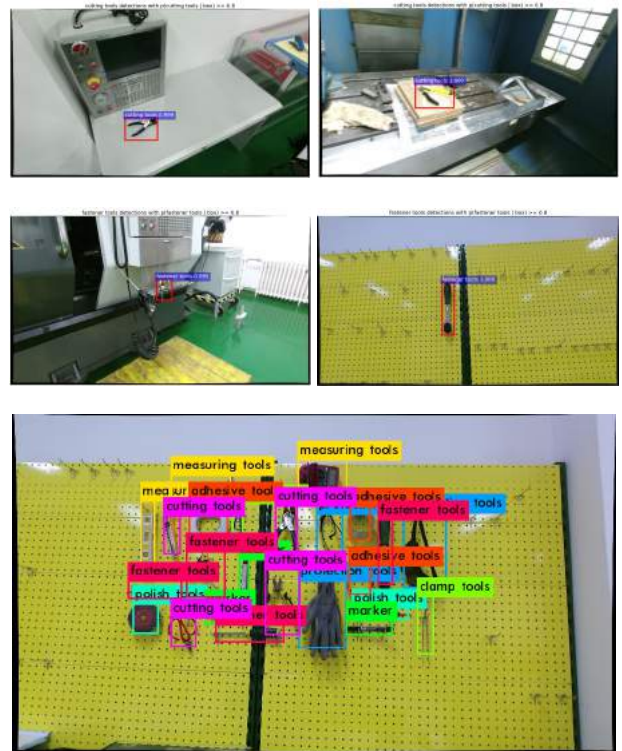


Fig. 6. Detection results in single and multi tools scenes. The conclusion presents that tools features can be easily affected by clutter background and dynamic environmental illumination. The image blur caused by the worker moving also make the performance fall short. This implies the defects of current detection methods and extensive efforts have to be dedicated according to the industrial requirements.

Table 2. The categories and usages of the tools in ITD Dataset.









Category	Sample Image	Name	Affordance
Cutting Tools		Scissor Utility Knife Puncher Nipper plier	Cut or separate small amounts of a material from the work piece by means of shear deformation This can be accomplished by single-point tools or multi-point tools
Fastener Tools		Open-end Wrench Torque wrench Hex wrench Screw driver	Provide grip and mechanical advantage in applying torque to turn objects or affixes multi objects together and the joints can be dismantled without damaging the joining components
Adhesive Tools		Pressure-sensitive tape Water activated tape Heat sensitive tape	Bind items together and resists their separation through non metallic substance applied to one surface
Measuring Tools		Multi-meter Vernier scale Air level	Measure a physical quantity This may require one-hand or two-hand operation
Clamp Tools		Plastic tweezers Flap tip clamp	Hold or pick-up items tightly together to be easily handled with the fingers
Marker		Permanent marker pen Waterproof marker	Draw or highlight notices on items They can be water-proof, dry-erase, or permanent.
Polish Tools		Machine file Sand paper	Smooth a workpiece's surface by rubbing it or using a chemical action
Protection Tools		Safety goggle Weld eye protector Glove	Enclose or protect body from injury or harmful contacts include physical, electrical, heat or chemicals

Table 3. Numerical results of baseline models evaluated with ground truth on Faster R-CNN, R-FCN, YOLO V3 and SSD methods over the ITD dataset.

	Faster	R-FCN	YOLO	SSD
Cutting Tools	70.12	72.65	85.56	69.51
Fastener Tools	58.61	63.64	71.81	53.43
Adhesive Tools	81.75	81.80	88.43	80.93
Measuring Tools	60.53	63.64	83.41	61.66
Clamp Tools	61.31	63.64	66.26	60.65
Marker	62.45	63.58	67.83	60.92
Polish Tools	50.76	54.52	73.18	48.32
Protection Tools	69.41	72.71	89.11	68.18
mAP	64.37	67.02	78.20	62.95

Region Proposal Network (RPN) and an object detection network [30]. The quality of detector is improved by using sparse object proposals. The whole image will be processed through conventional and max polling layers in order to produce a conventional feature map. A fixed length feature vector will be extracted by the region of interest pooling layer from the feature map. The features can be used for faster inference by classification and bounding-box regression.

4.1.2. R-FCN

The detection strategy of R-FCN consists of region proposal and region classification [8]. The candidate regions are extracted by the Region Proposal Network. R-FCN ends with a position-sensitive region of interest pooling layer. By cropping features from this last layer prior to prediction, R-FCN model could achieve similar accuracy to Faster R-CNN with less running time.

4.1.3. YOLO V3

YOLO applies an end-to-end single convolutional neural network that divides the image into regions, bounding boxes and region probabilities [29]. By examining the entire image during the training procedure, it get the contextual information and the knowledge of surroundings.

4.1.4. SSD

Single Shot MultiBox Detector (SSD) approach uses a single feed forward convolutional network that procedures bounding boxes collection and anchor offsets without requiring a pre-proposal classification [26].

4.2. Protocol

4.2.1. Protocol for holdout validation

We spitted the dataset by categories into training (50%), validation (25%) and testing (25%) sets randomly. We adopted the PASCAL Visual Object Challenge mean average precision

(mAP) evaluation metrics [12]. The mAP is calculated by 2:

$$mAP = \frac{\sum_{n=1}^C AvgPrecision(n)}{C} \quad (2)$$

$$AvgPrecision = \sum_{l=1 \dots N} P(l) \Delta Recall(l),$$

where C denotes the number of categories, $P(l)$ and $\Delta Recall(l)$ denote the precision value at every threshold and change in the recall respectively.

A detection is marked correct when the intersection size of the bounding boxes of the trial and the ground truth is more than half the size of their union. The numerical results (AP) of baseline models evaluated with ground truths are shown in Table.3. For its performance in skewed datasets [9], the precision and recall (PR) curve is also used as a valuable analytical tool for assessment.

4.2.2. Protocol for 4-fold cross-validation test

To further validate the ITD dataset, the 4-fold cross-validation test were carried out, which ensures that every image is tested once to prevent any bias error [1]. The dataset is divided by categories into 4 subsets (25% each) randomly. Every subset will works as the test dataset once, while the other three subsets are used as training and validation dataset. To be specific, when the subdataset is seclcted to train the model, 30% of images in subset will be used as validation dataset to fine-tune the model hyperparameters. And every model will be trained and tested four times to validate the proposed ITD dataset.

4.3. Results

The experiments have been conducted on a PC with a 2.40GHz Intel(R) Xeon(R) CPU E5-2620 CPU, a GTX TITAN X GPU and 128GB memory. Fig.5 shows the PR curves for Faster R-CNN, R-FCN, YOLO V3 and SSD methods over the ITD dataset and Fig.6 shows the single and multi tools detection results in different industrial scenes.

4.3.1. Comparison between different tools

As we can see from the results exhibited in Table.3, performances in clamp tools, marker and measuring tools are suboptimal, which attribute to their relatively small and may easily blocked by tools holder and grippers. Items like cutting tools, adhesive tools and protection tools, present good results partly due to their large size and difficult to be covered. YOLOv3 leads to the best accuracy, followed by R-FCN. The mAP results of SSD is lower than the others. The random crop approach used by the SSD data augmentation method may cause the consequence.

4.3.2. Comparison between different methods

The curves demonstrated in Fig.5 indicate that YOLO V3 is superior to other approaches. It is probably due to the improvement of predication strategy. YOLO extracts features at 3 different scales [29]. The change allows the method to get more meaningful information from small size objects. However, speed results show the different trend, the R-FCN algorithm is 52,989s, while YOLO v3 algorithm is 771,072s. These

Table 4. The performance of Faster R-CNN, R-FCN, YOLO V3 and SSD over 4-fold cross validation on the ITD dataset.

Test fold	Method	AvgPrecision(%)							
		Cutting	Fasterner	Adhesive	Measuring	Clamp	Marker	Polish	Protection
1 st fold	Faster	70.32	59.31	81.66	59.78	61.01	62.01	51.34	68.88
	R-FCN	72.51	62.69	81.12	68.08	67.97	62.76	54.41	72.68
	YOLO	86.31	70.31	87.56	82.76	66.11	66.68	72.89	89.47
	SSD	70.08	54.21	79.78	60.66	59.88	60.07	48.76	69.12
2 nd fold	Faster	70.21	59.98	81.87	59.31	61.32	62.41	51.51	68.92
	R-FCN	72.33	62.97	81.65	68.45	68.02	62.56	54.62	72.72
	YOLO	86.75	70.87	87.43	82.81	66.15	66.72	72.81	89.71
	SSD	69.79	54.66	80.02	60.76	59.93	60.26	48.78	69.04
3 rd fold	Faster	69.93	59.21	81.32	59.82	60.89	61.93	51.46	68.56
	R-FCN	73.08	62.12	80.97	67.91	67.78	62.83	54.21	72.55
	YOLO	86.82	70.42	87.88	82.21	66.09	66.81	72.63	89.32
	SSD	70.43	54.68	79.87	60.79	59.12	59.89	48.61	69.53
4 th fold	Faster	71.13	60.08	82.02	59.44	61.12	62.30	51.21	68.31
	R-FCN	72.23	62.45	81.43	68.22	67.45	62.78	54.61	72.18
	YOLO	86.12	69.89	87.33	82.88	65.93	66.53	72.46	90.03
	SSD	69.92	54.13	79.91	60.43	60.09	59.87	48.80	68.95
Average	Faster	70.40	59.65	81.72	59.59	61.09	62.16	51.38	68.67
	R-FCN	72.54	62.56	81.29	68.17	67.81	62.73	54.46	72.52
	YOLO	86.50	70.37	87.55	82.67	66.07	66.69	72.70	89.63
	SSD	70.06	54.42	79.90	60.66	59.76	60.02	48.74	69.16

approaches will degrade in industrial tools detection for relatively small training instances. It figures that for tools detection in industrial environments, those methods should ameliorate accordingly.

4.3.3. Comparison between different scenes

By analyzing the detection results of each scene (examples shown in Fig.6), the conclusion presents that tools features can be easily affected by clutter background and dynamic environmental illumination. The image blur caused by the worker moving also make the performance fall short. This implies the defects of current detection methods and extensive efforts have to be dedicated according to the industrial requirements.

4.3.4. Comparison through 4-fold cross-validation test

By adopting the 4-fold cross-validation method, the performance of each model over the ITD dataset is demonstrated in Table.4. In general, YOLOv3 still outperforms the other three detection methods. The different results between each categories are mainly caused by tools with different features. And the same categories get similar results among different test folds. It can conclude that there is also no huge bias error in ITD dataset.

5. Conclusion

We build a large-scale dataset for tools detection in industrial environments which is much more specialized and suitable

than any other general datasets in this field. We also establish a benchmark for items detection in industrial scenes. We believe ITD will promote the development of tools detection algorithms in industry. We currently only label tools in general but labeling grasping places may also provide significant manipulation information that may be useful for industrial utilization. In the future, we intend to further extend the dataset in terms of categories and sample quantities.

Acknowledgment

The work of C. Luo was supported in part by the National Natural Science Foundation of China under Grant 61701541 and in part by the Shandong Provincial Natural Science Foundation, China under Grant ZR2017QF003. The work of H. Zhou was supported in part by UK EPSRC under Grants EP/N508664/1, EP/R007187/1, EP/N011074/1, and Royal Society-Newton Advanced Fellowship under Grant NA160342. The work of P. Ren was supported in part by the National Natural Science Foundation of China under Grant 61671481 and in part by the Qingdao Applied Fundamental Research under Grant 16-5-1-11-jch. The authors would like to thank the Oil Industry Training Center in China University of Petroleum for the industrial tools annotation support.

References

- [1] Al-antari, M.A., Al-masni, M.A., Choi, M.T., Han, S.M., Kim, T.S., 2018. A fully integrated computer-aided diagnosis system for digital x-

- ray mammograms via deep learning detection, segmentation, and classification. *International journal of medical informatics* 117, 44–54.
- [2] Bai, X., Yan, C., Yang, H., Bai, L., Zhou, J., Hancock, E.R., 2018a. Adaptive hash retrieval with kernel based similarity. *Pattern Recognition* 75, 136–148.
- [3] Bai, X., Zhang, H., Zhou, J., 2014. Vhr object detection based on structural feature extraction and query expansion. *IEEE Transactions on Geoscience and Remote Sensing* 52, 6508–6520.
- [4] Bai, X., Zhou, J., Robles-Kelly, A., 2018b. Pattern recognition for high performance imaging.
- [5] Borji, A., Sihite, D.N., Itti, L., 2012. Salient object detection: A benchmark. in: *Computer Vision–ECCV 2012*. Springer, pp. 414–429.
- [6] Breazeal, C., Scassellati, B., 2002. Robots that imitate humans. *Trends in cognitive sciences* 6, 481–487.
- [7] Bullock, I.M., Feix, T., Dollar, A.M., 2015. The yale human grasping dataset: Grasp, object, and task data in household and machine shop environments. *The International Journal of Robotics Research* 34, 251–255.
- [8] Dai, J., Li, Y., He, K., Sun, J., 2016. R-fcn: Object detection via region-based fully convolutional networks, in: *Advances in neural information processing systems*, pp. 379–387.
- [9] Davis, J., Goadrich, M., 2006. The relationship between precision-recall and roc curves, in: *Proceedings of the 23rd international conference on Machine learning*, ACM. pp. 233–240.
- [10] Deng, J., Dong, W., Socher, R., Li, L.J., Li, K., Fei-Fei, L., 2009. ImageNet: A Large-Scale Hierarchical Image Database, in: *CVPR09*.
- [11] Everingham, M., Eslami, S.M.A., Van Gool, L., Williams, C.K.I., Winn, J., Zisserman, A., 2015. The pascal visual object classes challenge: A retrospective. *International Journal of Computer Vision* 111, 98–136.
- [12] Everingham, M., Van Gool, L., Williams, C.K.I., Winn, J., Zisserman, A., 2010. The pascal visual object classes (voc) challenge. *International Journal of Computer Vision* 88, 303–338.
- [13] Fang, K., Zhu, Y., Garg, A., Kurenkov, A., Mehta, V., Fei-Fei, L., Savarese, S., 2018. Learning task-oriented grasping for tool manipulation from simulated self-supervision. *arXiv preprint arXiv:1806.09266*.
- [14] Fermüller, C., Wang, F., Yang, Y., Zampogiannis, K., Zhang, Y., Barranco, F., Pfeiffer, M., 2018. Prediction of manipulation actions. *International Journal of Computer Vision* 126, 358–374.
- [15] Gao, F., Fei, M., Jun, W., Jinping, S., Erfu, Y., Huiyu, Z., 2018a. Visual saliency modeling for river detection in high-resolution sar imagery. *IEEE Access*, 1000–1014.
- [16] Gao, F., Huang, T., Jinping, S., Jun, W., Amir, H., Erfu, Y., 2018b. A new algorithm of sar image target recognition based on improved deep convolutional neural network. *Cognitive Computation*.
- [17] Girshick, R., 2015. Fast r-cnn, in: *Proceedings of the IEEE international conference on computer vision*, pp. 1440–1448.
- [18] Girshick, R., Donahue, J., Darrell, T., Malik, J., 2014. Rich feature hierarchies for accurate object detection and semantic segmentation, in: *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 580–587.
- [19] Heitz, G., Koller, D., 2008. Learning spatial context: Using stuff to find things, in: *European conference on computer vision*, Springer. pp. 30–43.
- [20] Jamone, L., Ugur, E., Cangelosi, A., Fadiga, L., Bernardino, A., Piater, J., Santos-Victor, J., 2016. Affordances in psychology, neuroscience and robotics: a survey. *IEEE Transactions on Cognitive and Developmental Systems*.
- [21] Koppula, H.S., Gupta, R., Saxena, A., 2013. Learning human activities and object affordances from rgb-d videos. *The International Journal of Robotics Research* 32, 951–970.
- [22] Lachat, E., Macher, H., Mittet, M., Landes, T., Grussenmeyer, P., 2015. First experiences with kinect v2 sensor for close range 3d modelling. *The International Archives of Photogrammetry, Remote Sensing and Spatial Information Sciences* 40, 93.
- [23] Lai, K., Bo, L., Ren, X., Fox, D., 2011. A large-scale hierarchical multi-view rgb-d object dataset, in: *Robotics and Automation (ICRA), 2011 IEEE International Conference on, IEEE*. pp. 1817–1824.
- [24] Lenz, I., Lee, H., Saxena, A., 2015. Deep learning for detecting robotic grasps. *The International Journal of Robotics Research* 34, 705–724.
- [25] Lin, T.Y., Maire, M., Belongie, S., Hays, J., Perona, P., Ramanan, D., Dollár, P., Zitnick, C.L., 2014. Microsoft coco: Common objects in context, in: *European conference on computer vision*, Springer. pp. 740–755.
- [26] Liu, W., Anguelov, D., Erhan, D., Szegedy, C., Reed, S., Fu, C.Y., Berg, A.C., 2016. SSD: Single shot multibox detector, in: *ECCV*.
- [27] Liu, Z., Yuan, L., Weng, L., Yang, Y., 2017. A high resolution optical satellite image dataset for ship recognition and some new baselines., in: *ICPRAM*, pp. 324–331.
- [28] Myers, A., Teo, C.L., Fermüller, C., Aloimonos, Y., 2015. Affordance detection of tool parts from geometric features., in: *ICRA*, pp. 1374–1381.
- [29] Redmon, J., Farhadi, A., 2018. Yolov3: An incremental improvement. *arXiv preprint arXiv:1804.02767*.
- [30] Ren, S., He, K., Girshick, R., Sun, J., 2015. Faster R-CNN: Towards real-time object detection with region proposal networks. *arXiv preprint arXiv:1506.01497*.
- [31] Russakovsky, O., Deng, J., Su, H., Krause, J., Satheesh, S., Ma, S., Huang, Z., Karpathy, A., Khosla, A., Bernstein, M., Berg, A.C., Fei-Fei, L., 2015. ImageNet Large Scale Visual Recognition Challenge. *International Journal of Computer Vision (IJCV)* 115, 211–252.
- [32] Stasse, O., Flayols, T., Budhiraja, R., Giraud-Esclasse, K., Carpentier, J., Mirabel, J., Del Prete, A., Souères, P., Mansard, N., Lamiroux, F., et al., 2017. Talos: A new humanoid research platform targeted for industrial applications, in: *Humanoid Robotics (Humanoids), 2017 IEEE-RAS 17th International Conference on, IEEE*. pp. 689–695.
- [33] Sünderhauf, N., Brock, O., Scheirer, W., Hadsell, R., Fox, D., Leitner, J., Upcroft, B., Abbeel, P., Burgard, W., Milford, M., et al., 2018. The limits and potentials of deep learning for robotics. *The International Journal of Robotics Research* 37, 405–420.
- [34] Tikhonoff, V., Pattacini, U., Natale, L., Metta, G., 2013. Exploring affordances and tool use on the icub, in: *Humanoid Robots (Humanoids), 2013 13th IEEE-RAS International Conference on, IEEE*. pp. 130–137.
- [35] Wang, C., Bai, X., Wang, S., Zhou, J., Ren, P., 2019. Multiscale visual attention networks for object detection in vhr remote sensing images. *IEEE Geoscience and Remote Sensing Letters* 16, 310–314.
- [36] Xia, G.S., Bai, X., Ding, J., Zhu, Z., Belongie, S., Luo, J., Datcu, M., Pelillo, M., Zhang, L., 2018. Dota: A large-scale dataset for object detection in aerial images, in: *Proc. CVPR*.
- [37] Xiao, B., Hancock, E.R., Wilson, R.C., 2009. Graph characteristics from the heat kernel trace. *Pattern Recognition* 42, 2589–2606.
- [38] Yi, S.J., McGill, S.G., Vadakedathu, L., He, Q., Ha, I., Han, J., Song, H., Rouleau, M., Zhang, B.T., Hong, D., et al., 2015. Team thor’s entry in the darpa robotics challenge trials 2013. *Journal of Field Robotics* 32, 315–335.
- [39] Yue, Z., Fei, G., Qingxu, X., Jun, W., Teng, H., Erfu, Y., Huiyu, Z., 2019. A novel semi-supervised convolutional neural network method for synthetic aperture radar image recognition. *Cognitive Computation*, 1–12.
- [40] Zhang, H., Bai, X., Zhou, J., Cheng, J., Zhao, H., 2013. Object detection via structural feature selection and shape model. *IEEE transactions on image processing* 22, 4984–4995.
- [41] Zhu, Y., Zhao, Y., Chun Zhu, S., 2015. Understanding tools: Task-oriented object modeling, learning and recognition, in: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 2855–2864.