



Smithsonian

Using machine learning tools to quantify mercuric chloride staining across digitized herbarium specimens

Isabella Schrader^{1,2}, Mike Trizna², Ashlyn Powell³, Paul Frandsen^{2,3}, Alex White^{1,2}, Rebecca Dikow², Eric Schuettpelz¹

¹ Smithsonian Institution, National Museum of Natural History, Botany Department

² Smithsonian Institution, Office of the Chief Information Officer, Data Science Lab

³ Brigham Young University, Department of Plant and Wildlife Sciences



REU Site, OCE-1560088

Introduction

Mercuric chloride was commonly used in the past to prevent insect damage to botanical specimens in herbaria. Because this substance is toxic to humans, knowing the number and location of contaminated specimens in a collection is important. Fortunately, the staining is visible upon specimen inspection, because mercuric chloride crystallizes over time. Because the U.S. National Herbarium contains more than 5 million specimens, manual inspection of every specimen is not tractable. In 2017, NMNH and OCIO scientists built a machine learning model to identify mercuric chloride staining on digitized herbarium sheets, but only applied the model to a portion of the digitized herbarium (Schuettpelz et al, 2017). This project seeks to update the model and apply it to all digitized specimens.

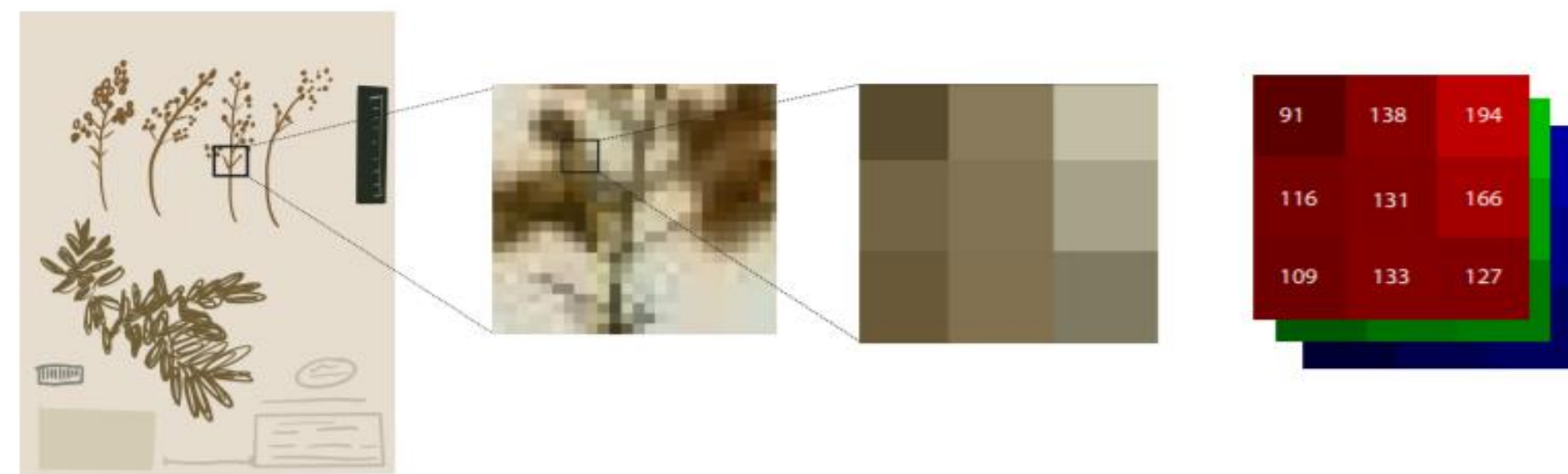


Figure 1: Process of herbarium specimen converted into inputs for machine learning (Borowiec et al., 2021)

Methods

Throughout this project, Jupyter notebooks were used to work in Python. First, Pixplot (<https://dhlabs.yale.edu/projects/pixplot/>) was used to visualize digitized specimens using unsupervised machine learning. Figure 2 shows a PixPlot of 1000 sample images.

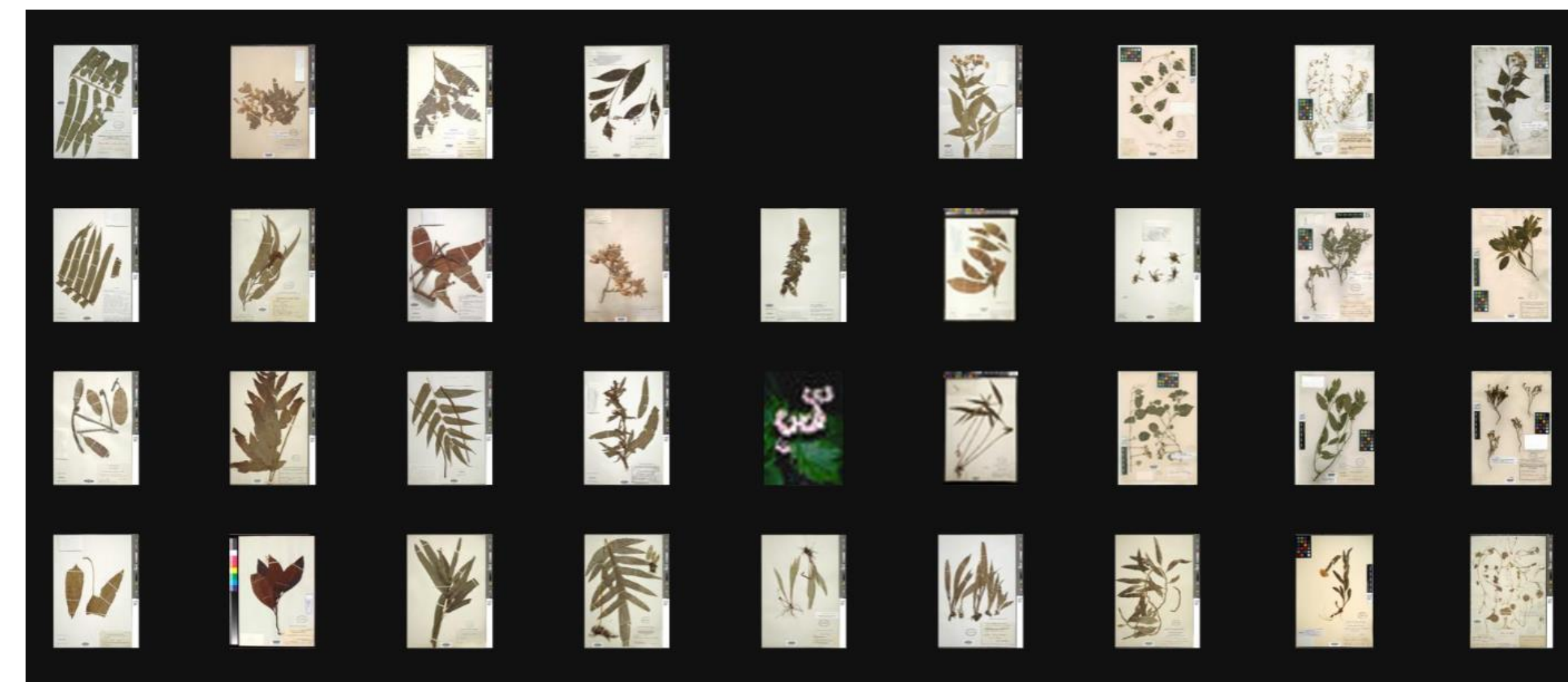


Figure 2: 1000 image PixPlot

We used the Smithsonian Institution High Performance Computing Cluster, also known as Hydra, to access files and store data. During the summer, Data Carpentry workshops were attended, which aided in developing skills used to group the data into smaller samples. We also used Google Colaboratory (<https://colab.research.google.com>) to train a new model to identify mercuric chloride staining, taking advantage of recent advances in machine learning technology.



Figure 3: Stained specimen



Figure 4: Unstained specimen

Acknowledgments / References

- Elizabeth Cottrell, Gene Hunt, and Virginia Power, for their roles in the NHRE program
- Borowiec, M. L., Frandsen, P., Dikow, R., McKeeken, A., Valentini, G., & White, A. E. (2021). Deep learning as a tool for ecology and evolution. [10.32942/osf.io/nt3as](https://doi.org/10.32942/osf.io/nt3as)
- YaleDHLab. (n.d.). YaleDHLab/pix-plot: A WebGL viewer for UMAP or TSNE-clustered images. GitHub.
- Schuettpelz, E., Frandsen, P. B., Dikow, R. B., Brown, A., Orli, S., Peters, M., Metallo, A., Funk, V. A., Dorr, L. J. 2017. Applications of deep convolutional neural networks to digitized natural history collections. Biodiversity Data Journal. doi: 10.3897/BDJ.5.e21139
- Funding from NSF OCE-1560088



edan_id	title	timestam	lastupdat	Barcode	specimen	media_c	media_g	media_g	media_g	media_a	aws_me	USNM	Other	timestam	lastupdat
		p_unix	e_unix		_guid	ount	uid	uid_list	ws_id	lis	dia_coun	Number	Numbers	p_dt	e_dt
edanmd	Peperomi				http://n2			http://n2	http://n2						
m	hermandii	1619514	1619514	1.0102e+	/65665/3			/65665/	/65665/	NMNH-				2021-04-	2021-04-
0	nmnhbot	folia	194	189	06	886114c	1.0	m334175	m334175	0101020	4	1.0	1884734	NaN	27
any_104	(Vahl) A.				e-5e7f-			257-	257-					09:03:14	09:03:09
89462	Dietr.				43c9...			8640-	8640-						
								4e70...	4e70...						

Figure 5: Sample DataFrames

Results and Discussion

More than two million digitized herbarium specimens are available on the Smithsonian Open Access platform. Python code to download and resize images is available on GitHub (https://github.com/sidatasciencelab/mercury_sheets). This code also has the potential to be modified to access Smithsonian Open Access images from other museums so that researchers can more easily assemble datasets for

other applications. Sixteen metadata fields are recorded for each specimen (Figure 5; two specimens shown). The machine learning model we are developing could also be applied to specimens from herbaria around the world to identify mercuric chloride staining.



Figure 6: Open Access