Data Sheet for Maize Image Data Collected in Tanzania Under the Lacuna Project

We present the Lacuna Maize data sheet created by a group of researchers from The Nelson Mandela African Institution of Science and Technology and the Tanzania Agricultural Research Institute (TARI) in Tanzania. We follow the datasheet for dataset framework created by (Gebru et al. 2021).

	Motivation
For what purpose was the	The dataset was created to provide an open, well-labelled,
data set created?	sufficiently curated and accessible maize image dataset. Data
	scientists, researchers, and the broader machine learning
	community can use the dataset for various machine learn-
	ing experiments to build maize crop disease diagnosis and
	spatial analysis solutions.
Was there a specific task in	Although the agricultural sector is a national economic de-
mind?	velopment priority in sub-Saharan Africa, crop pests and dis-
	eases have been the challenge affecting major food security
	crops like Maize. The images target the diagnostics of Maize
	Lethal Necrosis (MLN) and Maize Streak Virus (MSV) dis-
	eases. We are motivated in developing end to end tools to
	help farmers diagnose diseases and also help to improve the
	production of maize in the country. The current state of data
	collection and crop pest and disease diagnosis is transitioning
	from disease identification using visible symptoms to using
	data-driven solutions applying machine learning and com-
	puter vision techniques. The image data previously collected
	has not been sufficiently curated, prepared, and shared with
	the broader community.
Who created the dataset?	The dataset was created by Scientists and Masters students
	from the Nelson Mandela African Institution of Science and
	Technology and the Tanzania Agricultural Research Insti-
	tute (TARI) in Tanzania.
Who funded the creation of	This work was carried out with support from Lacuna Fund,
the dataset?	an initiative co-founded by The Rockefeller Foundation,
	Google.org, and Canada's International Development Re-
	search Centre. The views expressed herein do not necessarily
	represent those of Lacuna Fund, its Steering Committee, its
	funders, or Meridian Institute.: 0328-S-001.
	Composition

What do the instances that comprise the dataset represent?	Each instance in the dataset consists of a crop image with an image status, i.e., Healthy, Maize Lethal Necrosis, and Maize Streak Virus, crop variety, crop age, and location (district, sub-county).
How many instances are there in total (of each type, if appropriate)?	The dataset consists of 17,277 labeled images where Healthy are 5542, Maize Lethal Necrosis are 5068, and Maize Streak Virus are 6667.)
Does the dataset contain all possible instances or is it a sample (not necessarily random) of instances from a	The dataset has image samples collected from significant maize growing districts selected with the guidance of agricultural experts to obtain a representative dataset.
larger set? What data does each instance consist of? "Raw" data or features?	The data consists of raw image data.
Is there a label or target associated with each instance? If so, please provide a description.	Each instance is associated with a class label based on the status of the crop: healthy or diseased. The given labels per image are: Healthy, Maize Streak Virus and Maize Lethal Necrosis as shown in Figure 1.
Is any information missing from individual instances?	None
Are relationships between individual instances made explicit?	There are no relationships between the different image instances in the dataset.
Are there recommended data splits (for example, training, devel-opment/validation, testing)?	We do not specify any data splits.
Are there any errors, sources of noise, or redundancies in the dataset? If so, please provide a description.	None
Is the dataset self-contained, or does it link to or otherwise rely on external resources?	The dataset is self-contained.
Does the dataset contain data that might be considered confidential?	No.
Does the dataset contain data that, if viewed directly, might be of- fensive, insulting, threatening, or might otherwise cause anxiety?	No.
ounci wise cause anxiety:	Collection Process

with each instance acquired?	ing cameras on mobile phones from farmer's gardens. The
	AdSurv mobile application installed on Samsung phones was used to take photos of maize leaves.
7771	-
What mechanisms or proce-	The data was collected using the Adsurv application, which
dures were used to collect the	is a mobile application that enables crowdsourcing of crop
data?	disease data from farmers' gardens. Adsurv application was
	installed on mobile phones/tablets used during the data col-
	lection process.
If the dataset is a sample from	The dataset is not from a larger set.
a larger set, what was the	
sampling strategy?	
Who was involved in the data	The data was collected by a team of researchers and stu-
collection process?	dents from the The Nelson Mandela African Institution of
	Science and Technology and Tanzania Agricultural Research
	Institute (TARI) .
Over what timeframe was the	The data was collected in the range of six months, period-
data collected?	ically between February 2021 and July 2021 from farms in
	Arusha, Kilimanjaro and Manyara regions in Tanzania.
Were any ethical review pro-	No.
cesses conducted (for exam-	
ple, by an institutional review	
board)?	
Preprocessing, cleaning, and labelling	

Was any preprocess-	The following steps were taken to process the data:	
ing/cleaning/labeling of the data done (for example, discretization or bucketing, tokenization, part-of-speech tagging, SIFT feature extrac- tion, removal of instances, processing of missing val- ues)?	1. Gathering raw images: First the maize leaves images were collected using the AdSurv application installed on smartphones.	
	2. Eliminating duplicate images: The identified duplicate images were removed but a very small number (that were not initially identified) might still exist in the dataset.	
	3. The number of remaining duplicates should be small enough so as not to significantly impact training/testing. The dataset contains distinct images that are not defined to be duplicates but are extremely similar.	
	4. Labeling: The images were labeled to indicate the belonging class (healthy, Maize Lethal Necrosis, Maize Streak Virus).	
	5. Curation: The images were annotated for various computer vision tasks such as image classification, image object detection and image segmentation.	
	6. Renaming: The images in each class were renamed to comprise image number	
Was the "raw" data saved in addition to the preprocessed/cleaned/ labeled data (for example, to support unanticipated future uses)? If so, please provide a link or other access point to the "raw" data.	The raw unprocessed data is stored locally on data storage servers in the Makerere Artificial Intelligence Lab.	
Is the software that was used to preprocess/clean/label the data available? If so, please provide a link or other access point.	The link to the annotation tool is available: https://github.com/AI-Lab-Makerere/web-annotation-tool	
Uses		
Has the dataset been used for any tasks already? If so, please provide a description.	Yes, we have used the dataset to build baseline disease classification models.	
Is there a repository that links to any or all papers or systems that use the dataset?	No.	

What (other) tasks could the dataset be used for?	The dataset can be used for building object detection, segmentation, and time-series analysis models.
Is there anything about	No.
v S	NO.
the composition of the dataset or the way it was	
collected and preprocessed/	
cleaned/labeled that might	
impact future uses?	
impact future uses:	Distribution
Will the dataset be dis-	Yes, the dataset will be made publicly available.
tributed to third parties out-	,
side of the entity (for exam-	
ple, company, institution, or-	
ganization) on behalf of which	
the dataset was created? If so,	
please provide a description.	
How will the dataset be dis-	The dataset and the associated metadata are stored on the
tributed (for example, tarball	Harvard DataVerse which is an open-source data repositor
on website, API, GitHub)?	The dataset is assigned a Digital Object Identifier: https:/
Does the dataset have a dig-	doi.org/10.7910/DVN/GDONSQ.
ital object identifier (DOI)?	doi.org/10.7910/DVN/GDONOQ.
When will the dataset be dis-	The deteget is excilable under the specified DOI
tributed?	The dataset is available under the specified DOI.
Will the dataset be dis-	The dataset is licensed under the CC BY license that allow
tributed under a copyright	users to share and adapt the dataset so long as they give
or other intellectual property	credit to data set creators.
(IP) license, and/or under ap-	
plicable terms of use (ToU)?	
Have any third parties im-	No.
posed IP-based or other re-	
strictions on the data associ-	
ated with the instances?	
Do any export controls or	No.
other regulatory restrictions	
apply to the dataset or to in-	
dividual instances?	
XX71 11 1	Maintenance
Who will be support-	The dataset will be maintained by the research team at the
ing/hosting/maintaining	Makerere Artificial Intelligence Lab and The Nelson Mande
the dataset?	African Institution of Science and Technology . The tea
	will support, host, and maintain the dataset.
How can the owner/curator/	The dataset manager can be contacted via email.
manager of the dataset be con-	
tacted (for example, email ad-	
dress)?	
Is there an erratum?	No.

Will the dataset be updated	All updates to the dataset will be documented and commu-
(for example, to correct label-	nicated through the Makerere AI Lab GitHub repository.
ing errors, add new instances,	
delete instances)?	
Will older versions of the	Yes, the older versions will be stored locally on data storage
data- set continue to be sup-	servers in the Makerere Artificial Intelligence Lab and on
ported/hosted/ maintained?	remote data storage buckets on the Google cloud.
If so, please describe how.	
If others want to ex-	Interested researchers can send an email to data managers
${ m tend/augment/build}$	manager one and manager two to discuss the dataset exten-
on/contribute to the dataset,	sion and contribution.
is there a mechanism for them	
to do so?	

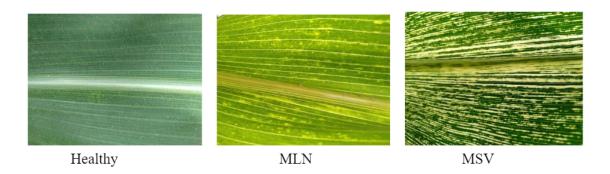


Figure 1: Maize Data Labels.

References

Gebru, T., Morgenstern, J., Vecchione, B., Vaughan, J. W., Wallach, H., Iii, H. D., Crawford, K. (2021). Datasheets for datasets. Communications of the ACM, 64(12), 86-92.