

Methods

StomataCounter: a neural network for automatic stomata identification and counting

Karl C. Fetter^{1,2*} , Sven Eberhardt^{3*} , Rich S. Barclay² , Scott Wing²  and Stephen R. Keller¹ 

¹Department of Plant Biology, University of Vermont, Burlington, VT 05405, USA; ²Department of Paleobiology, Smithsonian Institution, National Museum of Natural History, Washington, DC 20560, USA; ³Amazon.com, Inc., Seattle, WA 98121, USA

Author for correspondence:

Karl C. Fetter

Tel: +1 802 656 2930

Email: kf@uvm.edu

Received: 5 March 2019

Accepted: 21 April 2019

New Phytologist (2019)

doi: 10.1111/nph.15892

Key words: computer vision, convolutional deep learning, neural network, phenotyping, stomata.

Summary

- Stomata regulate important physiological processes in plants and are often phenotyped by researchers in diverse fields of plant biology. Currently, there are no user-friendly, fully automated methods to perform the task of identifying and counting stomata, and stomata density is generally estimated by manually counting stomata.
- We introduce StomataCounter, an automated stomata counting system using a deep convolutional neural network to identify stomata in a variety of different microscopic images. We use a human-in-the-loop approach to train and refine a neural network on a taxonomically diverse collection of microscopic images.
- Our network achieves 98.1% identification accuracy on *Ginkgo* scanning electron microscopy micrographs, and 94.2% transfer accuracy when tested on untrained species.
- To facilitate adoption of the method, we provide the method in a publicly available website at <http://www.stomata.science/>.

Introduction

Stomata are important microscopic organs that mediate important biological and ecological processes and function to exchange gas with the environment. A stoma is composed of a pair of guard cells forming an aperture that, in some species, are flanked by a pair of subsidiary cells (Bergmann & Sack, 2007). Regulation of the aperture pore size, and hence the opening and closing of the pore, is achieved through changing turgor pressure in the guard cells (Shimazaki *et al.*, 2007). Stomata have evolved to permit exchange between internal and external sources of gases, most notably CO₂ and water vapour, through the impervious layer of the cuticle (Kim *et al.*, 2010).

Because of their importance in regulating plant productivity and response to the environment, stomata have been one of the key functional traits of interest to researchers working across scales in plant biology. At the molecular level, regulation of stomata has been the subject of numerous genetic studies (see Shimazaki *et al.*, 2007 and Kim *et al.*, 2010 for detailed reviews) and crop improvement programme have modified stomata phenotypes to increase yield (Fischer *et al.*, 1998). Stomata also mediate trade-offs between carbon gain and pathogen exposure that are of interest to plant ecophysiologicalists and pathologists. For example,

foliar pathogens frequently exploit the aperture pore as a site of entry. In *Populus*, some species, and even populations within species, have evolved growth strategies that maximise carbon fixation through increased stomatal density and aperture pore size on the adaxial leaf surface. This adaptation results in a cost of increased infection by fungal pathogens that have more sites of entry to the leaf (McKown *et al.*, 2014). Stomata, as sites of water vapour exchange, are also implicated in driving environmental change across biomes (Hetherington & Woodward, 2003) and variation of stomatal density and aperture pore length are linked to changes in ecosystem productivity (Wang *et al.*, 2015). Stomatal traits are of particular interest to paleoecologists and paleoclimatologists due to the relationship between stomatal traits and gas exchange. Measurements of stomatal density from fossil plants have been proposed as an indicator of paleoatmospheric CO₂ concentration (Royer, 2001), and measuring stomatal traits to predict paleoclimates has become widely adopted (McElwain & Steinthorsdottir, 2017).

Researchers across a wide variety of disciplines in plant biology will phenotype stomatal traits for decades to come due to their physiological importance. A typical stomatal phenotyping pipeline consists of collecting plant tissue, creating a mounted tissue for imaging, imaging the specimen, and manual phenotyping of a trait of interest. This last step can be the most laborious, costly, and time-consuming task, reducing the efficiency of the

*These authors contributed equally to this work.

data acquisition and analysis pipeline. This is especially important in large-scale plant breeding and genome-wide association studies, in which phenotyping has been recognised as the new data collection bottleneck, in comparison with the relative ease of generating large genome sequence datasets (Hudson, 2008).

Here, we seek to minimise the burden of high-throughput phenotyping of stomatal traits by introducing an automated method to identify and count stomata from micrographs. Although automated phenotyping methods using computer vision have been developed (Higaki *et al.*, 2014; Laga *et al.*, 2014; Duarte *et al.*, 2017; Jayakody *et al.*, 2017), these highly specialised approaches require feature engineering specific to a collection of images. These methods do not transfer well to images created with novel imaging and processing protocols. Additionally, tuning hand-crafted methods to work on a general set of conditions is cumbersome and often impossible to achieve. Recent applications of deep learning techniques to stomata counting have demonstrated success (Bhugra *et al.*, 2018; Aono *et al.*, 2019), but are not easily accessible to the public and use training data sets sampled from few taxa.

Deep convolutional neural networks (DCNN) circumvent specialised approaches by training the feature detector along with the classifier (LeCun *et al.*, 2015). The method has been widely successful for a range of computer vision problems in biology such as medical imaging (see Shen *et al.* (2017) for a review) or macroscopic plant phenotyping (Ubbens & Stavness, 2017). The main caveat of deep learning methods is that large numbers of parameters have to be trained in the feature detector. Improvements in network structure (He *et al.*, 2016) and training procedure (Simonyan & Zisserman, 2014) have helped training of large networks that incorporate gradient descent learning methods and prove to be surprisingly resilient against overfitting (Poggio *et al.*, 2017). Nevertheless, a large number of correctly annotated training images is still required to allow the optimiser to converge to a correct feature representation. Large labelled training sets such as the ImageNet database exist (Deng *et al.*, 2009), but for a highly specialised problems, such as stomata identification, publicly available datasets are not available at the scale required to train a typical DCNN.

We solve these problems by creating a large and taxonomically diverse training dataset of plant cuticle micrographs and by creating a network with a human-in-the-loop approach. Our

development of this method is provided to the public as a web-based tool called StomataCounter, which allows plant biologists to rapidly upload plant epidermal image datasets to pre-trained networks and then annotate stomata on cuticle images when desired. We applied this tool to a the training dataset and achieved robust identification and counts of stomata on a variety of angiosperm and pteridosperm taxa.

Materials and Methods

Biological material

Micrographs of plant cuticles were collected from four sources: the Cuticle database (<https://cuticledb.eesi.psu.edu/>, Barclay *et al.*, 2007); a *Ginkgo* common garden experiment (Barclay & Wing, 2016); a new intraspecific collection of balsam poplar (*Populus balsamifera*); and a new collection from living material at the Smithsonian, National Museum of Natural History (USNM) and the United States Botanic Garden (USBG) (Table 1). Specimens in the cuticle database collection were previously prepared by clearing and staining leaf tissue and then imaged. The entire collection of the cuticle database was downloaded on 16 November 2017. Downloaded images contained both the abaxial and adaxial cuticles in a single image, and were automatically separated with a custom bash script. Abaxial cuticle micrographs were discarded if no stomata were visible or if the image quality was so poor that no stomata could be visually identified by a human. *Ginkgo* micrographs were prepared by chemically separating the upper and lower cuticles and imaging with an environmental scanning electron microscopy (SEM) Barclay & Wing (2016). To create the poplar dataset, dormant cuttings of *P. balsamifera* genotypes were collected across the eastern and central portions of its range in the United States and Canada by S. R. Keller *et al.* and planted in common gardens in Vermont and Saskatchewan in 2014. Fresh leaves were sampled from June to July 2015 and immediately placed in plastic bags and then a cooler for temporary storage, up to 3 h. In a laboratory, nail polish (Sally Hansen, big kwik dry top coat no. 42494) was applied to a 1 cm² region of the adaxial and abaxial leaf surfaces, away from the mid-vein, and allowed to dry for *c.* 20 min. The dried cast of the epidermal surface was lifted with clear tape and mounted onto a glass slide. The collection at the USNM and

Table 1 Description of the datasets used for training and testing the network.

Dataset	Training		Test		Preparation method	Imaging method	Magnification	Z-stacks applied	Citation
	<i>N</i> _{images}	<i>N</i> _{species}	<i>N</i> _{images}	<i>N</i> _{species}					
Poplar	3123	1	175	1	Nail polish	DIC	400	No	This paper
Ginkgo	408	1	200	1	Lamina peel	SEM	200	No	Barclay & Wing (2016)
Cuticle DB	678	613	696	599	Clear & stain	Brightfield	400	No	Barclay <i>et al.</i> (2007)
USNM/USBG	409	124	696	132	Nail polish	DIC, SEM	100, 200, 400	Yes	This paper
Aorta	—	—	116	1	Elastic-Van Gieson	Brightfield	40	No	Johnson & Cipolla (2017)
Breast cancer	—	—	58	1	Hematoxylin & eosin stain	Brightfield		No	Gelasca <i>et al.</i> (2008)
Total	4618	739	1941	735					

DIC, Differential interference contrast; SEM, scanning electron microscopy. 1Training SEM *n* = 15; 2Training set, *n* = 4.

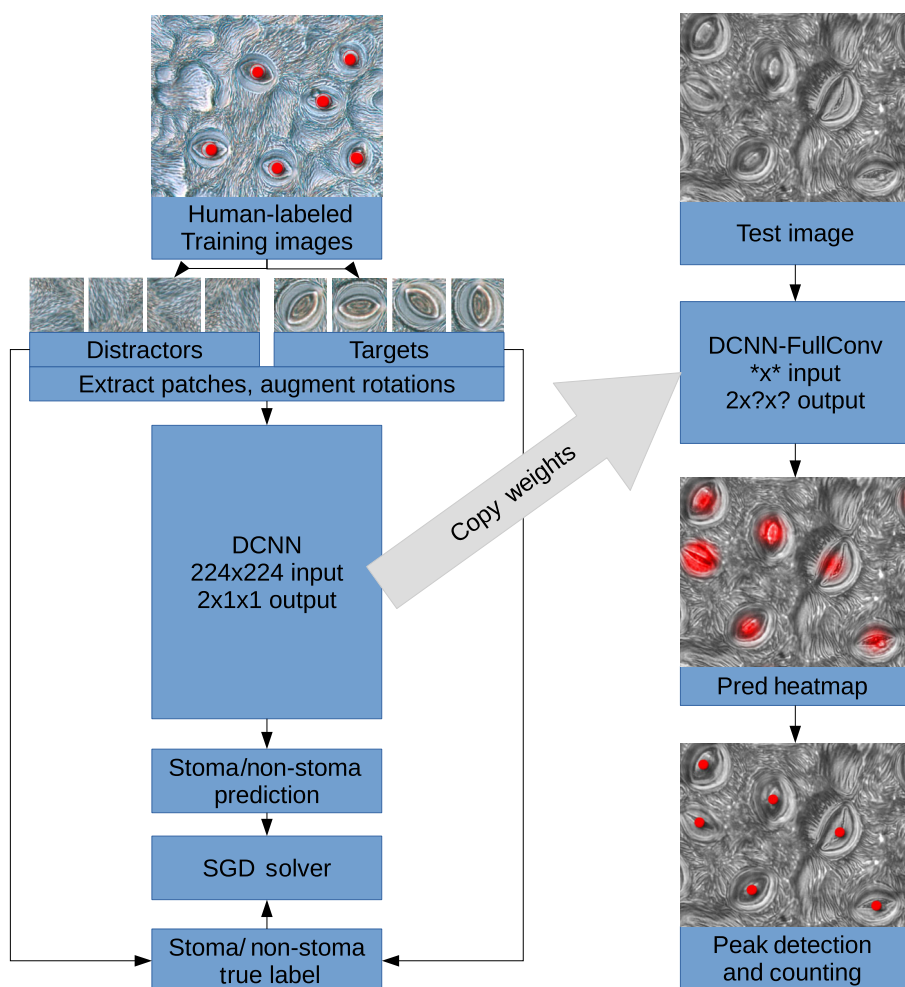


Fig. 1 Architecture of the deep convolutional neural network (DCNN) and classification tasks. Left: training and testing procedure. First column: target patches were extracted and centered around human-labeled stomata center positions; distractor patches were extracted in all other regions. A binary image classification network was trained. Second column: The image classification network was applied fully convolutional to the test image to produce a prediction heatmap. On the thresholded heatmap, peaks were detected and counted.

USBG was made similarly from opportunistically sampled tissues of the gardens around the National Museum of Natural History building. USBG collections were focused on the Hawaiian, southern exposure, and medicinal plant, and tropical collections. The taxonomic identity of each specimen was recorded according to the existing label next to each plant. The USNM/USBG collection was imaged with an Olympus BX-63 microscope using differential interference contrast (DIC). Some specimens had substantial relief and z-stacked images were created using cellSens image stacking software (Olympus Corp., Tokyo, Japan). Slides were imaged in three non adjacent areas per slide. The poplar collection was imaged with an Olympus BX-60 using DIC and each slide was imaged in two areas. Mounted material for the USNM/USBG and *Ginkgo* collections are deposited in the Smithsonian Institution, National Museum of Natural History in the Department of Paleobiology. Material for the poplar collection is deposited in the Keller laboratory in the Department of Plant Biology at the University of Vermont, USA. The training dataset totaled 4618 images (Table 1).

Deep convolutional neural network

We used a DCNN to generate a stomata likelihood map for each input image, followed by thresholding and peak detection to

localise and count stomata (Fig. 1). Because dense per-pixel annotations of stoma vs non stoma are difficult to acquire in large quantity, we trained a simple image classification DCNN based on the AlexNet structure instead (Krizhevsky *et al.*, 2012), and copied the weights into a fully convolutional network to allow per-location assessment of stomata likelihood. Although this method does not provide dense per-pixel annotations, the resulting resolution proved to be high enough to differentiate and count individual stomata.

We used pre trained weights for the lowest five convolutional layers from conv1 to conv5. The weights were taken from the ILSVRC image classification tasks (Russakovsky *et al.*, 2015) made available in the caffe net distribution (Jia *et al.*, 2014). All other layers were initialised using Gaussian initialisation with scale 0.01.

Training was performed using a standard stochastic gradient descent solver as in Krizhevsky *et al.* (2012), with learning rate 0.001 for pre trained layers and 0.01 for randomly initialised layers with a momentum of 0.9. Because the orientation of any individual stomata does not hold information for identification, we augmented data by rotating all training images into eight different orientations, applied random flipping and randomly positioned crop regions of the input size within the extracted 256×256 image patch. For distractors, we sampled patches from random image regions on human-annotated images that were at least 256 pixels distant from any labelled stoma. The

trained network weights were transferred into a fully convolutional network (Long *et al.*, 2015), which replaces the final fully connected layers by convolutions. To increase the resolution of the detector slightly, we reduced the stride of layer *pool5* from two to one, and added a dilation (Yu & Koltun, 2015) of two to layer *fc6_{conv}* to compensate. Due to margins (96 pixels) and stride (32 over all layers), application of the fully convolutional network to an image of size s yielded an output probability map p of size $s - (96 \times 2)/32$ along each dimension.

To avoid detecting low-probability stomata within the noise, the probability map p was thresholded and all values below the threshold $p\text{-thresh} = 0.98$ were set to zero. Local peak detection was run on a 3×3 pixel neighbourhood on the thresholded map, and each peak, excluding those located on a 96-pixel width border, was labelled as a stoma centre. This intentionally excludes stomata for which the detection peak is found near the border within the model margin to match the instructions given to human annotators. Resulting stoma positions were projected back onto the original image (see Fig. 2).

We built a user-friendly web service, StomataCounter, freely available at <http://stomata.science/>, to allow the scientific community easy access to the method. We are using a flask/jquery/bootstrap stack. Source codes for network training, as well as the webserver are available at <http://stomata.science/source>. To use

StomataCounter, users upload single.jpg images or a zip files of their .jpg images containing leaf cuticles prepared. Z-stacked image sets should be combined into a single image before uploading. A new dataset is then created where the output of the automatic counts, image quality scores and image metadata are recorded and can easily be exported for further analysis.

In addition to automatic processing, the user can manually annotate stomata and determine the empirical error rate of the automatic counts through a straightforward and intuitive web interface. The annotations can then be reincorporated into the training dataset to improve future performance of StomataCounter by contacting the authors and requesting retraining of the DCNN.

Statistical analyses

We tested the performance of the DCNN with a partitioned set of images from each dataset source. Whenever possible, images from a given species were used in either the training or test set, but not both, and only seven out of 1467 species are included in both. In total, 1941 images were used to test the performance of the network (Table 1). After running the test set through the network, stomata were manually annotated. If the centre of a stoma intersected the bounding box around the perimeter of the image, it was not counted.

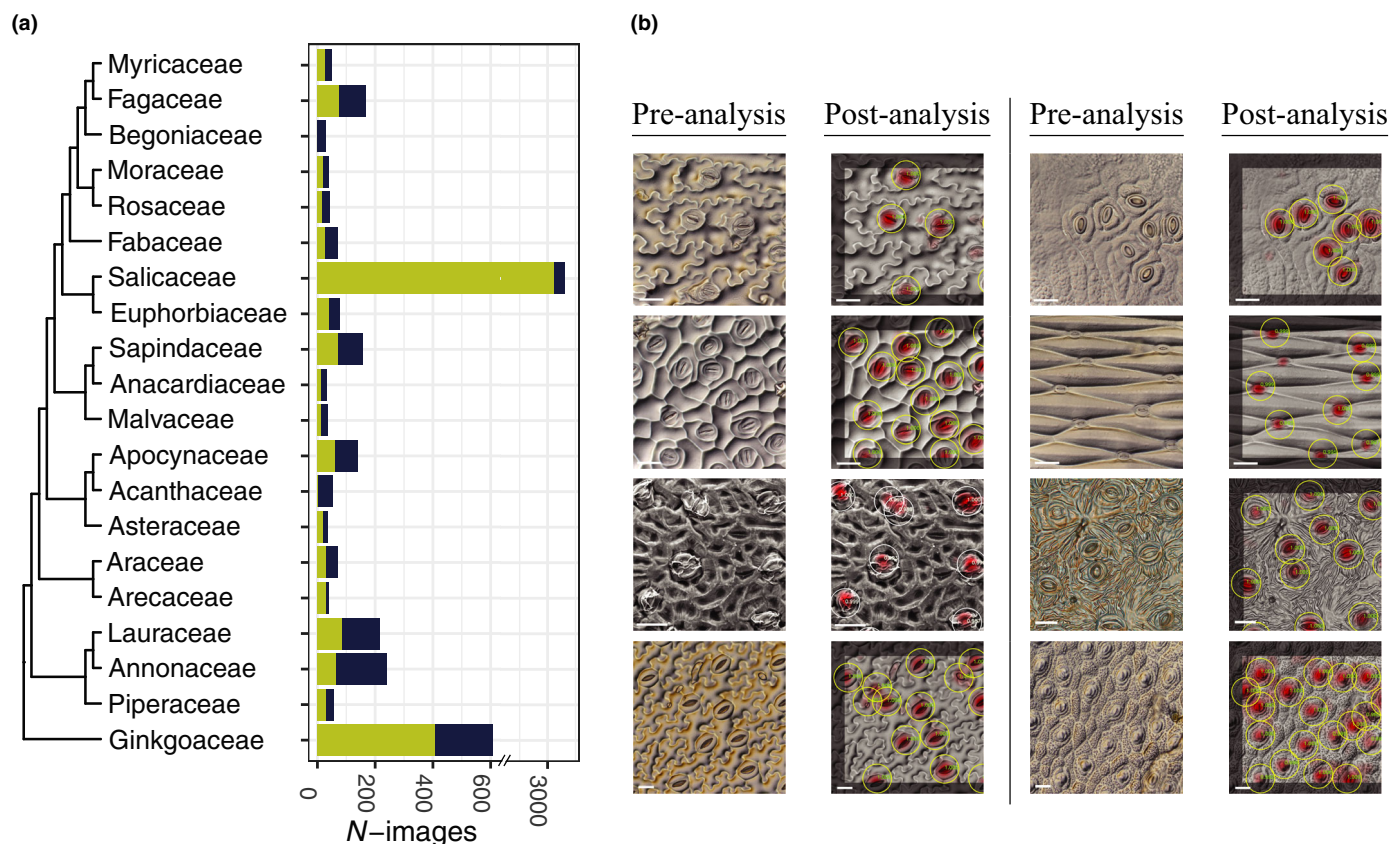


Fig. 2 Sample sizes of images for the top 20 families represented in the training (yellow) and test (blue) datasets (a). Examples of pre- and post-analysis images. A probability heatmap map is overlain onto the input image in the red channel. Detected stomata marked with circles with peak values given in green (b). Species in image pairs: *Adiantum peruvianum*, *Begonia loranthoides* (top row); *Chaemaeranthum gadacardii*, *Echeandia texensis* (second row); *Ginkgo biloba*, *Populus balsamifera* (third row); *Trillium luteum*, *Pilea libanensis* (bottom row). Bars: (top and second rows) 20 μ m, (row three) 25 μ m; (bottom row) 50 μ m.

To evaluate the DCNN, we first determined if it could identify stomata when they are known to be present and fail to identify them when they are absent. To execute this test, a set of 25 randomly selected abaxial plant cuticle micrographs containing stomata was chosen from each of the four datasets for a total of 100 images. To create a set of test images known to lack stomata, 100 adaxial cuticle micrographs were randomly sampled from the cuticle database. Visual inspection confirmed that none of the adaxial images contained stomata. Micrographs of thoracic aorta from an experimental rat model of preeclampsia (Johnson & Cipolla, 2017), and breast cancer tissue micrographs (Gelasca *et al.*, 2008) were used as negative controls from non-plant material. As a second test, we determined how well the method could identify stomata from stomatal patterning mutants in *Arabidopsis* that violate the single cell spacing rule. We sampled images from Papanatsiou *et al.* (2016) including 88 *mm1* and 90 wild-type micrographs, and used a paired *t*-test to determine if counts from manual or automatic were statistically different. We then determined if the precision was different between mutant and wild type accessions with ANOVA.

Identification accuracy is tested by applying the DCNN to a small image patch either centered on a stoma (target), or taken from at least 256 pixels distance to any labeled stomata (distractor). This yields true positive (N_{TP}), true negative (N_{TN}), false positive (N_{FP}) and false negative (N_{FN}) samples. We defined the classification accuracy A as:

$$A = \frac{N_{TP} + N_{TN}}{N_{TP} + N_{TN} + N_{FP} + N_{FN}}$$

We defined classification precision P as:

$$P = \frac{N_{TP}}{N_{TP} + N_{FP}}$$

We assumed the human counts contain only true positives and the automatic count contain true positives and false positives. As such, we defined precision as:

$$\text{Precision} = \log\left(\frac{\text{Human count}}{\text{Automatic count}}\right)$$

This definition of precision identifies over-counting errors as negative values and under-counting as positive values. This measure of precision is undefined if either manual or automatic count is zero, and 30 of the 1772 observations were discarded. These samples were either out of focus, lacked stomata entirely, or too grainy for human detection of were stomata. Classification recall R was defined as:

$$R = \frac{N_{TP}}{N_{TP} + N_{FN}}$$

Classification accuracy, precision, and recall, were calculated from groups of images constructed to span the diversity of imaging capture methods (that is brightfield, DIC, and SEM) and

magnification (that is $\times 200$ and $\times 400$) to determine how well the training set from one group transfers to identifying stomata in another group. Groups used for transfer assessment were the cuticle database, the *Ginkgo* collection, micrographs imaged at $\times 200$, at $\times 400$, and the combined set of images.

We use linear regression to understand the relationship between human and automatic stomata counts. For the purpose of calculating error statistics, we considered deviation from the human count attributable to error in the method. Images were partitioned for linear models by collection source, higher taxonomic group (that is Eudicot, Monocot, Magnoliid, Gymnosperm, Fern, or Lycophyte), and magnification. To understand how different sources of variance contribute to precision variance, we collected data on the taxonomic family, magnification, imaging method, and three measures of image quality. The taxonomic family of each image was determined using the open tree of life taxonomy accessed with the *rotl* package (Michonneau *et al.*, 2016). Two image quality measures used were based on the properties of an image's power-frequency distribution tail and described the standard deviation (fSTD) and mean (fMean) of the tail. Low values of fSTD and fMean indicate a blurry image, while high values indicate non blurry images. The third image quality measure, tEntropy, is a measure an image's information content. High entropy values indicated high contrast/noisy images while low values indicate lower contrast. These image quality measures were created with *PyImq* (Koho *et al.*, 2016) and standardised between zero and one. Random effects linear models were created with the R package *LME4* (Bates *et al.*, 2014) by fitting log precision to taxonomic family, magnification, and imaging method as factors. Linear models were fit for the scaled error and image quality scores. We used the root mean square error (RMSE) of the model residuals to understand how the factors and quality scores described the variance of log precision. Higher values of RMSE indicated larger residuals. Statistical analyses were conducted in PYTHON and R (R Core Team, 2013). Stomata counts, image quality scores, and taxonomy of training and test set micrographs are provided in Supporting information Table S1. Images used for the training and test sets are available for download from Dryad (<https://doi.org/10.5061/dryad.kh2gv5f>).

Data availability

Training and test set micrographs, model weights, and caffeNet model definition protocol are available as downloads from Dryad (<https://doi.org/10.5061/dryad.kh2gv5f>).

Results

Stomata detection

StomataCounter was able to accurately identify and count stomata when they were present in an image. False positives were detected in the adaxial cuticle, aorta, and breast cancer cell image sets at low frequency (Fig. S1). The mean number of stomata detected in the adaxial, aorta, and breast cancer image sets was

1.5, 1.4 and 2.4, respectively, while the mean value of the abaxial set containing stomata was 24.1.

Stomata measured from micrographs of mutant and wild-type *Arabidopsis* by human and automatic counting showed few differences between counting method (Fig. S2). A paired T-test of human and automatic counts failed to reject the null hypothesis of no difference in means ($t = -0.91$, $df = 175$, $P\text{-value} = 0.37$), but precision was different between wild-type and mutant genotypes (ANOVA: $F\text{-value} = 28.82$, $P\text{-value} = 2.51e^{-07}$). The variance of precision in Col-0 accessions ($\sigma^2 = 0.189$) was higher than the mutant accessions ($\sigma^2 = 0.153$), possibly due to image capture settings on the microscope.

Correspondence between automated and human stomata counting varied among the respective sample sets. There was close agreement among all datasets to the human count, with the exception of some of the samples imaged at $\times 200$ magnification (Fig. 3a–c). In these samples, the network tended to under count relative to human observers. Despite the variation among datasets and the large error present in the $\times 200$ dataset, the slopes of all models were close to 1 (Table 2). The $\times 400$, *Ginkgo*, and cuticle database sets all performed well at lower stomata counts, as indicated by their proximity to the expected one-to-one line and decreased in precision as counts increased.

Precision has a non-linear relationship with image quality, and changes in variance of precision are correlated with variation of image quality. With a ratio of 17 : 1 images in the training vs test sets, the poplar dataset had the best precision, followed by the *Ginkgo*, Cuticle database, and USNM/USBG datasets (Fig. 3d). Among the different imaging methods, SEM had the best precision, followed by DIC, and finally and brightfield microscopy (Fig. 3e). The variance of precision is higher for images with $\times 200$ magnification (Fig. 3f). RMSE values were lowest for taxonomic family and the family : magnification interaction, suggesting these factors contributed less to deviations between human and automated stomata counts than image quality or imaging method (Table 3).

Classification accuracy

The peak accuracy (94.2%) on the combined test sets is achieved when all training sets are combined (Fig. 6, see later). The combined dataset performs best on all test subsets of the data; that is, adding additional training data – even from different sets – is always beneficial for the generalisation of the network. Accuracy from train to test within a single species is higher (for example *Ginkgo* training for *Ginkgo* test at 97.4%) than transfer within datasets with a large number species across families ($\times 400$ training to $\times 400$ test: 85.5%).

The network does not generalise well between vastly different scales, that is the $\times 200$ dataset, which contains images down-scaled to half the image width and height. In this case, only training within the same scale achieved high accuracy (97.3%), while adding additional samples from the larger scale reduces the performance (to 90.7%).

Precision values are generally higher than recall (0.99 precision on the combined training and test sets; 0.93 recall, Fig. S3),

which shows that we mostly miss stomata rather than misidentifying non-existing stomata.

Increased training size is correlated with increased accuracy (Fig. 4), and providing a large number of annotated images is beneficial, as it lifts training accuracy from 72.8% with a training set of 10 images to 94.2% with the complete set of training images.

Discussion

Stomata are an important functional trait to many fields within plant biology, yet manual phenotyping of stomata counts is a laborious method that has few controls on human error and reproducibility. We created a fully automatic method for counting stomata that is both highly sensitive and reproducible, allows the user to quantify error in their counts, but is also entirely free of parameter optimisation from the user. Furthermore, the DCNN can be iteratively retrained with new images to improve performance and adjust to the needs of the community. This is a particular advantage of this method for adjusting to new taxonomic sample sets. However, new users are not required to upload new image types or images from new species for the method to work on their material. The pre loaded neural network was specifically trained on diverse set of image types and from many species ($n = 739$).

As the complexity of processing pipelines in biological studies increases, reproducibility of studies increasingly becomes a concern (Vieland, 2001). Apart from the reduced workload, automated image processing provides better reproducibility than manual stomata annotations. For instance, if multiple experts count stomata, they may not agree, causing artificial differences between compared populations. This includes how stomata at the edge of an image are counted, and what to do with difficult to identify edge cases. Automatic counters will have an objective measure, and introduce no systematic bias between compared sets as long as the same model is used. Additionally, our human counter-missed stomata that the machine detected (Fig. 5).

Our method is not the first to identify and count stomata. However, previous methods have not been widely adopted by the community and a survey of recent literature indicates manual counting is the predominant method (Takahashi *et al.*, 2015; Peel *et al.*, 2017; Liu *et al.*, 2018; Morales-Navarro *et al.*, 2018; Sumathi *et al.*, 2018). Previous methods have relied on substantial image pre-processing to generate images for thresholding to isolate stomata for counting (Oliviera *et al.*, 2014; Duarte *et al.*, 2017). Thresholding can perform well in a homogenous collection of images, but quickly fails when images collected by different preparation and microscopy methods are provided to the thresholding method (K. Fetter, pers. obs.). Some methods also require the user to manually segment stomata and subsequently process those images to generate sample views to supply to template matching methods (Laga *et al.*, 2014). Object-oriented methods (Jian *et al.*, 2011) also require input from the user to define model parameters. These methods invariably requires the user to participate in the counting process to tune parameters and monitor the image processing, and are not fully autonomous.

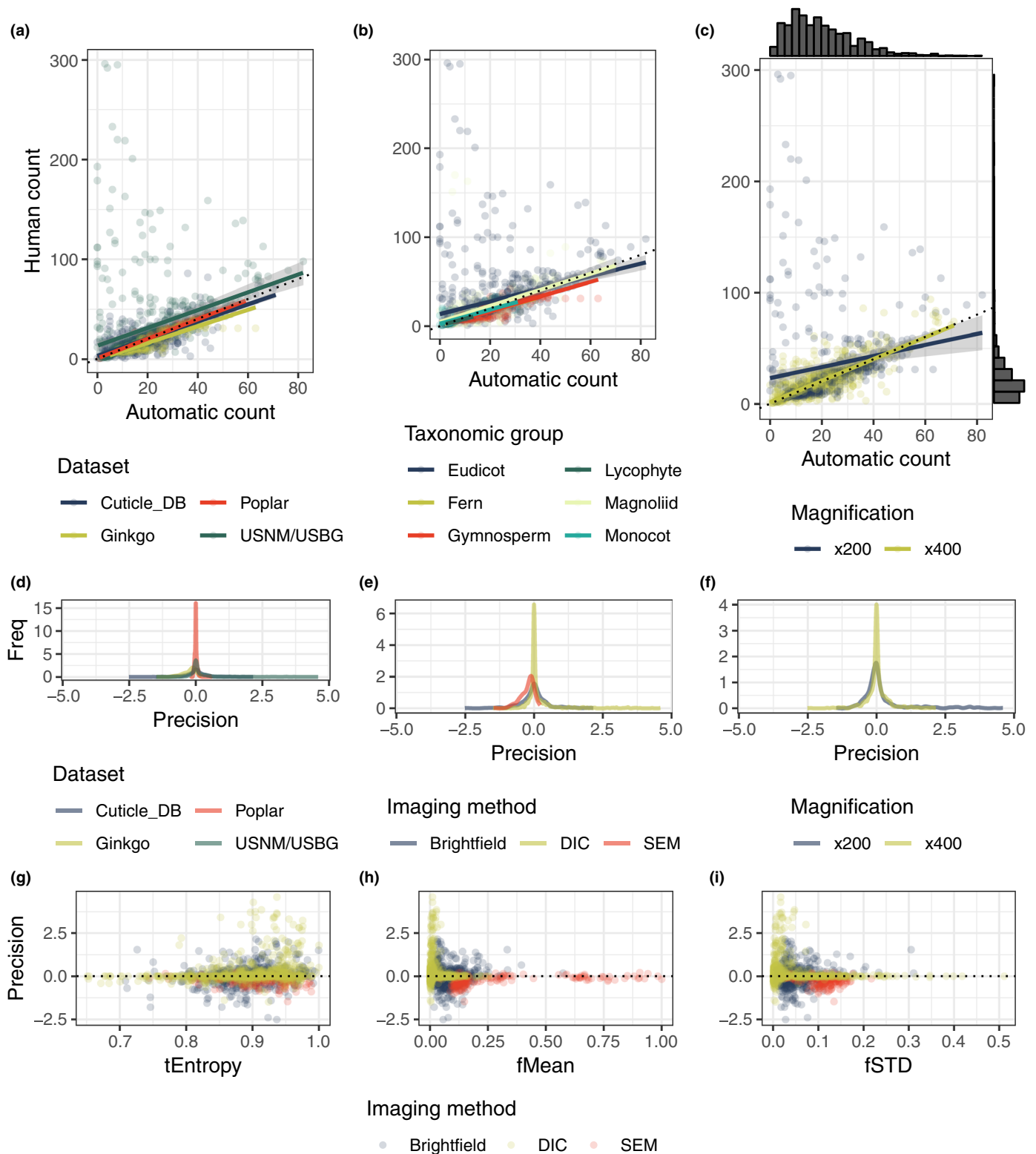


Fig. 3 Results of human and automatic counts and precision for each dataset organized by collection source (a, d), taxonomic group (b), and magnification (c), and imaging method (e). Precision and image quality scores have non-linear relationships, and changes in variance correlate to image quality (g–i). Positive values of precision indicate undercounting of stomata relative to human counts, while negative values indicate overcounting. The dotted black line is the 1 : 1 line, indicating where perfect automatic counts occur in a plot. See Table 2 for model summaries. tEntropy, image entropy; fMean, mean power-frequency tail distribution; fSTD, standard deviation of power-frequency tail distribution; DIC, Differential interference contrast; SEM, scanning electron microscopy.

Table 2 Summary of linear model fit parameters in Fig. 3 for different test datasets.

(a) Dataset	Cuticle DB	Poplar	<i>Ginkgo</i>	USNM/USBG	USNM/USBG [†]
<i>a</i>	3.315 ^{1***}	0.823	−0.364	10.068 ^{1***}	−0.695
SE _a	(0.635)	(0.456)	(0.677)	(2.041)	(0.5)
<i>s</i>	0.853 ^{1***}	0.978 ^{1***}	0.832 ^{1***}	0.997 ^{1***}	1.167 ^{1***}
SE _s	(0.029)	(0.014)	(0.028)	(0.089)	(0.026)
<i>r</i> ²	0.578	0.964	0.812	0.157	0.835

(b) Taxonomic group	Eudicot	Monocot	Magnoliid	Gymnosperm	Fern	Lycophyte
<i>a</i>	10.405 ^{1***}	2.360 ^{1***}	4.360 ^{1*}	0.004	0.624	0.771
SE _a	(1.643)	(0.598)	(1.844)	(0.735)	(0.474)	(1.414)
<i>s</i>	0.808 ^{1***}	0.808 ^{1***}	0.912 ^{1***}	0.830 ^{1***}	0.799 ^{1***}	0.751 ^{1**}
SE _s	(0.063)	(0.056)	(0.086)	(0.031)	(0.039)	(0.169)
<i>r</i> ²	0.141	0.617	0.272	0.778	0.938	0.739

(c) Magnification	×200		×400	
<i>a</i>	18.346 ^{1***}		1.446 ^{1***}	
SE _a	(3.192)		(0.381)	
<i>s</i>	0.654 ^{1***}		0.970 ^{1***}	
SE _s	(0.123)		(0.017)	
<i>r</i> ²	0.055		0.740	

Dataset definitions given in text. *a*, *y*-intercept; SE, standard error. Significance indicated by: *, $P < 0.05$; **, $P < 0.01$; ***, Root mean square error (RMSE) of the model residuals. $P < 0.001$. [†] ×200 images removed from this USNM/USBG set.

Table 3 Root mean square error (RMSE) of the model residuals.

Model	RMSE
~Family	0.463
~Magnification	0.523
~Family : magnification	0.399
~Imaging method	0.519
~fMean	0.523
~fSTD	0.523
~tEntropy	0.521
~tEntropy : imaging method	0.515

Lower RMSE values suggest a better fit of the model. The response in each model was precision.

fMean, mean power-frequency tail distribution; fSTD, standard deviation of power-frequency tail distribution.

More recently, cascade classifier methods have been developed which perform well on small collections of test sets (Violet-Chabrand & Brendel, 2014; Jayakody *et al.*, 2017). Additionally, most methods rely on a very small set of images (50–500) typically sampled from just a few species or cultivars to create the training and/or test set (for example Bhugra *et al.*, 2018).

Apart from generalisation concerns, several published methods require the user to have some experience coding in PYTHON or C++, a requirement likely to reduce the potential pool of end users. Our method resolves these issues by being publicly available, fully autonomous of the user, who is only required to upload jpeg formatted images, is free of any requirement for the user to code, and is trained on a relatively large and taxonomically diverse

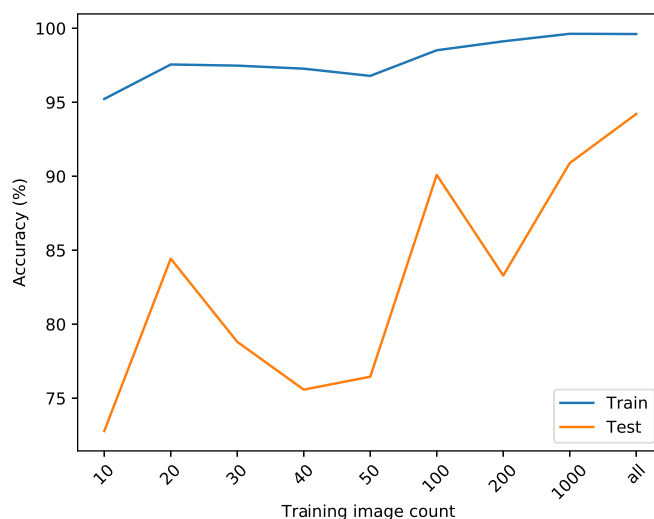


Fig. 4 Accuracy by training image count and the effect of increasing the training set size on classification accuracies. Since classification is a binary task, chance level is at 50%. Training image count 'all' includes all 4618 annotated training images.

set of cuticle images. Users may wish to set up StomataCounter locally on their own servers, and we have integrated the PYTHON scripts into an offline program users can run at the command line. The command line version is available in the github repository: <https://github.com/SvenTwo/epidermal>. Users can generate model weights from a custom set of training images and supply it to the image processing script for stomatal identification, or use the predefined weights generated from this work.

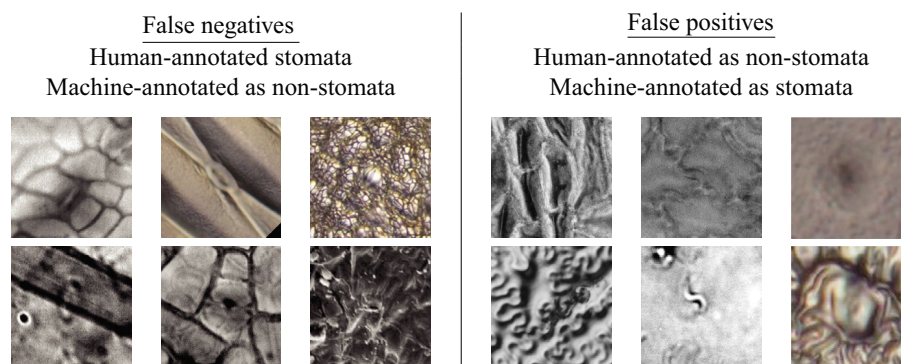


Fig. 5 Samples that were mislabeled with high confidence by either the machine or human. False negatives (left panel) from stomata that were detected by the human, but not by the machine. Image features typically generating false negatives are blurry images, artifacts obstruction the stomata, low contrast between epidermal and guard cells, and very small scale stomata. False positives (right panel) are image features that are labeled by the machine as stomata, but not by the human counter, or are missed by the human counter (e.g. top and bottom right images in the panel). Errors typically generating false positives are image artifacts that superficially resemble stomata, particularly shapes mimicking interior guard cell structure.

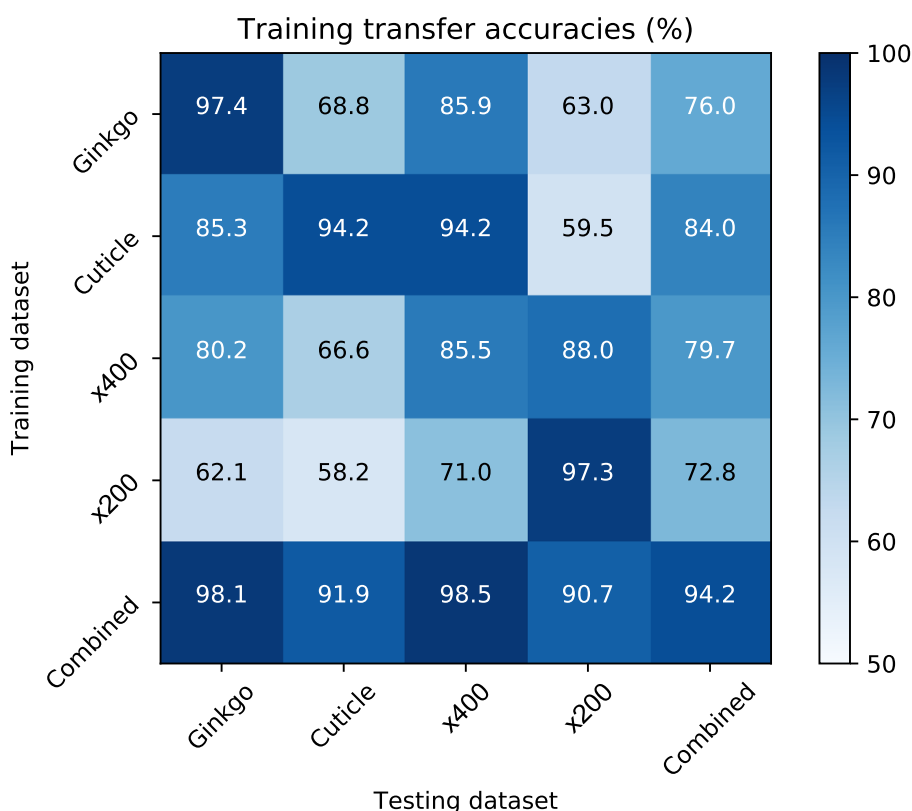


Fig. 6 Accuracies for models trained on different training datasets (vertical), tested on different test datasets. Combined is a union of all training and test sets. For precision and recall values, see Supporting information Fig. S3.

We have demonstrated that this method is capable of accurately identifying stomata when they are present, but false positives may still be generated by shapes in images that approximate the size and shape of stomata guard cells. Conversely, false negatives are generated when a stomata is hidden by a feature of the cuticle or if poor sample preparation/imaging introduces blur. This issue is likely to be avoidable through increased sample size of the training image set and good sample preparation and microscopy techniques by the end user. The method is also able to detect stomata from mutant *Arabidopsis* genotypes that violate the single spacing rule.

The importance of having a well-matched training and testing image set was apparent at $\times 200$, where there was a subset of observations with low transfer accuracy (Fig. 6), and StomataCounter consistently under counted relative to human observation (Fig. 3). We argue that transferring architecture between scales is not advisable and images should be created by the user to match the predominant size (2048×2048) and magnification ($\times 400$) of images in the training network. Our training set of images spanned 82 different families and was over-represented by angiosperms. Stomata in gymnosperms are typically sunken into pores that make it difficult to obtain good nail polish casts. Models tested to explain variation in scaled error

revealed that taxonomic family and its interaction with magnification were the factors that had the best explanatory power for scaled error. Future collections of gymnosperm cuticles could be uploaded to the DCNN to retrain it in order to improve the performance of the method for gymnosperms.

More generally however, this highlights how users will need to be thoughtful about matching of training and test samples for taxa that may deviate in stomata morphology from the existing reference database. We therefore recommend that users working with new or morphologically divergent taxa first run several pilot tests with different magnification and sample preparation techniques to find optimal choices that minimize error for their particular study system. SEM micrographs had the least amount of error, followed by DIC, and finally brightfield (Fig. 3g–i). Lastly, image quality was strongly related to log precision; predictably, images that are too noisy (that is high entropy) and out-of-focus (low fMean or fSTD) will generate higher error. Obtaining high quality, in-focus images should be a priority during data acquisition. We provide these guidelines for using the method, and recommend that users read these guidelines before collecting a large quantity of images:

- Collect sample images using different microscopy methods from the same tissues. We recommend an initial collection of 25–100 images before initiating a new large-scale study.
- Run images through StomataCounter.
- Establish a ‘true’ stomata count using the annotation feature.
- Regress image quality scores (automatically provided in output csv file) against log precision.
- Regress human vs automatic counts and assess error.
- Choose the microscopy method that minimises error and image the remaining samples.

Different microscopy methods can include using DIC or phase contrast filters, adjusting the aperture to increase contrast, or staining tissue and imaging under fluorescent light (see Eisele *et al.*, 2016 for more suggestions). If a large collection of images is already available and re-imaging is not feasible, we recommend the users take the following actions:

- Randomly sample 100 images.
- Upload the images to StomataCounter.
- Annotate images to establish the ‘true’ count.
- Explore image quality scores with against the log(precision) to determine a justifiable cut-off value for filtering images.
- Discard images below the image quality cut-off value.

New users can also contact the authors through the StomataCounter web interface and we may retrain the model to include the 100 annotated images.

Fast and accurate counting of stomata increases productivity of workers and decreases the time from collecting a tissue to analyzing the data. Until now, assessing measurement error required phenotyping a reduced set of images multiple times by, potentially, multiple counters. With StomataCounter, users can instantly phenotype their images and annotate them to create empirical error rates. The open source code and flexibility of using new and customized training sets will make

StomataCounter and important resource for the plant biology community.



Acknowledgements

We thank Kyle Wallick at the United States Botanic Garden for facilitating access to the living collections of the Garden. Scott Whitaker aided in imaging cuticle specimens with SEM and DIC microscopy in the Laboratories of Analytical Biology at the Smithsonian Institution, National Museum of Natural History. Rat aorta images were provided by Dr A. Chapman and Dr M. Cipolla at the University of Vermont. Arabidopsis stomatal patterning mutants were provided by Dr Maria Papanatsiou and Dr Michael Blatt. Dr Terry Delaney at the University of Vermont kindly allowed KCF to use his microscope for imaging. Funding to create Stomata Counter was provided by an NSF grant to SRK (IOS-1461868) and a Smithsonian Institution Fellowship to KCF.

Author contributions

The research was conceived and performed by KCF and SE. The website and PYTHON scripts were written by SE. Data were collected and analysed by KCF. *Ginkgo* Images were submitted by RSB and SW. All authors interpreted the results. The manuscript was written by KCF, SE, and SRK. All authors edited and approved the manuscript. SE and KCF contributed equally to this work.

ORCID

Rich S. Barclay  <https://orcid.org/0000-0003-4979-6970>
 Sven Eberhardt  <https://orcid.org/0000-0002-7535-3551>
 Karl C. Fetter  <https://orcid.org/0000-0002-9234-9300>
 Stephen R. Keller  <https://orcid.org/0000-0001-8887-9213>
 Scott Wing  <https://orcid.org/0000-0002-2954-8905>

References

- Aono A, Nagai J, Dickel G, Marinho R, Oliveira P, Faria F. 2019. A stomata classification and detection system in microscope images of maize cultivars. *bioRxiv*: 538165.
- Barclay R, McElwain J, Dilcher D, Sageman B. 2007. The cuticle database: developing an interactive tool for taxonomic and paleoenvironmental study of the fossil cuticle record. *Courier-Forschungsinstitut Senckenberg* 258: 39.
- Barclay R, Wing S. 2016. Improving the *Ginkgo* CO₂ barometer: implications for the early Cenozoic atmosphere. *Earth and Planetary Science Letters* 439: 158–171.
- Bates DM, Maechler M, Bolker BM, Walker S. 2014. Fitting linear mixed-effects models using lme4. *Journal of Statistical Software* arXiv: 1406.5823.
- Bergmann DC, Sack FD. 2007. Stomatal development. *Annual Review of Plant Biology* 58: 163–181.
- Bhugra S, Mishra D, Anupama A, Chaudhury S, Lall B, Chugh A, Chinnusamy V. 2018. Deep convolutional neural networks based framework for estimation of stomata density and structure from microscopic images. In: *European Conference on Computer Vision*. Munich, Germany: (ECCV) Workshops, 412–423.
- Deng J, Dong W, Socher R, Li L-J, Li K, Li F-F. 2009. ImageNet: a large-scale hierarchical image database. In: *2009 IEEE Conference on Computer Vision and Pattern Recognition*, 248–255.
- Duarte KT, de Carvalho MAG, Martins PS. 2017. *Segmenting high-quality digital images of stomata using the wavelet spot detection and the watershed*

- transform. In: *VISIGRAPP (4: VISAPP)*. Setubal, Portugal: Science and Technology Publications, Lda, 540–547.
- Eisele JF, Fafiler F, Burgel PF, Chaban C. 2016. A rapid and simple method for microscopy-based stomata analyses. *PLoS ONE* 11: 1–13.
- Fischer R, Rees D, Sayre K, Lu ZM, Condon A, Saavedra AL. 1998. Wheat yield progress associated with higher stomatal conductance and photosynthetic rate, and cooler canopies. *Crop Science* 38: 1467–1475.
- Gelasca ED, Byun J, Obara B, Manjunath BS. 2008. Evaluation and benchmark for biological image segmentation. In: *2008 15th IEEE International Conference on Image Processing*. 1816–1819.
- He K, Zhang X, Ren S, Sun J. 2016. Deep residual learning for image recognition. In: *Proceedings of the IEEE conference on computer vision and pattern recognition*. Los Alamitos, CA, USA: IEEE Computer Society, 770–778.
- Hetherington AM, Woodward FI. 2003. The role of stomata in sensing and driving environmental change. *Nature* 424: 901–908.
- Higaki T, Kutsuna N, Hasezawa S. 2014. CARTA-based semi-automatic detection of stomatal regions on an *Arabidopsis* cotyledon surface. *Plant Morphology* 26: 9–12.
- Hudson ME. 2008. Sequencing breakthroughs for genomic ecology and evolutionary biology. *Molecular Ecology Resources* 8: 3–17.
- Jayakody H, Liu S, Whitty M, Petrie P. 2017. Microscope image based fully automated stomata detection and pore measurement method for grapevines. *Plant Methods* 13: 94.
- Jia Y, Shelhamer E, Donahue J, Karayev S, Long J, Girshick R, Guadarrama S, Darrell T. 2014. Caffe: convolutional architecture for fast feature embedding. *arXiv preprint arXiv 1408.5093*.
- Jian S, Zhao C, Zhao Y. 2011. Based on remote sensing processing technology estimating leaves stomatal density of *Populus euphratica*. 2011 *IEEE International Geoscience and Remote Sensing Symposium*. IEEE, 547–550.
- Johnson AC, Cipolla MJ. 2017. Altered hippocampal arteriole structure and function in a rat model of preeclampsia: potential role in impaired seizure-induced hyperemia. *Journal of Cerebral Blood Flow & Metabolism* 37: 2857–2869.
- Kim TH, Bohmer M, Hu H, Nishimura N, Schroeder JI. 2010. Guard cell signal transduction network: advances in understanding abscisic acid, CO₂, and Ca²⁺ signaling. *Annual Review of Plant Biology* 61: 561–591.
- Koho S, Fazeli E, Eriksson JE, Hanninen PE. 2016. Image quality ranking method for microscopy. *Scientific Reports* 6: 28962.
- Krizhevsky A, Sutskever I, Hinton GE. 2012. Imagenet classification with deep convolutional neural networks. *Advances in Neural Information Processing Systems* 00: 1097–1105.
- Laga H, Shahinnia F, Fleury D. 2014. Image-based plant stomata phenotyping. In: *2014 13th International Conference on Control Automation Robotics and Vision, ICARCV 2014*, 217–222.
- LeCun Y, Bengio Y, Hinton G. 2015. Deep learning. *Nature* 521: 436.
- Liu C, He N, Zhang J, Li Y, Wang Q, Sack L, Yu G. 2018. Variation of stomatal traits from cold to temperate to tropical forests and association with water-use efficiency. *Functional Ecology* 32: 20–28.
- Long J, Shelhamer E, Darrell T. 2015. Fully convolutional networks for semantic segmentation. In: *Proceedings of the IEEE conference on computer vision and pattern recognition*. Los Alamitos, CA, USA: IEEE Computer Society, 3431–3440.
- McElwain JC, Steinthorsdottir M. 2017. Paleoecology, ploidy, paleoatmospheric composition, and developmental biology: a review of the multiple uses of fossil stomata. *Plant Physiology* 174: 650–664.
- McKown AD, Guy RD, Quamme L, Klapset J, La Mantia J, Constabel C, El-Kassaby YA, Hamelin RC, Zifkin M, Azam M. 2014. Association genetics, geography and ecophysiology link stomatal patterning in *Populus trichocarpa* with carbon gain and disease resistance trade-offs. *Molecular Ecology* 23: 5771–5790.
- Michonneau F, Brown JW, Winter DJ. 2016. rotl: an R package to interact with the Open Tree of Life data. *Methods in Ecology and Evolution* 7: 1476–1481.
- Morales-Navarro S, Perez-Diaz R, Ortega A, de Marcos A, Mena M, Fenoll C, Gonzalez-Villanueva E, Ruiz-Lara S. 2018. Overexpression of a SDD1-Like gene from wild tomato decreases stomatal density and enhances dehydration avoidance in *Arabidopsis* and cultivated tomato. *Frontiers in Plant Science* 9: 940.
- Oliviera MWdS, da Silva NR, Casanova D, Pinheiro LFS, Kolb RM, Bruno OM. 2014. Automatic counting of stomata in epidermis microscopic images. *X Workshop de Visao Computacional* 3: 253–257.
- Papanatsiou M, Amtmann A, Blatt MR. 2016. Stomatal spacing safeguards stomatal dynamics by facilitating guard cell ion transport independent of the epidermal solute reservoir. *Plant Physiology* 172: 254–263.
- Peel JR, Mandujano Sanchez MC, Lopez Portillo J, Golubov J. 2017. Stomatal density, leaf area and plant size variation of *Rhizophora mangle* (Malpighiales: Rhizophoraceae) along a salinity gradient in the Mexican Caribbean. *Revista de Biologia Tropical* 65: 701–712.
- Poggio T, Kawaguchi K, Liao Q, Miranda B, Rosasco L, Boix X, Hidary J, Mhaskar H. 2017. Theory of deep learning III: explaining the non-overfitting puzzle. *arXiv:1801.00173*.
- R Core Team. 2013. *R: a language and environment for statistical computing*. Vienna, Austria: R Foundation for Statistical Computing.
- Royer DL. 2001. Stomatal density and stomatal index as indicators of paleoatmospheric CO₂ concentration. *Review of Palaeobotany and Palynology* 114: 1–28.
- Russakovsky O, Deng J, Su H, Krause J, Satheesh S, Ma S, Huang Z, Karpathy A, Khosla A, Bernstein M et al. 2015. Imagenet large scale visual recognition challenge. *International Journal of Computer Vision* 115: 211–252.
- Shen D, Wu G, Suk HI. 2017. Deep learning in medical image analysis. *Annual Review of Biomedical Engineering* 19: 221–248.
- Shimazaki K, Doi M, Assmann SM, Kinoshita T. 2007. Light regulation of stomatal movement. *Annual Review of Plant Biology* 58: 219–247.
- Simonyan K, Zisserman A. 2014. Very deep convolutional networks for large-scale image recognition. *arXiv:1409.1556*.
- Sumathi M, Bachpai V, Deeparaj B, Mayavel A, Dasgupta MG, Nagarajan B, Rajasugunasekar D, Sivakumar V, Yasodha R. 2018. Quantitative trait loci mapping for stomatal traits in interspecific hybrids of *Eucalyptus*. *Journal of Genetics* 97: 323–329.
- Takahashi S, Monda K, Negi J, Konishi F, Ishikawa S, Hashimoto-Sugimoto M, Goto N, Iba K. 2015. Natural variation in stomatal responses to environmental changes among *Arabidopsis thaliana* ecotypes. *PLoS ONE* 10: e0117449.
- Ubbens JR, Stavness I. 2017. Deep plant phenomics: a deep learning platform for complex plant phenotyping tasks. *Frontiers in Plant Science* 8: 1190.
- Viale-Chabrand S, Brendel O. 2014. Automatic measurement of stomatal density from microphotographs. *Trees* 28: 1859–1865.
- Vieland VJ. 2001. The replication requirement. *Nature Genetics* 29: 244.
- Wang R, Yu G, He N, Wang Q, Zhao N, Xu Z. 2015. Latitudinal variation of leaf stomatal traits from species to community level in forests: linkage with ecosystem productivity. *Nature Scientific Reports* 5: 1–11.
- Yu F, Koltun V. 2015. Multi-scale context aggregation by dilated convolutions. *arXiv:1511.07122*.

Supporting Information

Additional Supporting Information may be found online in the Supporting Information section at the end of the article.

Fig. S1 Stomata identification test results.

Fig. S2 Stomata identification in *Arabidopsis* cell-patterning mutants.

Fig. S3 Precision and recall of transferred training and testing datasets.

Table S1 Stomata counts, image quality scores, and taxonomy of training and test set images.

Please note: Wiley Blackwell are not responsible for the content or functionality of any Supporting Information supplied by the authors. Any queries (other than missing material) should be directed to the *New Phytologist* Central Office.