

COCONET 2019

by Silvia Rostianingsih

Submission date: 13-Jan-2020 03:26PM (UTC+0700)

Submission ID: 1241393105

File name: Revised_COCONET_2019.docx (737.9K)

Word count: 2035

Character count: 10989

COCO (Creating Common Object in Context) Dataset for Laboratory Apparatus

Silvia Rostianingsih¹, Alexander Setiawan¹, Christopher Imantaka Halim
Informatics Departments, Petra Christian University, Indonesia

Abstract. — In order to create machine learning, we need to build a model. The model is created from a process called training. The goal of training is to develop an accurate model that answers some questions and in order to train a model, we need to collect a dataset. The quality and quantity of the data gathered will determine how good the predictive model can be. Helping the model to understand datasets like humans do is one of the important processes of machine learning. Datasets need to be constructed and transformed correctly. In this research, we compare the difference between creating a COCO dataset manually and creating a synthetic COCO dataset. Creating datasets for laboratory apparatus is not as difficult as creating a human object. The laboratory has a specific shape and form, thus the dataset had to have a limited number. As a result, we create both a dataset both manually and synthetically. The synthetic dataset helps to gain more datasets by combining some objects with different backgrounds.

Keywords— COCO dataset, synthetic, annotate, laboratory apparatus

1. Introduction

One important process of machine learning is helping the model to understand a dataset like humans do. Dataset are required to be constructed and transform correctly. However, some problems arise when dealing with data. First, one limitation to developing a good model is limited data [1]. Second, an imbalanced dataset presents difficulties for learning models to achieve high performances [2]. Third, some privacy issues occur when dealing with information about the individual [3] [4]. To balance the dataset, we need to collect data equally. Various techniques can be used to collect datasets such as manually, using API, doing web scraping, using a dataset platform or generating synthetic dataset.

To generate a synthetic dataset, many algorithms are used depending on the data type. For an online social network dataset, graph topology can be used [3]. If it comes with images, some methods include UCSD Anomaly Detection [1] [5], procedural randomization and manual modelling [6] or image masking and vertices [7].

In this research, we are going to create COCO dataset of laboratory apparatus, which COCO has not provided it yet. The laboratory apparatus has a certain forms, thus most of them have a transparent form. In order to enrich the dataset, we also create synthetics dataset. We are going to make a comparison between creating the COCO dataset manually and creating the synthetic COCO dataset.

2. Literature

2.1. Common Objects in Context (COCO) dataset

COCO is a large-scale object detection, segmentation and captioning dataset [8]. They use several sources to collect object categories, which are grouped into three types, namely iconic-object images, iconic-scene images and non-iconic images. Iconic-object images have a single large object that appeals

to the image. Iconic-scene images are recognized as a scene without any object being obvious. A non-iconic image is an image with some objects. It showed that datasets with more non-iconic images are better to use for generalizing. Three stages are used to annotate the image collections. First, category labelling is used to determine which object categories are present in each image. Second, instance spotting is used to label each instance of a specific category found in the previous stage. The final stage is segmenting each object instance.

COCO is a common dataset format used by Microsoft, Google, and Facebook. It is used as a benchmark to measure machine learning algorithm performance. The other common datasets are PASCAL VOC and ImageNet.

2.2. Synthetic dataset

Synthetic data generation is an alternative method for creating a large number of datasets. There are two techniques to create synthetic dataset [9]. First, we create fully synthetic data by making data synthetically. Second, we can combine the synthetic dataset with real data, by replacing some sensitive attribute values. Beside the privacy issue, synthetic dataset is also effective to evaluate algorithms and their usefulness. Mayer found points that need to be considered in order to create a synthetic dataset, namely (1) diversity is important, (2) realism is overrated, (3) learning schedules matters and (4) camera knowledge helps [6].

2.3. Laboratory apparatus

For our collections, we are creating datasets for laboratory apparatus such as a beaker, Erlenmeyer flask, Florence flask, graduated cylinder and petri dish. Each contains 200-300 datasets. We gathered the data from taking pictures at chemistry laboratories located at Sepuluh Nopember Institute of Technology, Surabaya. We are focusing on chemistry laboratories for first to second-year students, because they use laboratory apparatus a lot. We were visiting four laboratories from three faculties, which are Department of Chemistry, Department of Chemical Engineering and Department of Industrial Chemical Engineering. The example of images can be seen in Figure 1.

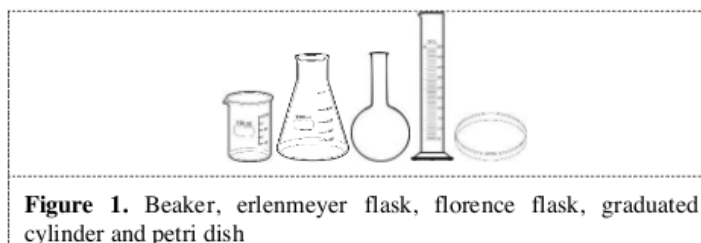


Figure 1. Beaker, erlenmeyer flask, florence flask, graduated cylinder and petri dish

The beaker, Erlenmeyer flask, Florence flask, graduated cylinder and petri dish are some examples of the laboratory apparatus. The beaker can be made of plastic or glass. We only gathered data from a beaker glass, which can be heated. The common sizes are 50 mL, 100 mL, 250 mL and 400 mL. The erlenmeyer flask is made of glass, which can be heated and used in titration. The common sizes are 100 mL or 250 mL. The florence flask is made of glass, which can be heated and used in making or storing solution. The common sizes are 125 mL, 250 mL and 500 mL. A graduated cylinder can be made of glass or plastic, which is used to measure approximate values; however, it must not be heated. The common sizes are 10 mL, 50 mL and 100 mL used to measure approximate volumes, must not be heated.

3. Methodology

3.1. Creating COCO dataset manually

In order to create the COCO dataset manually (Figure 2), first we specify the laboratory apparatus labels. As mentioned previously, labels would be Erlenmeyer flask, Florence flask, graduated cylinder and petri dish.

Second, we collect the images from laboratories at Sepuluh Nopember Institute of Technology, Surabaya. For each category, there are approximately 200-300 images. The images are taken with the original background. Some of the objects had an overlap position. We also scraped the Google Image Search [10], using a python helper software. Herewith the script:
 googleimagesdownload --keywords " beaker, erlenmeyer flask, florence flask, graduated cylinder, petri dish"

Last, we annotate the images as COCO format. We are using an open-source annotation software, which can automatically produce a COCO formatted data [11]. This step is the most time-consuming.

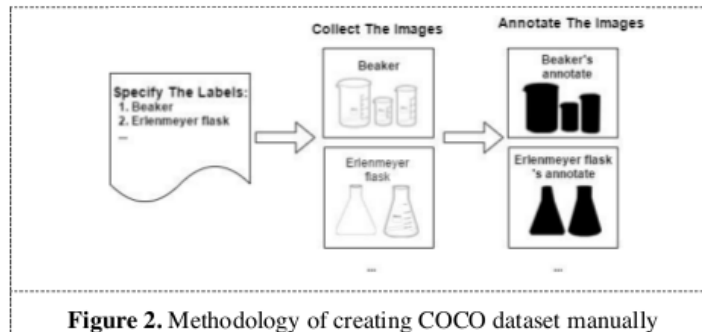


Figure 2. Methodology of creating COCO dataset manually

3.2. Creating synthetic COCO dataset

In order to create a synthetic COCO dataset (Figure 3), first we specify the laboratory apparatus labels. As mentioned earlier, labels would be Erlenmeyer flask, Florence flask, graduated cylinder and petri dish. Second, we choose the foreground and background image, and then we combine with the laboratory apparatus which had already been annotated.

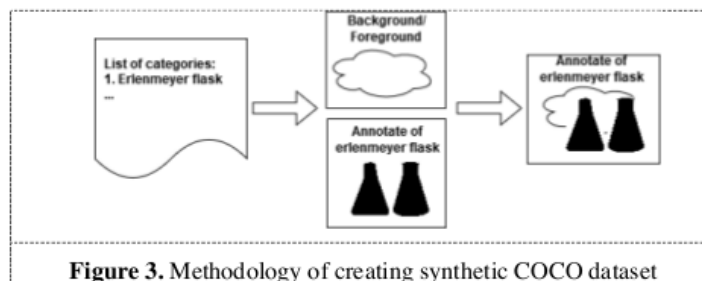


Figure 3. Methodology of creating synthetic COCO dataset

4. Implementation

4.1. Creating COCO dataset manually

COCO is a standard dataset format for annotating the image collection, which is used to for data preparation in machine learning. Annotate means to create metadata for an image. From Figure 2, we know there is an image named beaker and we know the position (x and y coordinate, width, length, and polygon area).

The labels from Figure 2 will be the category of the COCO dataset (Figure 4). The number described in each category is the number of objects already annotated.

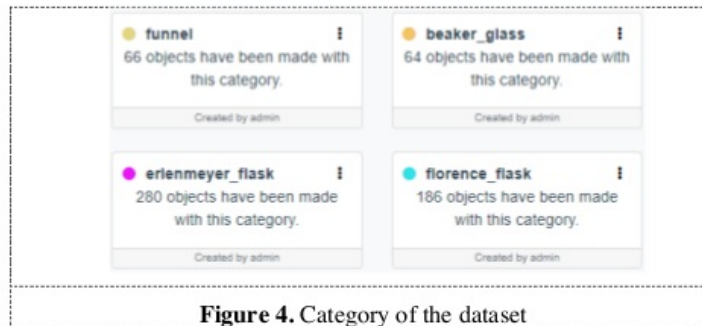
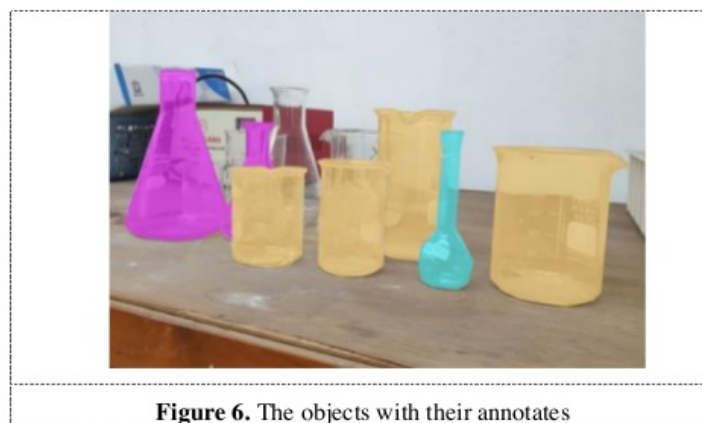


Figure 5 contains many laboratory apparatus objects, where more than 10 objects are displayed. For this example, we are only annotating seven objects with three categories (Figure 6). The categories are the beaker glass, Erlenmeyer flask and Florence flask (Figure 7). For the beaker glass, there are four objects that represent with yellow color. For the Erlenmeyer flask, there are two objects that represent with the magenta color, For the Florence flask, there is one object that represents with cyan color.



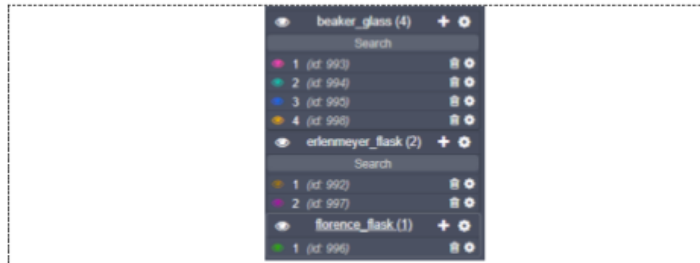


Figure 7. The categories and their objects

After creating the COCO dataset by using the tools from Justin Brooks, we download the json. The json structure contains images, categories and annotations. The images contain an id (auto increment), dataset_id (the folder id), file name, width, height, and file location (Figure 8).

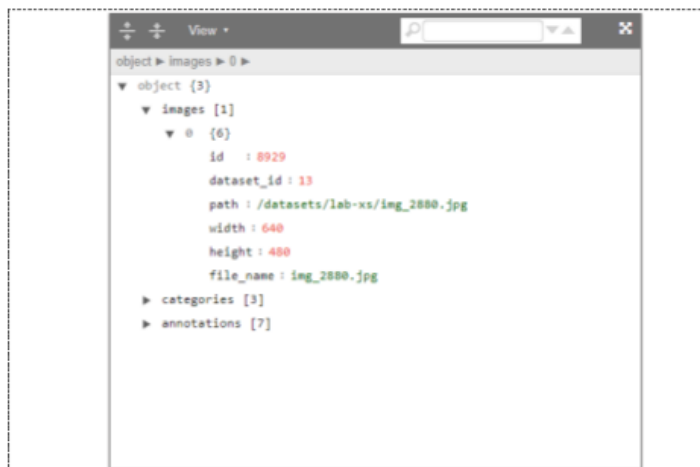


Figure 8. The structure of json (1)

Categories are the information of the labels and the color (Figure 9). Annotations contain segmentation (Figure 10), bounding box and color. Segmentation contains the position of each vertex from the annotation. The bounding box is the maximum x and y of the image. The result from this manually annotate the objects can be seen in Figure 11. We had a set of datasets that annotate manually. Average time to manually annotate per objects is three and a half minutes.



Figure 9. The structure of json (2)



Figure 9. The structure of json (3)



Figure 10. The result of manual annotate

4.2. Creating Synthetic COCO dataset

Creating a synthetic dataset is challenging; we need to combine many backgrounds and objects. The background for synthetic dataset would be the laboratory itself. We choose ten backgrounds for five objects and we generate 500 annotations for each objects. The json format will be same as Figure 7 – Figure 11. We also used the tools from Brooks [11]. The result from the synthetic dataset can be seen in Figure 12. With the same background, we put the funnel, beaker, Erlenmeyer flask and Florence flask in a different place. The average time to create a synthetic dataset per object is 0.05 minute .

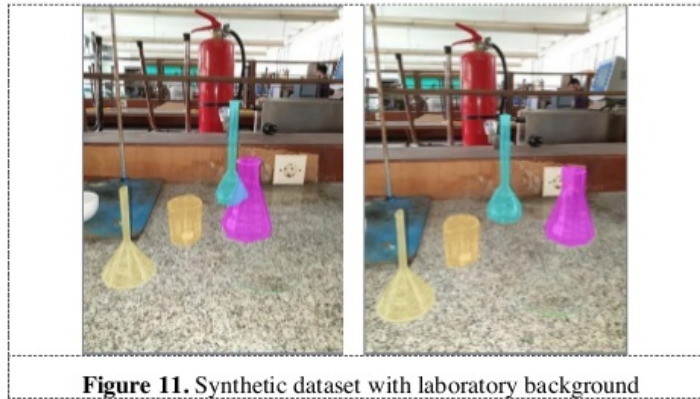


Figure 11. Synthetic dataset with laboratory background

5. Conclusion

Laboratory apparatus has a specific shape and form; therefore, the dataset had a limited number. In order to enrich the dataset, we need to collect images with many backgrounds or foregrounds. A synthetic dataset helps to collect more datasets by combining background and object. We placed some overlap objects in one background to get different viewpoints. For further research, the COCO dataset will be used for recognize laboratory apparatus using machine learning.

- [1] H. K. Ekbatani, O. Pujol and S. Segui, "Synthetic Data Generation for Deep Learning in Counting Pedestrians," in *Proceedings of the 6th International Conference on Pattern Recognition Applications and Methods (ICPRAM)*, Porto, 2017.
- [2] D. C. Li, S. C. Hu, L. S. Lin and C. W. Yeh, "Detecting representative data and generating synthetic samples to improve learning accuracy with imbalanced data sets," *PLoS One*, vol. 12, no. 8, 2017.
- [3] D. F. Nettleton, "A Synthetic Data Generator for Online Social Network Graphs," *Social Network Analysis and Mining*, vol. 6, no. 1, 2016.
- [4] S. H and M. H. S, "A Review Of Synthetic Data Generation Methods For Privacy Preserving Data Publishing," *International Journal of Scientific & Technology Research*, vol. 6, no. 3, pp. 95-101, 2017.
- [5] W. X. Li, N. Vasconcelos and V. Mahadevan, "Anomaly Detection and Localization in Crowded Scenes," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 36, no. 1, pp. 18-32, 2013.
- [6] N. Mayer, E. Ilg, P. Fischer, C. Hazirbas, D. Cremers, A. Dosovitskiy and T. Brox, "What Makes Good Synthetic Training Data for Learning Disparity and Optical Flow Estimation?," *International Journal of Computer Vision*, vol. 126, no. 9, pp. 942-960, 2018.
- [7] T. Remez, J. Huang and M. Brown, "Learning to Segment via Cut-and-Paste," in *European Conference on Computer Vision*, Munich, 2018.

- [8] T. Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár and C. L. Zitnick, "Microsoft COCO: Common Objects in Context," in *European Conference on Computer Vision*, Zurich, 2014.
- [9] A. Dandekar, R. A. M. Zen and S. ´. Bressan, "A comparative study of synthetic dataset generation techniques," National University of Singapore, Singapore, 2018.
- [10] H. Vasa, "Google Images Download," 2019. [Online]. Available: <https://github.com/hardikvasa/google-images-download>.
- [11] J. Brooks, "COCO Annotator," 2019. [Online]. Available: <https://github.com/jsbroks/coco-annotator/>.

Acknowledgments

Special thanks to our Ministry of Research, Technology and Higher Education for supporting this research financially. We also thank Sepuluh Nopember Institute of Technology, Surabaya for allowing us to gather dataset information from a chemistry laboratory.

Coconet

ORIGINALITY REPORT

11%

SIMILARITY INDEX

8%

INTERNET SOURCES

8%

PUBLICATIONS

3%

STUDENT PAPERS

PRIMARY SOURCES

1	Lecture Notes in Computer Science, 2014. Publication	3%
2	www.flashcardmachine.com Internet Source	1%
3	www.ijraset.com Internet Source	1%
4	Mejdi Dallel, Vincent Havard, Yohan Dupuis, David Baudry. "Digital twin of an industrial workstation: A novel method of an auto-labeled data generator using virtual reality for human action recognition in the context of human-robot collaboration", Engineering Applications of Artificial Intelligence, 2023 Publication	1%
5	neptune.ai Internet Source	1%
6	www.answers.com Internet Source	1%
7	core.ac.uk Internet Source	1%

8

journals.plos.org

Internet Source

1 %

9

Dan-Gabriel Anton. "Evaluating the effort of building a Machine Learning model for malware detection from ground zero", 2022 24th International Symposium on Symbolic and Numeric Algorithms for Scientific Computing (SYNASC), 2022

Publication

1 %

10

Hasan Ozgur Kapici, Hakan Akcay, Ece Ebrar Koca. "Comparison of the Quality of Written Scientific Arguments in Different Laboratory Environments", International Journal of Science and Mathematics Education, 2021

Publication

1 %

Exclude quotes On

Exclude matches < 1%

Exclude bibliography On