

1 Predicting photosynthetic pathway from anatomy using machine learning
2
3 Ian S. Gilman¹, Karolina Heyduk², Carlos A. Maya-Lastra², Lillian P. Hancock², and Erika J.
4 Edwards²
5

6 ¹Michigan State University, East Lansing, Michigan; ²The University of Connecticut, Storrs,
7 Connecticut; ³Yale University, New Haven, Connecticut
8

9 Author for correspondence:

10 *Ian S. Gilman*

11 *Email: gilmania@msu.edu*

12

13 ORCIDs:

14 Ian S. Gilman: 0000-0002-0390-9370

15 Karolina Heyduk: 0000-0002-1429-6397

16 Carlos A. Maya-Lastra: 0000-0002-0550-3331

17 Lillian P. Hancock: 0000-0002-1394-4970

18 Erika J. Edwards: 0000-0003-0515-2778

19

20 **Main text**

21 Introduction word count: 781

22 Materials and Methods word count: 1763

23 Results word count: 837

24 Discussion word count: 1783

25 Color figures: 6

26 Grayscale figures: 0

27 Tables: 0

28

29 **Supporting Information**

30 Supporting Figures: 9

31 Supporting Tables: 10

32 Supporting Methods: 1

33 Supporting Datasets: 1

34

35 SUMMARY

36 - Plants with Crassulacean acid metabolism (CAM) have long been associated with a specialized
37 anatomy, including succulence and thick photosynthetic tissues. Firm, quantitative boundaries
38 between non-CAM and CAM plants have yet to be established – if they indeed exist.
39 - Using novel computer vision software to measure anatomy, we combined new measurements
40 with published data across flowering plants. We then used machine learning and phylogenetic
41 comparative methods to investigate relationships between CAM and anatomy.
42 - We found significant differences in photosynthetic tissue anatomy between plants with
43 differing CAM phenotypes. Machine learning based classification was over 95% accurate in
44 differentiating CAM from non-CAM anatomy, and had over 70% recall of distinct CAM
45 phenotypes. Phylogenetic least squares regression and threshold analyses revealed that CAM
46 evolution was significantly correlated with increased mesophyll cell size, thicker leaves, and
47 decreased intercellular airspace.
48 - Our findings suggest that machine learning may be used to aid the discovery of new CAM
49 species and that the evolutionary trajectory from non-CAM to strong, obligate CAM requires
50 continual anatomical specialization.

51

52 **Keywords: (4-10)**

53 Asparagaceae, Crassulacean acid metabolism, machine learning, photosynthesis, Portullugo

54 INTRODUCTION

55 Carbon concentrating mechanisms increase the efficiency of photosynthesis by raising the
56 concentration of CO₂ inside photosynthetic tissues relative to the ambient environment. The most
57 common carbon concentrating mechanism, Crassulacean acid metabolism (CAM), was first
58 discovered because of marked physiological differences between succulent and nonsucculent
59 plants (de Saussure, 1804). Generally, CAM species conduct gas exchange at night to reduce
60 transpirational water loss; the nocturnally fixed carbon is stored as malic acid overnight and
61 released the next day behind closed stomata, thereby saturating photosynthetic tissues with CO₂
62 (Osmond, 1978). Although the co-occurrence of CAM and succulent anatomy is so consistent
63 that botanists have used it as a guide to find new CAM plants (Coutinho, 1969), quantitative
64 relationships between anatomy and CAM remain elusive.

65 CAM and succulence may be correlated because they are co-selected as adaptations to
66 water limitation. CAM species can be up to eightfold as water use efficient as C₃ species (Winter
67 *et al.*, 2005) and the water stored in succulent plants is essential for drought avoidance (Males,
68 2017). Although such a correlation does not necessarily imply that derived anatomy is a
69 prerequisite of, or is caused by, CAM evolution, there are at least two hypothesized direct
70 functional links between CAM and succulent anatomy. First, storage of nocturnally fixed CO₂ as
71 malic acid in mesophyll vacuoles may require large vacuoles in photosynthetic cells and
72 therefore larger, succulent mesophyll cells (Zambrano *et al.*, 2014; Töpfer *et al.*, 2020). Second,
73 increased succulence in mesophyll cells may lower intercellular air space (IAS) and therefore
74 mesophyll CO₂ conductance (g_m) (Nelson *et al.*, 2008; Cousins *et al.*, 2020); thus, increased
75 succulence may increase selection for CAM by lowering the efficiency of C₃ photosynthesis
76 (Nelson *et al.*, 2008; Edwards, 2019). It is also possible that the evolution of CAM does not
77 entail selection on succulence *per se*, but that the use of CAM reduces constraints on succulence
78 evolution by removing g_m limitations due to carbon concentration.

79 Quantitative studies of CAM and anatomy have generally been restricted to relatively few
80 taxa at the extremes of the CAM phenotypic spectrum, but have generally found positive
81 correlations between CAM and succulence. Individual studies have reported that CAM species
82 tend to have greater leaf thickness (LT) and larger mesophyll cell area (MA), but mixed trends

83 have been observed for IAS (Nelson *et al.*, 2005 2008; Zambrano *et al.*, 2014; Earles *et al.*,
84 2018; Luján *et al.*, 2022); however, a recent meta-analysis of these relationships found
85 inconsistent trends across clades (Herrera, 2020). Recently, hybrids between species with
86 different photosynthetic types have been used to study the relationships between CAM activity
87 and anatomical traits. In both *Yucca* (Agavoideae, Asparagaceae) (Heyduk *et al.*, 2020) and
88 *Cymbidium* (Orchidaceae) (Yamaga-Hatakeyama *et al.*, 2022), hybrids of C₃ x CAM crosses
89 possessed intermediate anatomical phenotypes and CAM activity. Within *Yucca* hybrid
90 genotypes, however, the correlations between CAM activity and anatomy decreased in
91 magnitude or disappeared entirely (Heyduk *et al.*, 2020).

92 The mosaic of past research provides limited insight into the evolution of CAM and
93 photosynthetic tissue anatomy because it has focused on the extremes of CAM phenotypes (i.e.,
94 non-CAM species and species that use CAM as their primary metabolism). However, there are
95 many recognized CAM phenotypes that differ in pattern and magnitude of CAM activity
96 (Winter, 2019). Here, we use term “CAM” to refer to all species capable of CAM, regardless of
97 strength or pattern of expression, and “minority-CAM” and “primary-CAM” to refer to species
98 that fix the minority and majority of CO₂ with CAM, respectively. Primary-CAM (pCAM) is
99 consistent with past definitions of “CAM plant” (Winter, 2019) and “strong CAM” (Edwards,
100 2019), while minority-CAM (mCAM) encompasses species that can facultatively use CAM or
101 constitutively use CAM at low levels, but primarily use C₃ or C₄ photosynthesis for CO₂
102 assimilation (mCAM = “C₃+CAM” of Edwards (2019), but with the inclusion of C₄+CAM
103 species). It is generally assumed that the evolution of pCAM requires transitioning through
104 mCAM (Hancock and Edwards, 2014; Yang *et al.*, 2015; Edwards, 2019), but the relative
105 timings of anatomical shifts during the evolution of mCAM and pCAM – and whether or not
106 mCAM species possess a specialized anatomy – remain open questions.

107 Here, we combined anatomical measurements from thousands of angiosperm species
108 from over 200 families to draw anatomical boundaries between non-CAM, mCAM, and pCAM
109 phenotypes. Using supervised machine learning models, we were able to classify CAM
110 phenotypes from anatomical measurements with moderate to high accuracy. Finally, in a detailed
111 study of the Portullugo clade (Fig. 1), we reconstructed the evolution of CAM and used
112 phylogenetic comparative methods to establish significant relationships between anatomy and

113 CAM evolution. Our findings support the hypothesis that CAM evolution entails anatomical
114 evolution and reveal nuances about the earliest stages of CAM evolution.

115

116 MATERIALS AND METHODS

117 Public anatomical data sets, taxon sampling, and specimen imaging

118 Publicly available data were gathered from the TRY (Fraser *et al.*, 2020) and BROT2
119 (Tavşanoğlu and Pausas, 2018) plant trait databases and individual studies of CAM anatomy in
120 Orchidaceae (Silvera *et al.*, 2005), Bromeliaceae (Males, 2018), Asparagaceae (Heyduk *et al.*,
121 2016), Caryophyllales (Ogburn and Edwards, 2012, 2013), Papua New Guinean epiphytes
122 (Earnshaw *et al.*, 1987), Clusiaceae (Luján *et al.*, 2022), and across angiosperms (Nelson *et al.*,
123 2008) (Supporting Information Table S1). These data contained observations of mesophyll cell
124 area (MA), leaf thickness (LT), mesophyll intercellular air space (IAS), leaf dry matter content
125 (LDMC), and specific leaf area per unit dry mass (SLA). We generated two new datasets of MA,
126 IAS, and LT for members of the Asparagaceae (subfamilies Agavoideae and Nolinoideae) and
127 Portullugo (including families Anacampserotaceae, Cactaceae, Didiereaceae, Montiaceae,
128 Molluginaceae, and Portulacaceae) (Supporting Information Table S2). In 2017, leaf cross
129 sections were taken from 15 Portullugo species grown at Brown University, Providence, RI.
130 Tissue sections were immediately placed in 10% neutral buffered formalin and sent to the
131 Veterinary Diagnostic Laboratories in the College of Veterinary Medicine at the University of
132 Georgia (Athens, GA) for fixation, embedding, and sectioning and staining with toluidine blue.
133 In the spring of 2019, we collected leaf or stem cross sections of 41 species of Asparagaceae and
134 38 species of Portullugo growing at the Desert Botanical Garden, Phoenix, AZ; fixed specimens
135 were created as above and imaged on an Olympus BX51 microscope (Evident Corporation,
136 Toyko, Japan) with an Infinity3-3UR camera (Teledyne Lumenera, Ottawa, Canada). To
137 supplement our sampling, we were provided high resolution images of leaf cross sections of 13
138 *Portulaca* (Portulacaceae) species used in Ocampo *et al.* (2013) by the authors.

139 The multiple data sets had some taxonomic overlap and some included multiple
140 measurements from multiple accessions of the same species. To reduce our data set to one

141 observation per species, we took the mean of each feature where multiple accessions were
142 measured; these mean species values were used as the basis for analysis throughout. We binned
143 each taxon into three CAM phenotypes based on Gilman *et al.* (2023) and references therein: C₃,
144 C₃-C₄, and C₄ taxa were coded as “non-CAM”; taxa that fix the minority of their daily CO₂ with
145 CAM (C₃+CAM, C₃-C₄+CAM, and C₄+CAM) were coded as minority CAM (mCAM); and taxa
146 that primarily use CAM to fix CO₂ (i.e., over 50%, resulting in $\delta^{13}\text{C}$ ratios $\geq -18.7\text{\textperthousand}$; Winter and
147 Holtum, 2002) as primary CAM (pCAM). The final data set contained observations from 5,316
148 non-CAM, 207 mCAM, and 222 pCAM taxa (Supporting Information Dataset S1).

149 **Measuring plant anatomy**

150 Automated analyses of plant tissues can be difficult because many or most cells are in direct
151 contact with other cells around much of their perimeter, rather than being separated by clear
152 boundaries. We developed a lightweight image segmentation tool built in Python 3 with OpenCV
153 v4.5.2 (Bradski, 2000) called MiniContourFinder to facilitate measurement of histology slides.
154 Segmentation in MiniContourFinder is accomplished through a combination of thresholding,
155 gradient, and morphological operations (Supporting Information Figure S1). MiniContourFinder
156 was designed to allow users with minimal experience on the command line or image processing
157 to quickly generate accurate and reproducible contours, particularly from plant histology images.
158 MiniContourFinder can be run through the command line or a graphical user interface to tune
159 contours in real time. We used MiniContourFinder to measure MA in our new Asparagaceae and
160 Portulugo data sets. We used ImageJ v1.53 (Schneider *et al.*, 2012) to calculate LT (for leafy
161 species) and IAS (in roughly 300 μm x 300 μm areas of mesophyll).

162 **Statistical analysis**

163 We investigated group differences in anatomical measurements by assessing normality and
164 homoscedasticity, comparing raw and transformed data, testing for group differences, and finally
165 using post-hoc tests to identify group differences. We first assessed assumptions of normality
166 using D’Angostino and Pearson’s test (D’Angostino and Pearson, 1973) and homoscedasticity
167 using Bartlett’s test (Bartlett, 1937) of raw and \log_{10} -transformed data. None of the features were
168 normal when raw or transformed, but \log_{10} -transformation substantially decreased
169 heteroscedasticity: all transformed features were homoscedastic except SLA, which was much

170 less heteroscedastic (Supporting Information Fig. S2 and Table S3). We therefore continued with
171 Kruskal-Wallis (KW) tests for group differences (Kruskal and Wallace, 1952) with the
172 transformed data, and Dunn's post-hoc tests (Dunn, 1964) where KW tests revealed significant
173 group differences. We tested for correlations between transformed features using Pearson's r
174 (Pearson, 1895). All statistical analyses were performed using Python v3.7.12, scipy v1.5.3
175 (Virtanen *et al.*, 2020), and scikit-posthocs v0.6.4 (Terpilowski, 2019).

176 **Supervised classification**

177 We attempted to classify species' CAM phenotypes based on anatomy using the supervised
178 learning method gradient boosting implemented in XGBoost via the Python package 'xgboost'
179 v.1.5.0 (Chen and Guestrin, 2016). XGBoost implements gradient tree boosting algorithms
180 (Friedman *et al.*, 2000; Friedman, 2001) that use greedy learning over an ensemble of regression
181 trees to train classification models. XGBoost is rare in that it can accept observations with
182 missing values without the need for data imputation. We conducted multiclass classification of
183 non-CAM, mCAM, and pCAM taxa and a simpler, binary classification of non-CAM and CAM
184 taxa, where mCAM and pCAM taxa were combined. We explored a variety of alternative
185 parameterizations: changing the default booster (gbtree) to DART (Rashmi and Gilad-Bachrach,
186 2015), which can reduce overfitting by randomly dropping decision trees; changing the objective
187 function (softmax or softprob for multiclass classification; logistic probability, logistic raw score,
188 or hinge loss for binary classification); and changing the evaluation metric (multiclass logloss,
189 AUC, or multiclass error rate for multiclass classification; error rate for binary classification)
190 (the AUC evaluation metric required a softprob objective function). In all cases we randomly
191 divided our data set between training (80%) and testing (20%).

192 We also tried several strategies to reduce the effects of highly imbalanced classes and
193 sparsity. We attempted to reduce class imbalance by adjusting the parameter 'max_delta_step'
194 (MDS), by random over- or under-sampling our training data, and by merging mCAM and
195 pCAM into a binary classification model. Increasing MDS above its default (0) creates an
196 additional penalty that reduces splitting within trees, or the addition of trees entirely, in highly
197 imbalanced data sets. Random over-sampling (ROS) resamples minority classes until all class
198 labels are equal (augmenting training data), while random under-sampling (RUS) subsamples

199 classes until all class labels are equal (reducing training data). Our data were also quite sparse
200 (67% missing data) because we merged data from largely non-overlapping studies. We evaluated
201 three data imputation strategies: median (missing features were imputed with the median),
202 iterative (missing features were imputed by regression of present features), and K -nearest
203 neighbors (Knn; missing features were imputed using the nearest neighbors in a Knn
204 embedding).

205 **Phylogenetic tree inference**

206 The Portullugo, the clade inclusive of the Portulacineae (families Anacampserotaceae,
207 Basellaceae, Cactaceae, Didiereaceae, Montiaceae, Portulacaceae, and Talinaceae) and its sister
208 clade (Molluginaceae) is well-suited for large, comparative phylogenetic studies because of
209 recent sequence data, its diversity of CAM phenotypes, and the overlap between anatomical data
210 and extant phylogenies. We constructed a new phylogeny of the Portullugo by merging two
211 previously published sequence matrices that were obtained using different techniques. The first
212 dataset consisted of 841 loci from transcriptomic data used to study the evolution of
213 Portulacineae and its adaptation to harsh environments (Wang *et al.*, 2019). The second dataset
214 was a targeted enrichment of 83 gene families, primarily with roles in plant respiration and
215 photosynthesis (Goolsby *et al.*, 2018; Hancock *et al.*, 2018; Moore *et al.*, 2018). To find common
216 loci between the datasets, we independently called consensus sequences for each locus and
217 mapped them against the sugar beet genome (assembly version EL10_1.0; McGrath *et al.*, 2022)
218 using Blast v.2.13.0 (Camacho *et al.*, 2009). Mapping consensus sequences for each locus
219 proved more accurate than using random representative sequences for a given locus due to high
220 sequence variation. If consensus loci hit multiple reference scaffolds, we retained the reference
221 locus with the highest bitscore. We used the resulting mapping coordinates to search for potential
222 overlapping loci between datasets and aligned them using MAFFT v.7.508 (Katoh and Sandley,
223 2013). Loci showing considerable dataset overlap were concatenated to create an initial matrix of
224 loci represented by both datasets, and then flanked with 7 randomly selected non-overlapping
225 loci from each dataset to increase the number of taxa included and overall matrix size.

226 We concatenated all loci and constructed a maximum likelihood-based tree using IQ-
227 TREE v2.2.0.3 (Minh *et al.*, 2020). Within IQ-TREE, a model of sequence evolution was

228 selected using the automated model finder (Kalyaanamoorthy *et al.*, 2017) constrained to the
229 GTR family of models; node support was assessed using ultrafast bootstrap approximation
230 (Hoang *et al.*, 2017), and the tree space was constrained by a guide tree in which all families
231 were monophyletic. We time calibrated the tree using the fast least squares dating method (To *et*
232 *al.*, 2016) included in IQ-TREE using the entire concatenated sequence matrix the 13 secondary
233 node calibrations used by Wang *et al.* (2019) from Arakaki *et al.* (2011) (Supporting Information
234 Table S4). Confidence intervals were inferred by 100 resamplings of branch lengths by drawing
235 new clock rates (log-normal distribution with mean 1 and standard deviation 0.2), tip dates were
236 set to 0, and a GTR+F substitution model was selected with the automated model finder.

237 **Phylogenetic trait analyses**

238 We reconstructed the evolutionary history of CAM in the Portullugo using stochastic character
239 mapping (Nielsen, 2002; Huelsenbeck *et al.*, 2003) implemented with the ‘make.simmap’
240 function of the R package ‘phytools’ v1.2-0 (Revell, 2012), to model CAM evolutionary history
241 assuming 1) an all rates different (ARD) model, and 2) a constrained ARD model without
242 reversions from pCAM to mCAM; both models assumed a root state of non-CAM. The
243 constraint of the latter model was informed by the lack of evidence for reversions from pCAM
244 throughout vascular plants. In all analyses, we pruned our tree to one sample per species and
245 node reconstructions were visualized as pie charts summarizing the state frequencies over 10000
246 stochastic maps.

247 To assess the relationships between CAM phenotypes and anatomical traits in the Portullugo, we
248 used a threshold model of trait evolution (Wright, 1934; Felsenstein, 2005), implemented with
249 the ‘threshBayes’ function (Revell, 2014) of ‘phytools’ v1.2-0, and phylogenetic least squares
250 (PGLS) regression (Grafen, 1989), implemented with the R package ‘nlme’ v3.1-162 (Pinheiro
251 *et al.*, 2023). We used PGLS regression to assess relationships between continuous anatomical
252 traits and between anatomical traits and discrete CAM phenotypes (as a predictor variable). We
253 used threshold models to measure the correlations between anatomical traits and CAM
254 phenotype. In all analyses, our tree was pruned to match the taxa with anatomical data and
255 reduced to one sample per taxon where necessary.

256

257 **RESULTS**

258 Non-phylogenetic analyses of anatomy across angiosperms demonstrated significant group
259 differences for all five anatomical features investigated (Supporting Information Table S5).
260 Dunn's post-hoc tests identified significant ($p < 0.05$), and generally consistent, differences
261 between CAM phenotypes for most features: the largest differences were observed between non-
262 CAM and pCAM phenotypes, with mCAM intermediate but not always significantly different
263 from both non-CAM and pCAM (Fig. 2). Where sufficient data were available, these trends were
264 supported within individual families (Supporting Information Fig. S3). We found significant
265 negative correlations between MA and leaf dry matter content (LDMC), between LT and specific
266 leaf area (SLA), LDMC, and IAS, and between LDMC and SLA; a significant positive
267 correlation was found between LT and MA (Supporting Information Fig. S4).

268 Multiclass classification with XGBoost yielded similar results regardless of evaluation
269 metric or objective function, with booster choice being the only source of variation (Supporting
270 Information Fig S5). Because of the similarity of those results, we only continued using models
271 with the softprob objective function and multiclass error rate (merror) evaluation metric
272 (hereafter, DART and gmtree 'base models'). The two base models had similar cross-validation
273 test accuracies ($96.0 \pm 1.1\%$) (Fig. 3a), precision and recall of non-CAM, mCAM, and pCAM
274 (Fig. 3b), and feature importance rankings ($LT > MA \geq IAS \geq LDMC > SLA$) (Supporting
275 Information Fig. S6 and Table S6). No imbalance-reduction sampling, imputation method, or
276 alternative parameterization increased overall accuracy (Fig. 3a); however, random over-
277 sampling (ROS) and random under-sampling (RUS) increased recall for mCAM and pCAM taxa
278 (Fig. 3b). Between models of similar accuracy, we prioritized improving mCAM recall (also
279 known as sensitivity in binary classification; true positives / true positives + false negatives)
280 because true negative rates of mCAM are not well known in most CAM-evolving clades. While
281 decreased non-CAM classification accuracy slightly decreased overall model accuracy, ROS
282 raised recall rates of mCAM and pCAM classification to 70% and >75%, respectively. Although
283 RUS further increased mCAM and pCAM recall (Fig. 3b), the substantial difference between
284 training and testing accuracy (Fig. 3a) suggested that these models were overfit. To further
285 address class imbalance, we combined mCAM and pCAM into a single "CAM" category and
286 attempted binary classification. Binary classification models had similar test accuracies (Fig. 3c),

287 but the hinge objective function yielded slightly higher CAM precision and recall. As in
288 multiclass classification, ROS greatly increased CAM recall, but the F1-score ($2 \times \text{precision} \times$
289 $\text{recall} / (\text{precision} + \text{recall})$) remained unchanged because of an equal magnitude drop in precision
290 (Fig. 3d).

291 Our preferred multiclass and binary classifiers both used gmtree boosters and ROS, and
292 the hinge object function for binary classification (Fig. 3e-f). Mean cross-validation accuracies
293 were $95.7 \pm 0.7\%$ and $96.1 \pm 0.6\%$ for multiclass and binary models, respectively (Fig. 3a,c). Most
294 non-CAM taxa incorrectly classified by multiclass models belonged to clades with diverse CAM
295 phenotypes (e.g., Bromeliaceae and Orchidaceae subfamily Epidendroideae), and mCAM taxa
296 were roughly equally classified as non-CAM or pCAM (Fig. 3e; Supporting Information Table
297 S7). Similarly, most incorrect predictions by the binary model were non-CAM species from
298 CAM-evolving lineages classified as CAM (Fig. 3e; Supporting Information Table S8);
299 generally, these taxa have not been thoroughly assessed for mCAM, and so it is possible that
300 they may actually have a facultative or very weak CAM cycle.

301 Our time calibrated species tree was congruent with those from which data were
302 compiled. Support was generally high, although multiple nodes along the backbone were
303 unresolved and left as polytomies in downstream analyses (Fig. 4; Supporting Information Fig.
304 S7). Stochastic character map reconstructions of CAM evolution suggested that mCAM evolved
305 at the base of Portulacineae, and that multiple transitions to pCAM occurred in the Cactaceae and
306 Didiereaceae, while multiple reversions to non-CAM occurred in the Montiaceae (Fig. 4).
307 Though similar, we preferred a constrained all rates different model of CAM evolution (Fig. 4;
308 Supporting Information Fig. 8) to an unconstrained model (Supporting Information Fig. S9)
309 because there is no strong empirical evidence of reversions from pCAM in any vascular plant
310 lineage. Significant phylogenetic signal was detected in all three traits measured across the
311 Portullugo (Supporting Information Table S9). Phylogenetic least squares (PGLS) regression
312 revealed multiple significant ($p < 0.05$) relationships among anatomical traits and between
313 anatomical traits and CAM phenotype (Fig. 5, Supporting Information Table S10). However,
314 AIC-based model selection favored a model between MA and IAS with a non-significant slope,
315 contrary to our expectation that greater mesophyll cell size would lead to lower IAS (Fig. 5a).
316 Greater MA was a significant ($p < 0.0001$) predictor of greater LT (Fig. 5b), and we found no

317 relationship between IAS and LT (Fig. 5c). CAM phenotype was a significant predictor of MA,
318 LT, and IAS (Fig. 5d-f). We next used phylogenetic threshold analyses to estimate the
319 correlations between CAM phenotype and anatomical traits under the hypothesis that there may
320 be anatomical boundaries between CAM phenotypes. Threshold analyses mostly supported
321 PGLS results, and recovered significant positive correlations between CAM phenotype and both
322 MA and LT (Fig. 6a-b). However, the posterior distribution of correlation coefficients between
323 CAM phenotype and IAS narrowly included 0 (Fig. 6c).

324

325 **DISCUSSION**

326 From the beaks of Galapagos finches (Darwin, 1839) to unique inflorescence architectures (Waal
327 *et al.*, 2012), the links between form and function have always inspired biologists. Fixed in place,
328 with passive mechanisms for carbon and water acquisition, plants rely on anatomical innovations
329 to adapt to different environments. Succulence has long been understood as a drought avoidance
330 adaptation, but its relationship with CAM has not been resolved as causal or merely coincidental.
331 Through our broad survey of angiosperms and detailed study of the Portulugo, we found support
332 for previous hypotheses of CAM and photosynthetic tissue anatomy co-evolution. Furthermore,
333 we demonstrate that the presence or absence of CAM may be predicted using only a handful of
334 anatomical measurements.

335 Anatomical measurements from over 200 angiosperm families revealed significant
336 differences in photosynthetic tissue anatomy of non-CAM, mCAM, and pCAM species. The
337 larger mesophyll cell area (MA) of both mCAM and pCAM species suggests some anatomical
338 specialization is required to perform CAM in any capacity, and the reduction in intercellular
339 airspace (IAS) of pCAM species indicates that further specialization is required to use CAM for
340 primary carbon metabolism. We also showed significant increases in leaf thickness (LT) and
341 decreases in leaf dry matter content (LDMC) from non-CAM to mCAM to pCAM, as well as
342 significantly lower specific leaf area (SLA) in pCAM species, which support past anatomical
343 studies that found thicker and more succulent leaves to be positively associated with strong CAM
344 activity within individual clades (Teeri *et al.*, 1981; Winter *et al.*, 1983; Nelson *et al.*, 2005,
345 2008; Zambrano *et al.*, 2014; Luján *et al.*, 2022).

346 Because lineage-specific organismal detail will surely influence physiology-anatomy
347 relationships, analyses of anatomy and CAM evolution are best evaluated using phylogenetic
348 comparative methods. PGLS regression and phylogenetic threshold analysis supported the
349 correlated evolution of larger mesophyll cells and thicker leaves. Although PGLS regression
350 further showed a continuous decrease in IAS from non-CAM to pCAM species, we found no
351 significant relationship between IAS and MA. That IAS and MA may evolve independently of
352 one another provides an important nuance to the co-evolution of succulence and CAM.
353 Decreased IAS in CAM species has often been discussed as an adaptation to reduce CO₂ efflux
354 during malate decarboxylation (Nelson *et al.*, 2008) or as a consequence of increased succulence
355 restricting g_m , which would limit CO₂ fixation by Rubisco during the day (Zambrano *et al.*, 2014;
356 Earles *et al.*, 2018; Edwards, 2019). More recently, reduced IAS has been hypothesized to be an
357 indirect consequence of increased mesophyll cell volume used for malic acid storage (Leverett *et*
358 *al.*, 2023). While we found that succulence generally increased with CAM evolution, the
359 decoupling of the underlying traits may allow the evolution of intermediate photosynthetic and
360 anatomical phenotypes that efficiently utilize both CAM and C₃ or C₄ photosynthesis. These
361 conclusions are consistent with photosynthetic models that found increased vacuolar volume
362 (and therefore MA) necessary for CAM (Töpfer *et al.*, 2020) and empirical findings that the high
363 IAS in mCAM species may allow for C₃ (or C₄) photosynthesis when plants are not engaging
364 CAM (Nelson *et al.*, 2008; Zambrano *et al.*, 2014). Furthermore, lowest IAS values in the
365 Portullugo were observed in pCAM species, which reinforces the hypothesis that extremely low
366 IAS may reduce g_m and C₃ or C₄ efficiency.

367 In addition to providing support for a positive relationship between CAM and succulence,
368 our findings point towards interactions between life history, CAM, and succulence for those taxa
369 that do not neatly fall along regression lines. Phylogenetic analyses of the Portullugo showed
370 general increases in succulence and a tightening of the distributions of underlying traits for
371 pCAM species. In contrast, mCAM taxa had both the single largest MA and greatest IAS
372 observations, with values that mostly spanned the non-CAM to pCAM range. The eight largest
373 observed MA values in the Portullugo were from mCAM species; most are annual species, with
374 the exceptions of *Parakeelya flava* (a perennial geophyte with above ground tissue that regrows
375 annually) and *Grahamia bracteata* and *Talinopsis frutescens* (which have non-succulent woody
376 stems and drought-deciduous leaves). This suggests that the evolution of pCAM requires a shift

377 to a (functional) perennial life history with long-lived photosynthetic tissues (Hancock *et al.*,
378 2019); indeed, we are unaware of any annual pCAM species. The halophyte *Halophytum*
379 *ameghinoi* had the second largest observed MA in the Portullugo. While saline soils may select
380 for increased succulence to maintain cytosolic ion balance (Naidoo and Rughunanan, 1990;
381 Ogburn and Edwards, 2010), high salt concentrations inhibit the central CAM enzymes
382 phosphoenolpyruvate carboxylase (PEPC) and malic enzyme (ME) (Kluge and Ting, 1978), and
383 may therefore represent an ecological constraint on the evolution of pCAM.

384 Our ancestral state reconstruction of CAM in the Portullugo was the first to model CAM
385 as an ordered multistate trait, and supported an early- to mid-Eocene origin of mCAM – a time
386 when the Earth’s atmosphere had relatively high levels of CO₂ (Rae *et al.*, 2021). The
387 reconstruction of mCAM at the crown of the Portulacineae agrees with transcriptomic data that
388 suggest a single recruitment event of a PEPC ortholog for use in CAM (Christin *et al.*, 2014;
389 Goolsby *et al.*, 2018). All transitions to pCAM were found to be within the past 30 Ma (Sage *et al.*,
390 2023), congruent with shifts across angiosperms (including within Caryophyllales) to C₄
391 photosynthesis, as atmospheric CO₂ fell below 500ppm (Christin *et al.*, 2011). Despite declining
392 CO₂ in the Oligocene and Miocene, multiple lineages within the Montiaceae have lost the ability
393 to perform CAM. Although we expect more Montiaceae lineages to exhibit CAM upon
394 experimentation, multiple independent losses of CAM have been experimentally validated in
395 *Parakeelya* (Hancock *et al.*, 2019), a clade endemic to hot, dry areas of Australia. While life
396 history may constrain the evolution of pCAM, it remains unclear why some members of the
397 Portulacineae – occupying similarly semiarid environments – transitioned to C₃ photosynthesis,
398 while others simultaneously transitioned to pCAM, as CO₂ continued to decline. We suspect that
399 these losses of CAM may be linked to shifts in phenology; for example, C₃ *Parakeelya* tend to
400 germinate toward the end of the wet season, when temperatures are cooler and water is readily
401 available.

402 Most clades with CAM lineages show highly bimodal distributions of carbon isotope
403 ratios (Messerschmid *et al.*, 2021) that have been used for decades to identify pCAM species, but
404 are generally unable to distinguish between mCAM from non-CAM. Laborious controlled
405 experiments (e.g., of gas exchange or malic acid content) with live plants have been the only
406 ways to identify mCAM, but such experiments are not feasible for many long-lived, rare, or

407 difficult to cultivate species. We found that differences in photosynthetic anatomy across
408 angiosperms translated into moderate to high accuracy in predicting CAM phenotype. After
409 assessing a variety of machine learning models, we found that random-over sampling (ROS)
410 increased prediction of mCAM and pCAM species while not overfitting to training data. To our
411 knowledge, machine learning has not yet been applied to predict the presence or absence of
412 physiological traits from anatomical measurements, such as CAM phenotypes. We believe that
413 the accuracy we obtained represents a lower bound on the true accuracy of our models because
414 some of our misclassified non-CAM species have not been thoroughly investigated for mCAM.
415 For example, multiple orchid and bromeliad species labeled as non-CAM were predicted to be
416 mCAM, but have not been subjected to drought experiments that might induce CAM activity.
417 We predict that some misclassified species, such as *Nolina bigelovii* (Asparagaceae) will exhibit
418 CAM upon experimentation, resulting in more true positive predictions. Experimentation should
419 continue to be the gold-standard for determining CAM phenotypes, but machine learning
420 models, such as those developed here, could play a valuable role in prioritizing study species and
421 would only require small tissue sections for initial fixation and measurement.

422 Applications of machine learning in Plant Physiology and evolution are only just
423 beginning. Machine learning has been successful in predicting real-time photosynthetic status;
424 for example, deep learning using hyperspectral reflectance in wheat has been used to predict
425 electron transport rate, CO₂ assimilate rate, stomatal conductance, and more (Furbank *et al.*,
426 2021). Our machine learning models were limited in several ways; perhaps most by the degree of
427 missing data and class imbalance. Our greatest model improvements came when using ROS,
428 suggesting that measuring new mCAM and pCAM species to reduce class imbalance will
429 increase model accuracy. If missing data could be sufficiently reduced, imputation strategies may
430 facilitate the use of models beyond XGBoost, which allows missing data. In addition to our
431 machine learning models, we hope that the tools and methodology developed here for measuring
432 anatomy and merging sequence matrices will facilitate future studies of anatomical evolution.
433 Although software exists for taking measurements from images (e.g., ImageJ (Schneider *et al.*,
434 2012), which we used for portions of this study), making dozens or hundreds of measurements
435 needed for phylogenetic studies remains time consuming and the results are not easily
436 reproducible. Our image segmentation software, MiniContourFinder can be automated from the
437 command line, quickly segment and measure image features, and record associated metadata so

438 exact measurements can be reproduced. Finally, our strategy for combining reduced-genomic
439 sampling data types into a single phylogenetic analysis is flexible and in theory adaptable to any
440 sequencing strategy. Most clades have reference, or near-reference, quality genomes within ~75
441 Ma of their focal taxa (as in this study) (Cheng *et al.*, 2018) that can serve as common maps to
442 identify overlapping genomic regions, and high-quality transcriptomes (Matasci *et al.*, 2014;
443 Leebens-Mack *et al.*, 2019) or targeted sequencing data (Johnson *et al.*, 2019) for constructing
444 backbones in larger phylogenies.

445 In conclusion, with a broad sampling of anatomical traits from thousands of angiosperms
446 and a detailed phylogenetic study of the Portullugo clade, we provided support for hypotheses of
447 CAM anatomical evolution. Our findings suggest that even weakly expressed CAM is correlated
448 with larger mesophyll cells, and that decreased intercellular airspace in photosynthetic tissue is
449 associated with a transition to using CAM as the primary carbon fixation pathway. Furthermore,
450 our findings point towards possible evolutionary constraints on pCAM evolution, such as annual
451 life history and halophytism. We were able predict CAM phenotypes from a handful of
452 anatomical features, which represents a successful first application of machine learning to this
453 problem, but also highlights the paucity of anatomical data for species capable of weak or
454 facultative CAM. As data accumulate, we hope that these correlations will be continuously
455 evaluated across vascular plants with tools that may allow causal evolutionary inference, such as
456 phylogenetic path analysis (von Hardenberg and Gonzalez-Voyer, 2013). We expect that efforts
457 to quantify key anatomical parameters for a diversity of CAM phenotypes will more sharply
458 delineate the anatomical requirements of even a weak CAM cycle, and demonstrate the
459 anatomical and biochemical interplay during the evolutionary transition to a primary CAM
460 physiology.

461

462 ACKNOWLEDGMENTS

463 We would like to thank Dr. Jonathan Koss for advice on implementing computer vision
464 algorithms, Dr. Elena Voznesenskaya for providing images of *Portulaca* species, Dr. Eric Lazo-
465 Wasem and Lourdes Rojas for help imaging specimens, and Joni Ward, Rual Puente-Martinez,
466 and the rest of the staff at the Desert Botanical Garden (Phoenix, AZ) for assistance with the

467 living collections that contributed to this manuscript. This research was funded by the National
468 Science Foundation (IOS-1754662 to EJE).

469

470 AUTHOR CONTRIBUTIONS

471 ISG, KH, and EJE designed the research plan; ISG, KH, and LPH collected, fixed, and imaged
472 specimens; ISG developed image segmentation software, curated morphological data, time
473 calibrated the phylogeny, conducted statistical analyses, and wrote the first draft of the
474 manuscript; CAM-L designed the sequence merging strategy, generated sequence matrices, and
475 constructed the initial phylogeny; all authors contributed to editing and revising the manuscript.

476

477 DATA AVAILABILITY STATEMENT

478 The data that supports the findings of this study are available in the Supporting Information of
479 this article. All statistical analyses of this study can be found at
480 <https://github.com/isgilman/Predicting-CAM> and all installation and documentation for
481 MiniContourFinder can be found at <https://minicontourfinder.readthedocs.io/en/latest/>.

482 Additional Supporting Information may be found online in the supporting information
483 section at the end of the article.

484

485 COMPETING INTERESTS

486 The authors declare no competing interests.

487

488 REFERENCES

489 **Arakaki M, Christin P-A, Nyffeler R, Lendel A, Eggli U, Ogburn RM, Spriggs E, Moore**
490 **MJ, Edwards EJ.** 2011. Contemporaneous and recent radiations of the world's major

491 succulent plant lineages. *Proceedings of the National Academy of Sciences USA* **108**:
492 8379–8384.

493 **Bartlett MS. 1937.** Properties of sufficiency and statistical tests. *Proceedings of the Royal
494 Society of London. Series A: Mathematical and Physical Sciences* **160**: 268–282.

495 **Bradski G. 2000.** The OpenCV Library. *Dr. Dobb's J. Software Tools*.

496 **Camacho C, Coulouris G, Avagyan V, Ma N, Papadopoulos J, Bealer K, Madden TL. 2009.**
497 BLAST+: architecture and applications. *BMC Bioinformatics* **10**: 421.

498 **Chen T, Guestrin C. 2016.** XGBoost: A scalable tree boosting system. In *Proceedings of the
499 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data
500 Mining* 785–794.

501 **Cheng S, Melkonian M, Smith SA, Brockington S, Archibald JM, Delaux P-M, Li F-
502 W, Melkonian B, Mavrodiev EV, Sun W, et al. 2018.** 10KP: a phylogenetic genome
503 sequencing plan. *GigaScience* **7**: 1–9.

504 **Christin P-A, Osborne CP, Sage RF, Arakaki M, Edwards EJ. 2011.** C₄ eudicots are not
505 younger than C₄ monocots. *Journal of Experimental Botany* **62**: 3171–3181.

506 **Christin P-A, Arakaki M, Osborne CP, Bräutigam A, Sage RF, Hibberd JM, Kelly
507 S, Covshoff S, Wong GK-S, Hancock L, et al. 2014.** Shared origins of a key enzyme
508 during the evolution of C₄ and CAM metabolism. *Journal of Experimental Botany* **65**:
509 3609–3621.

510 **Cousins AB, Mullendore DL, Sonawane BV. 2020.** Recent developments in mesophyll
511 conductance in C₃, C₄, and Crassulacean acid metabolism plants. *The Plant Journal* **101**:
512 816–830.

513 **Coutinho LM. 1969.** Novas observações sobre a ocorrência do “efeito de de Sassure” e suas
514 relações com a suculência, a temperatura folhear e os movimentos estomáticos. *Boletim
515 Da Faculdade De Filosofia Ciências E Letras Universidade De São Paulo Botânica* **24**:
516 77–102.

517 **D'Agostino R, Pearson ES. 1973.** Tests for departure from normality. Empirical results for the
518 distributions of b_2 and $\sqrt{b_1}$. *Biometrika* **60**: 613–622.

519 **Darwin C. 1839.** *The Voyage of the Beagle*. D. New York: Appleton & Co.

520 **de Saussure T. 1804.** *Recherches chimiques sur la végétation*.

521 **Dunn OJ. 1964.** Multiple comparisons using rank sums. *Technometrics* **6**: 241–252.

522 **Earles JM, Theroux-Rancourt G, Roddy AB, Gilbert ME, McElrone AJ, Brodersen
523 CR. 2018.** Beyond porosity: 3D leaf intercellular airspace traits that impact mesophyll
524 conductance. *Plant Physiology* **178**: 148–162.

525 **Earnshaw MJ, Winter K, Ziegler H, Stichler W, Cruttwell NEG, Kerenga K, Cribb
526 PJ, Wood J, Croft JR, Carver KA, et al. 1987.** Altitudinal changes in the incidence of
527 Crassulacean acid metabolism in vascular epiphytes and related life forms in Papua New
528 Guinea. *Oecologia* **73**: 566–572.

529 **Edwards EJ. 2019.** Evolutionary trajectories, accessibility and other metaphors: the case of C₄
530 and CAM photosynthesis. *New Phytologist* **223**: 1742–1755.

531 **Felsenstein J. 2005.** Using the quantitative genetic threshold model for inferences between and
532 within species. *Philosophical Transactions of the Royal Society B: Biological Sciences*
533 **360**: 1427–1434.

534 **Fraser LH. 2020.** TRY—A plant trait database of databases. *Global Change Biol.* **26**, 189–190.

535 **Friedman JH. 2001.** Greedy function approximation: a gradient boosting machine. *The Annals
536 of Statistics* **29**: 1189–1232.

537 **Friedman J, Hastie T, Tibshirani R. 2000.** Additive logistic regression: a statistical view of
538 boosting (with discussion and a rejoinder by the authors). *The Annals of Statistics* **28**:
539 337–407.

540 **Furbank RT, Silva-Perez V, Evans JR, Condon AG, Estavillo GM, He W, Newman S, Poiré
541 R, Hall A, He Z. 2021.** Wheat physiology predictor: predicting physiological traits in

542 wheat from hyperspectral reflectance measurements using deep learning. *Plant Methods*
543 17: 108.

544 **Gilman IS, Smith JAC, Holtum JAM, Sage RF, Silvera K, Winter K, Edwards EJ. 2023.**
545 The CAM lineages of plant Earth. *Annals of Botany* (accepted).

546 **Goolsby EW, Moore AJ, Hancock LP, Vos JM de, Edwards EJ. 2018.** Molecular evolution
547 of key metabolic genes during transitions to C₄ and CAM photosynthesis. *American*
548 *Journal of Botany* 105: 602–613.

549 **Grafen A. 1989.** The phylogenetic regression. *Philosophical Transactions of the Royal Society*
550 *B: Biological Sciences* 326: 119–157.

551 **Hancock L, Edwards EJ. 2014.** Phylogeny and the inference of evolutionary trajectories.
552 *Journal of Experimental Botany* 65: 3491–3498.

553 **Hancock LP, Obbens F, Moore AJ, Thiele K, Vos JM de, West J, Holtum JAM, Edwards**
554 **EJ. 2018.** Phylogeny, evolution, and biogeographic history of *Calandrinia* (Montiaceae).
555 *American Journal of Botany* 105: 1021–1034.

556 **Hancock LP, Holtum JAM, Edwards EJ. 2019.** The evolution of CAM photosynthesis in
557 Australian *Calandrinia* reveals lability in C₃+CAM phenotypes and a possible constraint
558 to the evolution of strong CAM. *Integrative and Comparative Biology* 59: 517–534.

559 **Herrera A. 2020.** Are thick leaves, large mesophyll cells and small intercellular air spaces
560 requisites for CAM? *Annals of Botany* 125: 859–868.

561 **Heyduk K, McKain MR, Lalani F, Leebens-Mack J. 2016.** Evolution of a CAM anatomy
562 predates the origins of Crassulacean acid metabolism in the Agavoideae (Asparagaceae).
563 *Molecular Phylogenetics and Evolution* 105: 102–113.

564 **Heyduk K, Ray JN, Leebens-Mack J. 2020.** Leaf anatomy is not correlated to CAM function
565 in a C₃+CAM hybrid species, *Yucca gloriosa*. *Annals of Botany* 127: 437–449.

566 **Hoang DT, Chernomor O, Haeseler A von, Minh BQ, Vinh LS.** 2017. UFBoot2: improving
567 the ultrafast bootstrap approximation. *Molecular Biology and Evolution* **35**: 518–522.

568 **Huelsenbeck JP, Nielsen R, Bollback JP.** 2003. Stochastic mapping of morphological
569 characters. *Systematic Biology* **52**: 131–158.

570 **Johnson MG, Pokorny L, Dodsworth S, Botigue LR, Cowan RS, Devault A, Eiserhardt
571 WL, Epitalwage N, Forest F, Kim JT, et al.** 2019. A universal probe set for targeted
572 sequencing of 353 nuclear genes from any flowering plant designed using k-medoids
573 clustering. *Systematic Biology* **68**: 594–606.

574 **Kalyaanamoorthy S, Minh BQ, Wong TKF, Haeseler A von, Jermiin
575 LS.** 2017. ModelFinder: fast model selection for accurate phylogenetic estimates. *Nature
576 Methods* **14**: 587–589.

577 **Katoh K, Standley DM.** 2013. MAFFT multiple sequence alignment software version 7:
578 improvements in performance and usability. *Molecular Biology and Evolution* **30**: 772–
579 780.

580 **Kluge M, Ting IP.** 1978. *Crassulacean Acid Metabolism, Analysis of an Ecological Adaptation.*
581 Springer-Verlag.

582 **Kruskal WH, Wallis WA.** 1952. Use of ranks in one-criterion variance analysis. *Journal of the
583 American Statistical Association* **47**: 583–621.

584 **Leebens-Mack JH, Barker MS, Carpenter EJ, Deyholos MK, Gitzendanner MA, Graham
585 SW, Grosse I, Li Z, Melkonian M, Mirarab S, et al.** 2019. One thousand plant
586 transcriptomes and the phylogenomics of green plants. *Nature* **574**: 679–685.

587 **Leverett A, Borland AM, Inge EJ, Hartzell S.** 2023. Low internal air space in plants with
588 crassulacean acid metabolism may be an anatomical spandrel. *Annals of Botany*: mcad109.

589 **Luján M, Oleas NH, Winter K.** 2022. Evolutionary history of CAM photosynthesis in
590 Neotropical *Clusia*: insights from genomics, anatomy, physiology and climate. *Botanical
591 Journal of the Linnean Society* **199**: 538–556.

592 **Males J.** 2017. Secrets of succulence. *Journal of Experimental Botany* **68**: 2121–2134.

593 **Males J.** 2018. Concerted anatomical change associated with Crassulacean acid metabolism in
594 the Bromeliaceae. *Functional Plant Biology* **45**: 681–695 (2018).

595 **Matasci N, Hung L-H, Yan Z, Carpenter EJ, Wickett NJ, Mirarab S, Nguyen N, Warnow**
596 **T, Ayyampalayam S, Barker M, et al.** 2014. Data access for the 1,000 Plants (1KP)
597 project. *GigaScience* **3**: 17.

598 **McGrath JM, Funk A, Galewski P, Ou S, Townsend B, Davenport K, Daligault H, Johnson**
599 **S, Lee J, Hastie A, et al.** 2022. A contiguous *de novo* genome assembly of sugar beet
600 EL10 (*Beta vulgaris* L.). *DNA Research* dsac033.

601 **Messerschmid TFE, Wehling J, Bobon N, Kahmen A, Klak C, Los JA, Nelson DB, Santos P**
602 **dos, Vos JM de, Kadereit G.** 2021. Carbon isotope composition of plant photosynthetic
603 tissues reflects a Crassulacean acid metabolism (CAM) continuum in the majority of
604 CAM lineages. *Perspectives in Plant Ecology, Evolution and Systematics* **51**: 125619.

605 **Minh BQ, Schmidt HA, Chernomor O, Schrempf D, Woodhams MD, Haeseler A**
606 **von, Lanfear R.** 2020. IQ-TREE 2: new models and efficient methods for phylogenetic
607 inference in the genomic era. *Molecular Biology and Evolution* **37**: 1530–1534.

608 **Moore AJ, Vos JMD, Hancock LP, Goolsby E, Edwards EJ.** 2018. Targeted enrichment of
609 large gene families for phylogenetic inference: phylogeny and molecular evolution of
610 photosynthesis genes in the Portulugo Clade (Caryophyllales). *Systematic Biology* **67**:
611 367–383.

612 **Naidoo G, Rughunanan R.** 1990. Salt tolerance in the succulent, coastal halophyte,
613 *Sarcocornia natalensis*. *Journal of Experimental Botany* **41**: 497–502.

614 **Nelson EA, Sage TL, Sage RF.** 2005. Functional leaf anatomy of plants with Crassulacean acid
615 metabolism. *Functional Plant Biology* **32**: 409–419.

616 **Nelson EA, Sage RF.** 2008. Functional constraints of CAM leaf anatomy: tight cell packing is
617 associated with increased CAM function across a gradient of CAM expression. *Journal*
618 *of Experimental Botany* **59**: 1841–1850.

619 **Nielsen R.** 2002. Mapping mutations on phylogenies. *Systematic Biology* **51**: 729–739.

620 **Ocampo G, Koteyeva NK, Voznesenskaya EV, Edwards GE, Sage TL, Sage RF, Columbus**
621 **JT.** 2013. Evolution of leaf anatomy and photosynthetic pathways in Portulacaceae.
622 *American Journal of Botany* **100**: 2388–2402.

623 **Ogburn RM, Edwards EJ.** 2010. *Advances in Botanical Research*, Vol. 55 Ch. 4. Burlington:
624 Academic Press.

625 **Ogburn RM, Edwards EJ.** 2012. Quantifying succulence: a rapid, physiologically meaningful
626 metric of plant water storage. *Plant, Cell and Environment* **35**: 1533–1542.

627 **Ogburn RM, Edwards EJ.** 2013. Repeated origin of three-dimensional leaf venation releases
628 constraints on the evolution of succulence in plants. *Current Biology* **23**: 722–726.

629 **Osmond CB.** 1978. Crassulacean acid metabolism: a curiosity in context. *Annual Review of*
630 *Plant Physiology* **29**: 379–414.

631 **Pearson K.** 1895. Note on regression and inheritance in the case of two parents. *Proceedings of*
632 *the Royal Society of London* **58**: 240–242.

633 **Pinheiro JC, Bates D, R Core Team.** 2023. *nlme: Linear and Nonlinear Mixed Effects Models*.
634 R package version 3.1-163, <https://CRAN.R-project.org/package=nlme>.

635 **Rae JWB, Zhang YG, Liu X, Foster GL, Stoll HM, Whiteford RDM.** 2021. Atmospheric
636 CO₂ over the past 66 million years from marine archives. *Annual Review of Earth and*
637 *Planetary Sciences* **49**: 609–641.

638 **Rashmi KV, Gilad-Bachrach R.** 2015. DART: Dropouts meet multiple additive regression
639 trees. In: Lebanon G, Vishwanathan SVN, eds. *Proceedings of the Eighteenth*
640 *International Conference on Artificial Intelligence and Statistics*.

641 **Revell LJ.** 2012. phytools: an R package for phylogenetic comparative biology (and other
642 things). *Methods in Ecology and Evolution* **3**: 217–223.

643 **Revell LJ.** 2014. Ancestral character estimation under the threshold model from quantitative
644 genetics. *Evolution* **68**: 743–759.

645 **Sage RF, Gilman IS, Smith JAC, Silvera K, Edwards EJ.** 2023. Atmospheric CO₂ decline and
646 the timing of CAM plant evolution, *Annals of Botany*: mcad122.

647 **Schneider CA, Rasband WS, Eliceiri KW.** 2012. NIH Image to ImageJ: 25 years of image
648 analysis. *Nature Methods* **9**: 671–675.

649 **Silvera K, Santiago LS, Winter K.** 2005. Distribution of Crassulacean acid metabolism in
650 orchids of Panama: evidence of selection for weak and strong modes. *Functional Plant
651 Biology* **32**: 397–11.

652 **Tavşanoğlu C, Pausas JG.** 2018. A functional trait database for Mediterranean Basin plants.
653 *Scientific Data* **5**: 180135.

654 **Teeri JA, Tonsor SJ, Turner M.** 1981. Leaf thickness and carbon isotope composition in the
655 Crassulaceae. *Oecologia* **50**: 367–369.

656 **Terpilowski M.** 2019. scikit-posthocs: pairwise multiple comparison tests in Python. *The
657 Journal of Open Source Software* **4**: 1169.

658 **To T-H, Jung M, Lycett S, Gascuel O.** 2016. Fast dating using least-squares criteria and
659 algorithms. *Systematic Biology* **65**: 82–97.

660 **Töpfer N, Braam T, Shameer S, Ratcliffe RG, Sweetlove LJ.** 2020. Alternative CAM modes
661 provide environment-specific water-saving benefits in a leaf metabolic model. *Plant Cell*
662 **32**: 3689–3705.

663 **Virtanen P, Gommers R, Oliphant TE, Haberland M, Reddy T, Cournapeau D, Burovski
664 E, Peterson P, Weckesser W, Bright J, et al.** 2020. SciPy 1.0: fundamental algorithms
665 for scientific computing in Python. *Nature Methods* **17**: 261–272.

666 **von Hardenberg A, Gonzalez-Voyer A. 2013.** Disentangling evolutionary cause-effect
667 relationships with phylogenetic confirmatory path analysis. *Evolution* **67**, 378–387.

668 **Waal C, Barrett SCH, Anderson B. 2012.** The effect of mammalian herbivory on inflorescence
669 architecture in ornithophilous *Babiana* (Iridaceae): implications for the evolution of a
670 bird perch. *American Journal of Botany* **99**: 1096–1103.

671 **Wang N, Yang Y, Moore MJ, Brockington SF, Walker JF, Brown JW, Liang B, Feng**
672 **T, Edwards C, Mikenas J, et al.** 2019. Evolution of Portulacineae marked by gene tree
673 conflict and gene family expansion associated with adaptation to harsh environments.
674 *Molecular Biology and Evolution* **36**: 112–126.

675 **Winter K. 2019.** Ecophysiology of constitutive and facultative CAM photosynthesis. *Journal of*
676 *Experimental Botany* **2**: 16178–14.

677 **Winter K, Holtum JAM. 2002.** How closely do the $\delta^{13}\text{C}$ values of Crassulacean acid
678 metabolism plants reflect the proportion of CO_2 fixed during day and night? *Plant*
679 *Physiology* **129**: 1843–1851.

680 **Winter K, Wallace BJ, Stocker GC, Roksandic Z. 1983.** Crassulacean acid metabolism in
681 Australian vascular epiphytes and some related species. *Oecologia* **57**: 129–141.

682 **Winter K, Aranda J, Holtum JAM. 2005.** Carbon isotope composition and water-use
683 efficiency in plants with Crassulacean acid metabolism. *Functional Plant Biology* **32**:
684 381–388.

685 **Wright S. 1934.** An analysis of variability in number of digits in an inbred strain of guinea pigs.
686 *Genetics* **19**: 506–536.

687 **Yamaga-Hatakeyama Y, Okutani M, Hatakeyama Y, Yabiku T, Yukawa T, Ueno O. 2022.**
688 Photosynthesis and leaf structure of F_1 hybrids between *Cymbidium ensifolium* (C_3) and
689 *C. bicolor* subsp. *pubescens* (CAM). *Annals of Botany* **XX**: mcac157

690 **Yang X, Cushman JC, Borland AM, Edwards EJ, Wullschleger SD, Tuskan GA, Owen**
691 **NA, Griffiths H, Smith JAC, Paoli HCD, et al. 2015.** A roadmap for research on

692 Crassulacean acid metabolism (CAM) to enhance sustainable food and bioenergy
693 production in a hotter, drier world. *New Phytologist* **207**: 491–504.

694 **Zambrano VAB, Lawson T, Olmos E, Fernández-García N, Borland AM. 2014.** Leaf
695 anatomical traits which accommodate the facultative engagement of Crassulacean acid
696 metabolism in tropical trees of the genus *Clusia*. *Journal of Experimental Botany* **65**:
697 3513–3523.

698

699 **SUPPORTING INFORMATION**

700 **Figure S1** Image processing by MiniContourFinder

701 **Figure S2** Comparisons of raw and \log_{10} -transformed data.

702 **Figure S3** Results of Dunn's post-hoc tests for group differences been \log_{10} -transformed features
703 for select families.

704 **Figure S4** Correlations between \log_{10} -transformed features.

705 **Figure S5** Accuracies of base models.

706 **Figure S6** Relative feature importance scores.

707 **Figure S7** Time calibrated phylogeny of the Portullugo.

708 **Figure S8** Portullugo CAM constrained ARD reconstruction.

709 **Figure S9** Portullugo CAM ARD reconstruction.

710 **Table S1** Final anatomical data set information.

711 **Table S2** List of accessions sampled from for this study.

712 **Table S3** Results of D'Angostino and Pearson's test for normality and Bartlett's test for
713 homoscedasticity of raw and \log_{10} -transformed data.

714 **Table S4** Node calibrations used from Arakaki *et al.* (2011).

715 **Table S5** Results of Kruskal–Wallis tests for group difference between CAM phenotypes.

716 **Table S6** Relative feature importance of multiclass models.

717 **Table S7** Incorrect predictions of the best performing multiclass model.

718 **Table S8** Incorrect predictions of the best performing binary model.

719 **Table S9** Phylogenetic signal in anatomical features of the Portullugo.

720 **Table S10** Results of phylogenetic least squares (PGLS) regressions.

721 **Methods S1** Summary of MiniContourFinder image segmentation algorithm.

722 **Dataset S1** Species' anatomical data mean values.

723

724 **FIGURES**

725 **Figure 1.** Gross morphology (a-c) and photosynthetic anatomy (d-f) of species with varying
726 CAM phenotypes sampled for this study: (a,d) non-CAM *Claytonia lanceolata* Pursh
727 (Montiaceae) (b,e) minority CAM *Calyptridium umbellatum* (Torr.) Greene (Montiaceae), (c,f)
728 primary CAM *Ariocarpus retusus* Scheidw. (Cactaceae). Non-author photograph credits: (a) Dr.
729 Thomas Stoughton, (b) Anri Chomentowska, and (c) Desert Botanic Garden, Phoenix, AZ.

730 **Figure 2.** Results of Dunn's post-hoc tests for group differences been \log_{10} -transformed features.
731 Purple, yellow, and green box-and-whisker plots show non-CAM, minority CAM, and primary
732 CAM trait distributions; boxes represent the interquartile range (IQR) with a line representing
733 the median, whiskers show 1.5x the IQR, and points outside were considered outliers. MA,
734 mesophyll cell area; LT, leaf thickness; IAS, intercellular airspace; LDMC, leaf dry matter
735 content; SLA, specific leaf area.

736

737 **Figure 3.** Machine learning model accuracies. Classification error (a,c), precision and recall rates
738 (b,d), and best performing model confusion matrices (e-f) for multiclass and binary classifiers.
739 Multiclass models (a-b) varied in booster (DART or gmtree), sampling strategy (ROS or RUS),
740 imputation method (iterative, Knn, or median), and MDS (1, 2, 5, or 10); binary models (c-d)
741 varied in objective function (logistic, logitraw, or hinge) and sampling strategy (with or without
742 ROS). The columns of each confusion matrix (e-f) show the number of true CAM phenotypes in
743 the test data set and the rows show the model predictions. The diagonal in each matrix represents
744 correct model predictions and off-diagonal elements show incorrect predictions; for example, a
745 true pCAM species predicted to be non-CAM would be shown in the first row, third column of
746 (c). Knn, *K*-nearest neighbors; MDS, max_delta_step; mCAM, minority CAM; pCAM, primary
747 CAM. Base models are in bolded text and the best performing models are highlighted in red.
748

749 **Figure 4.** Time calibrated phylogeny of the Portullugo with inferred transitions between CAM
750 phenotypes. The Portullugo and Portulacineae nodes are highlighted, and color gradients indicate
751 transitions between non-CAM (purple), mCAM (yellow), and pCAM (green) based on the results
752 of our biologically-informed ancestral state reconstruction. Pie charts at nodes bracketing
753 inferred transitions show the fractions of stochastic maps supporting each ancestral state. This
754 tree has been pruned to show only those taxa with morphological data used in this study and
755 therefore not all transitions are shown; the full tree is shown in the inset and multiple ancestral
756 state reconstructions are available in the Supporting Information.
757

758 **Figure 5.** Results of phylogenetic least squares (PGLS) regression. Predictor and response
759 variables are shown on the horizontal and vertical axes, respectively. Points show trait values for
760 non-CAM (purple), mCAM (yellow), and pCAM (green) species. Solid and dashed grey lines
761 show the fitted regression lines using Brownian motion (BM) and Ornstein-Uhlenbeck (OU)
762 models of trait evolution, respectively. The best fit significant relationships are shown with bold
763 black lines, associated model coefficients, and grey shading to show standard error. MA,
764 mesophyll cell area; LT, leaf thickness; IAS, intercellular airspace.
765

766 **Figure 6.** Phylogenetic threshold model correlations. The distribution of correlation coefficients
767 (r) between CAM phenotype and \log_{10} -transformed mesophyll cell area (MA) (A), intercellular
768 airspace (IAS) (B), and leaf thickness (LT) (C). The grey histograms show the frequency of r
769 values visited by the MCMC sampler following a 20% burn-in period, red lines show the median
770 r values, and dashed black lines show the 95% credible interval.











