

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26

Flavor-Cyber-Agriculture: Optimization of plant metabolites in an open-source control environment through surrogate modeling

Arielle J. Johnson<sup>1¶</sup>, Elliot Meyerson<sup>2,3¶</sup>, Timothy L. Savas<sup>1</sup>, Risto Miikkulainen<sup>2,3&</sup>, Caleb B. Harper<sup>1&\*</sup>

<sup>1</sup>Media Lab, Massachusetts Institute of Technology, Cambridge, Massachusetts, United States of America

<sup>2</sup>Department of Computer Science, University of Texas at Austin, Austin, Texas, United States of America

<sup>3</sup>Sentient Technologies, San Francisco, California, United States of America

\* Corresponding Author

E-mail: [calebh@media.mit.edu](mailto:calebh@media.mit.edu) (CBH)

¶ These authors contributed equally to the work

& CBH and RM are Joint Senior Authors

## 27 **Abstract**

28 Food production in conventional agriculture faces numerous challenges such as reducing waste, meeting  
29 demand, maintaining flavor, and providing nutrition. Contained environments under artificial climate  
30 control, or cyber-agriculture, could in principle be used to meet many of these challenges. Through such  
31 environments, phenotypic expression of the plant---mass, edible yield, flavor, and nutrients---can be  
32 actuated through a “climate recipe,” where light, water, nutrients, temperature, and other climate and  
33 ecological variables are optimized to achieve a desired result. This paper describes a method for doing  
34 this optimization for the desired result of flavor by combining cyber-agriculture, metabolomic phenotype  
35 measurements, and machine learning. In a pilot experiment, (1) environmental conditions, i.e.  
36 photoperiod and ultraviolet (UV) light (known to affect production of flavor-active molecules in edible  
37 plants) were applied under different regimes to basil plants (*Ocimum basilicum*) growing inside a  
38 hydroponic farm with an open-source design; (2) flavor-active volatile molecules were measured in each  
39 plant using gas chromatography-mass spectrometry (GC-MS); and (3) symbolic regression was used to  
40 construct a surrogate model of this chemistry from the input environmental variables, and search in this  
41 model was used to discover new combinations of photoperiod and UV light to increase this chemistry.  
42 These new combinations, or climate recipes, were then implemented in the hydroponic farm, and several  
43 of them resulted in a marked increase in volatiles over control. The process demonstrated a “dilution  
44 effect”, i.e. a negative correlation between weight and desirable chemical species. It also discovered the  
45 surprising effect that a 24-hour photoperiod of photosynthetic-active radiation, the equivalent of all-day  
46 light, induces the most flavor molecule production in basil. In this manner, surrogate optimization through  
47 machine learning can be used to discover effective recipes for cyber-agriculture that would be difficult  
48 and time-consuming to find using hand-designed experiments.

## 50 **Introduction**

51

52 The so-called “dilution effect,” noted since the 1940’s and systematically reviewed since the early 1980’s  
53 [1], describes an inverse relationship between yield and nutrient concentration in food: For many  
54 nutritionally-important chemical constituents of food plants, such as minerals, protein, and vitamins, an  
55 increase in biomass is accompanied by a decrease in nutrient concentration. This effect has been  
56 systematically demonstrated in historical nutrient content studies over the last 50-70 years [2,3], as well  
57 as in controlled side-by-side trials that have shown a relationship between nutrient dilution and genetics  
58 [4], artificial fertilization [5], and elevated carbon dioxide levels related to climate change [6,7]. Flavor,  
59 known to be an important element of food and of eating behavior for organisms from insects to humans  
60 [8], has been declining alongside nutrients over approximately the last 50 years [9-11] in inverse  
61 proportion to rising yields. Flavor-active molecules in plants frequently have either positive health  
62 benefits (antioxidant, antimicrobial, anti-inflammatory) themselves or signal the presence of other  
63 beneficial or essential molecules, for example by being the enzymatic products of precursors (e.g. pro-  
64 vitamin A carotenoids, essential amino or fatty acids) necessary for human nutrition and health [9].

65

66 Vertical farming, or more generally cyber-agriculture, is a plant-growing format employing contained  
67 environments where light, water, nutrients, temperature, and other climate variables are provided  
68 artificially under computer control [12-14]. Data from environmental sensors informs the actuation of  
69 climatic conditions according to a “recipe” that is designed for best possible outcome, such as largest  
70 yield, best flavor, desired nutrients, and least cost. With cyber-agriculture, in principle it may be possible  
71 to increase food production quality and quantity, minimize waste and cost, and grow food with optimized  
72 climate recipes anywhere including locations otherwise unable to support agriculture. Conventional  
73 agriculture has been optimized for yield. What if it were optimized for quality and flavor?

74

75 This paper describes a proof-of-concept method aimed at optimizing flavor in a cyber-agricultural  
76 controlled environment, and a pilot experiment to validate this method. An experimental container, called

77 the Food Computer [12], was built at the MIT Media Lab with sensors, actuators, and computer control.

78 Basil (*Ocimum basilicum*) was chosen as a model organism because it has a fast growth cycle (five

79 weeks), and because the outcome can be readily measured in terms of fresh weight (quantity), and

80 chemical analysis of flavor (quality). A number of known growth recipes were implemented, together

81 with a broad range of their variations [15]. Machine learning technology [16-18] was then used to

82 optimize these recipes further. That is, based on these recipes and their associated outcomes, a surrogate

83 model was first constructed using symbolic regression. To keep the search problem manageable, the

84 optimization focused on the lighting conditions, keeping the other variables constant. The surrogate

85 model was then searched to discover potentially better lighting recipes, which were then tested in the

86 experimental container.

87

88 The light conditions had a large effect on the outcome, and the surrogate optimization method was able to

89 discover meaningful recipes. For instance, it discovered the well-known principle that flavor can be

90 traded off with mass, a version of the “dilution effect”: optimizing for flavor produced smaller plants,

91 while optimizing for mass produced less flavor. However, it also demonstrated how the approach can

92 discover new and surprising recipes, i.e. those that are counterintuitive but produce better outcomes. In

93 particular, the common-sense assumption that basil needs a few hours of darkness each day turned out to

94 be incorrect: The highest density of flavor molecules was produced through a 24-hr photoperiod, which

95 optimization discovered quickly and reliably. The results thus demonstrate that surrogate optimization and

96 machine discovery can be used to find growth recipes that are both effective and surprising.

97

## 98 **Measuring and optimizing flavor**

99

100 Flavor is largely a phenomenon of olfaction [19], and many aroma molecules are produced by the

101 secondary metabolism of plants. Plants have a particularly rich secondary or specialized metabolism [20],

102 a set of biosynthetic pathways synthesizing molecules that are not essential for the basic processes of life  
103 (cell division, reproduction, etc.) but rather confer fitness and adaptive advantage to the organism in its  
104 ecological niche [21], related to stress tolerance, defense, and communication [22]. Their expression and  
105 induction depends to various degrees on environmental and ecological conditions [23].

106

107 Cyber-physical agriculture methods such as the Food Computer (FC), where data from environmental  
108 sensors informs the actuation of climatic conditions according to a climate recipe [12-14] present unique  
109 opportunities for inducing plant phenotypic changes through environmental/ecological conditions alone.  
110 One example of this approach is to apply the ecological stresses to which adaptations have evolved as  
111 specific biosynthetic pathways.

112

113 *O. basilicum*, the basil plant, is typical of herb plants in that it produces many aromatic molecules,  
114 particularly the terpenoids 1,8-cineole, linalool, camphor, borneol, bergamotene, and farnesene, and the  
115 phenylpropenes eugenol, methyleugenol, and estragole [24]. These molecules are thought to play varying  
116 roles in stress adaptation and defense, and the production by the basil plant of aromatic molecules has  
117 been shown to increase upon exposure to these stresses, including water stress [25], ultraviolet and PAR  
118 light [26–28], heat [29], bacteria [30], chitosan (a compound derived from chitin, found in insect  
119 exoskeletons and fungal cell walls, [31]), and sodium and other minerals [32].

120

121 This paper explores methods for increasing flavor molecule production in *O. basilicum*, using: (1)  
122 ultraviolet light, PAR, and photoperiod as environmental and stress variables; (2) Gas Chromatography-  
123 Mass Spectrometry for semiquantitative analysis of volatiles; (3) surrogate optimization for discovering  
124 conditions that will maximize production of these volatiles.

125

## 126 **Materials and Methods**

127

128 This section describes the design of the Food Computer, i.e. the physical container environment used in  
129 the pilot experiment with basil. It also describes the process for growing basil in this environment, and  
130 methods for measuring the growth outcome in terms of weight and chemistry.

131

## 132 **Food Computer**

133 All basil plants were grown in a Food Server, a multi-tray, multi-rack hydroponic configuration of  
134 the OpenAg Food Computer™ (FC) environment [12]. Basil plants were germinated in engineered foam  
135 rooting cubes (Oasis Grower Solutions, Kent, OH), then transplanted to 36-position (4×9) food-grade  
136 resin floating lettuce rafts (Beaver Plastics, Acheson, AB, Canada) at 14 days of age. The plants were  
137 grown in a shallow water culture hydroponic system using 56.6-liter trays (Botanicare, Chandler, AZ)  
138 supplied by 75-liter reservoirs (Botanicare) and 700 gallon-per-hour rated Pondmaster magnetic drive  
139 pumps (Danner Manufacturing, Islandia, NY), with nutrient solutions (a “15-0-0” Calcium Nitrate  
140 solution and a “5-12-26” 5% Nitrate, 12% Phosphate, 26% soluble Potash solution combined with water  
141 for a final concentration of 150 ppm Nitrogen, 116 ppm Calcium, 52 ppm Phosphorus, 215 ppm  
142 Potassium; JR Peters, Allentown, PA) added by a water-powered proportional chemical injector  
143 (Dosatron, Clearwater, FL).

144 The FC was set up with trays in vertical stacks of three (denoted 0, 1, and 2) within a custom  
145 powder-coated steel frame (Indoor Harvest, Houston, TX). Each stack was thermally isolated with  
146 reflective foil captive-bubble insulation (Reflectix, Markleville, IN) and climate-controlled with a 10,000  
147 BTU air conditioning unit (AeonAir, Wilmington, DE). Three types of photosynthetic-active radiation  
148 (PAR) lights were used: Agrobrite high output T5 fluorescent fixtures (Hydrofarm, Fairless Hills, PA),  
149 Illumitex ES2 Eclipse red and blue LED fixtures (Illumitex, Austin, TX) and Phillips GreenPower deep  
150 red/blue LED production modules (Phillips, Somerset, NJ). Lights were fixed at a distance of 40 cm from

151 the foam float. Reptisun 10.0 UVB T5 High Output ultraviolet lights were added to treatment conditions  
152 (Zoo Med, San Luis Obispo, CA), also at a distance of 40 cm.



153

154 **Figure 1. Images inside the MIT Media Lab Food Server taken during the experiment.**

155

## 156 **Plant species and climate recipes**

157 Common Sweet Basil (*O. basilicum* var “Sweet”) seeds (Eden Brothers, Arden, NC) were used in  
158 the pilot experiment. From 14 days of age to harvest, they were grown in identical trays as described in  
159 “Food Computer” above, with one of three light types as the only source of PAR. Control conditions were  
160 grown with the PAR light; experimental treatment conditions had supplemental UV light. Treatment  
161 conditions, or “Climate Recipes”, in Rounds 2 and 3 of the experiment were determined based on  
162 suggestions from the surrogate optimization of chemscore from the previous round. The data from Round  
163 1 determined the conditions of Round 2, and the data from Round 2 determined the conditions of Round  
164 3).

165

## 166 **Harvest, weight and length measurement**

167 All plants in each round of the experiment were harvested on the same day. Four plants from each  
168 treatment condition were used for volatile analysis and the remaining 32 were used for height and weight  
169 measurements. Weight measurements were taken with roots removed.

170

## 171 **Sampling and sample preparation**

172 Immediately after harvesting, leaves were sampled from four plants from each treatment  
173 condition. Fifteen leaves from each plant were harvested: five from near the base, five from the middle,  
174 and five from the top, with each set selected randomly. Leaves were immediately frozen with dry ice or  
175 liquid nitrogen, homogenized into a powder, and kept frozen. The amount of 1 gram of frozen plant tissue  
176 was transferred to a 20 mL amber glass headspace vial (Supelco, Bellefonte, PA) and 2 mL of saturated,  
177 cold calcium chloride solution in distilled water was added to prevent enzymatic reactions. The vials were  
178 capped with magnetic, PTFE-lined silicone septa headspace caps (Supelco) and kept on ice before  
179 being transferred to GC-MS.

180

## 181 **Volatile Analysis**

182 The method of Johnson et al. [33] was adapted for the experiment. Sample vials were placed in  
183 the tray of the Gerstel MPS2 autosampler (Gertsel, Linthicum, MD), which performed the extraction and  
184 injection. One vial at a time was warmed to 40°C and agitated at 500 rpm for 5 minutes directly before  
185 extraction. A conditioned, 2-cm long 50/30  $\mu\text{m}$ -thick PDMS-DVB SPME fiber (Supelco) was introduced  
186 into the headspace of the vial for 45 minutes at 40°C with rotational shaking at 250 RPM. The fiber was  
187 removed from the headspace of the vial and immediately introduced into the inlet of an Agilent model  
188 7890 Gas Chromatograph- single quadrupole-MS (GC-MS) (Agilent Technologies) with a DB-5 column  
189 (30 meters long, 0.25 mm ID, 0.25  $\mu\text{m}$  film thickness, J&W Scientific, Folsom, CA). The inlet was held



190 at 250°C with a 2:1 split and had a 0.75mm i.d. SPME inlet liner installed (Agilent Technologies). The  
191 carrier gas was Helium, at a constant flow rate of 1 mL/minute. The starting oven temperature was 40°C,  
192 held for 3 minutes, followed by a 2°C/minute ramp until 180°C was reached, then the ramp was increased  
193 to 30°C/minute until 250°C was reached, and held for 3 minutes. The total runtime was 47 minutes. The  
194 transfer line to the mass spectrometer was held at 250°C, the source temperature was 230°C, and the  
195 quadrupole temperature was 150°C. The mass spectrometer had a 1.5-minute solvent delay and was run in  
196 scan mode with Electron Impact ionization at 70eV, from  $m/z$  40 to  $m/z$  300.

197 Compounds were identified and recorded based on a 90% or higher match using the NIST Mass  
198 Spectral Database and a signal to noise ratio above 10. Analyte peaks were integrated on the Total Ion  
199 Chromatogram.

200

## 201 **Optimization metric: Chemscore**

202 Optimizing the target metric should correspond to maximizing flavor in a general sense. The  
203 metric should also be robust to noise, since the number of evaluations is limited, and low-dimensional to  
204 make optimization easier.

205 Basil, like most foods, contains multiple molecules contributing to flavor. An average GC-MS  
206 chromatogram of basil contains around 30-40 different volatile molecules, with concentrations varying  
207 over several orders of magnitude. To construct a single metric to optimize, this GC-MS data is aggregated  
208 across samples and chemicals as the Chemscore. This score is a weighted average of the volatile profile  
209 compared to the control condition. It is a holistic placeholder for how flavorful a sample is, while  
210 normalizing for varying scales and distributions of different chemicals. Seventeen chemicals common  
211 across all GC-MS measurements were selected into the calculation of chemscore.

212

## 213 **Comparison metrics: R-score and Z-score**

214 For further comparison, an R-score and a Z-score, across all volatiles in a sample, were calculated  
215 for each treatment condition. The R-Score, the average ratio of volatiles in a treatment condition over  
216 their average in the control conditions in a round of the experiment, facilitates comparison of results  
217 across the three rounds of the experiment, under the assumption that uncontrollable environmental  
218 differences across rounds are captured in differing control results. The Z-score, which compares the  
219 abundances of each volatile molecule in a sample over or below its average in all samples in a round,  
220 gives a sense of the overall spread of results in the experiment.

221

## 222 **Methods: Surrogate modeling and optimization**

223 In optimization settings where the target function is expensive to evaluate (either temporally or  
224 financially), e.g., in the case of growing basil to maturity, surrogate-based optimization is a common  
225 method for minimizing the number of evaluations required to achieve an acceptable solution [34–36]. To  
226 choose the next samples to evaluate, surrogate methods build an explicit predictive model of the solution  
227 landscape and select the most promising samples according to this “surrogate model”. To implement such  
228 a method, input variables need to be defined, a class of regression models needs to be selected, and a  
229 method for discovering the next samples (recipes) from these models needs to be developed. This section  
230 details the development of these choices for the experiment in this paper, and notes methods for scaling  
231 up future work.

232

### 233 **Input Dimensions**

234 For this experiment, a recipe was defined by three input variables: photoperiod, UV period, and  
235 PAR (photosynthetically active radiation). Three input variables was an appropriate dimensionality for  
236 this pilot experiment, following the general rule-of-thumb that, for surrogate-based methods, the number  
237 of evaluations required to achieve reasonable results is around ten times the number of dimensions [34].

238 These variables were chosen because they are already known to increase the accumulation of volatiles  
239 [26–28], and are relatively simple to control in the described hardware setup.

240 Photoperiod is the number of hours the primary light panel is turned on each day. Recipes can  
241 thus have photoperiod values anywhere from 0 to 24hrs. Photoperiod is known to have significant effects  
242 on the accumulation of biomass and leaf area in plants [37], as well as the formation of trichomes, the  
243 structure that store flavor-active volatiles, in the *Thymus vulgaris* (thyme) plant [38], a botanical cousin,  
244 as they are both members of the *Lamiaceae* family, to basil. In addition, photoperiod has been shown to  
245 change the volatile profile of basil [39].

246 UV period is the number of hours per day plants receive supplemental UV-B radiation. Like  
247 photoperiod, UV period can take on values anywhere from 0 to 24hrs. UV has previously been shown to  
248 increase volatile content in basil [26]; it is included so that its effects can be validated and optimized in  
249 the Food Computer hardware setting.

250 PAR (Photosynthetically Active Radiation) is the amount of light available for photosynthesis. In  
251 the Food Computer setup, the PAR is determined by the primary light panel. There were nine light panels,  
252 each with a unique PAR value. To set PAR values for a batch of nine recipes, one light panel was assigned  
253 to each recipe. Thus, in contrast to photoperiod and UV period, each available PAR value can be used  
254 only once in each batch. This kind of hardware resource matching constraint is not common in either  
255 computer or physical experiments, so a custom optimization method must be developed.

256

## 257 **Regression Model**

258 Symbolic regression [40–42] was used for building surrogate models to predict chemscore from  
259 the input variables. Symbolic regression uses evolutionary optimization to discover nonlinear algebraic  
260 expressions that serve as surrogate models. For the experiment in this paper, a multi-objective Pareto  
261 optimization procedure was used [43,44]. The first objective is to minimize error, i.e., MSE with respect  
262 to predicting chemscore; the second objective is to maximize parsimony, i.e., minimize the size of the

263 algebraic expression (number of nodes). The fitting procedure then yields a Pareto front of models, from  
264 which a new batch of recipes can be selected.

265 For the flavor-optimization problem, symbolic regression has several advantages over other  
266 popular choices for surrogate models. First, by simultaneously optimizing for error and parsimony, search  
267 is biased towards the kinds of compact algebraic expressions that are desirable in the natural sciences  
268 [44]. These expressions are more interpretable than other regression models, because the relationships  
269 between variables can be read off directly from the expression. Such interpretability can lead to a better  
270 understanding of the search space, which helps in developing better models for future experiments.

271 Second, whereas surrogate models such as Gaussian processes can only interpolate, symbolic  
272 regression can extrapolate. Interpolation is sufficient when iterative incremental improvement can  
273 eventually lead to an optimal solution. However, in the experiment in this paper, only a single parallel  
274 batch of recipes is selected via surrogate optimization to be implemented in the Food Computer. So, it is  
275 advantageous to consider strong optimistic predictions a model makes about sparse regions in the recipe  
276 space. Note that if this process were used over multiple iterations, an inordinate amount of resources  
277 could be spent at the extremes of the recipe space.

278 Third, symbolic regression is robust to normalization of input and output variables: It  
279 automatically discovers reasonable scaling factors to use through optimized constants that are found to be  
280 useful in model expressions.

281 It is important to note that symbolic regression can have significant drawbacks as well [43]. First,  
282 it is computationally expensive compared to other regression methods; however in this paper,  
283 computation time is negligible compared to the time it takes to grow a batch of basil recipes. Second,  
284 surrogate optimization with symbolic regression models currently lacks theoretical convergence  
285 guarantees and performance bounds. Such convergence guarantees have potential practical benefits over  
286 many iterations of surrogate optimization; however, since only a single such iteration is performed in the  
287 experiment in this paper, such guarantees are unnecessary.

288

## 289 **Recipe Discovery**

290           There were three rounds of growing experiments. In each round, there are nine trays of basil  
291 growing in parallel. To ensure consistency across rounds, three of these nine trays are fixed to control  
292 recipes. This setup leaves six non-control recipes to be selected.

293           In the first round, recipes were selected by hand [15] to investigate the effects of UV supplement  
294 and choice of light panel. To add the photoperiod dimension, and create initial diversity in the recipe  
295 space, recipes in the second round were chosen by an unsupervised method: Six non-control recipes were  
296 found as centroids of Voronoi tessellation (CVT) given the first round of recipes [45]. Following a trust  
297 region approach [35], to implement the bias that good solutions are likely to be relatively close to expert  
298 hand-designed recipes, values for each dimension were constrained to be with a constant distance of  
299 previously evaluated values.

300           In the third round, recipes were selected from symbolic regression surrogate models [46]. Each  
301 run of symbolic regression yields a collection of models on the error-parsimony Pareto front. These  
302 models were clustered to determine an error threshold, above which models were underfitting. The six  
303 most parsimonious models not underfitting were then used to define a recipe to run in parallel. Since the  
304 recipe space has only three dimensions it is computationally efficient to use a dense grid search to select a  
305 recipe that maximizes expected chemscore. Greedy sequential selection is the most popular approach to  
306 constructing parallel batches from surrogates [47,48]. The recipes were thus selected sequentially in  
307 increasing order of model error. Such a selection handles the constraint that each available PAR value can  
308 be selected only once per round. If a variable is ignored by a model, the value of the variable is set to  
309 maximize exploration, since the model has indicated that exploitation of this variable is currently not  
310 useful.

311

## 312 **Results and Discussion**

313 The average weight, chemscore, total peak area of volatiles on the chromatogram, and chemscore as a  
 314 percentage of the control condition chemscore are presented in Table 1. Weight was recorded as the  
 315 weight of above-ground plant parts; roots were excluded.

316 **Table 1. Treatment conditions (UV and PAR photoperiod) and weight and chemical results.**

Round	Bay	Tray	UV Photoperiod	PAR Photoperiod	PAR	Weight (grams)	R score	Chemscore	Z score	R score
1	1	0	18	18	636.92	32.00	0.85	-0.77	0.65	
1	1	1	18	18	798.42	102.71	1.00	0.21	1.15	
1	1	2	18	18	832.58	133.59	1.06	0.44	1.37	
1	2	0	0	18	820.25	72.08	1.13	0.46	1.45	
1	2	1	0	18	1,098.75	235.44	0.81	-0.68	0.79	
1	2	2	0	18	403.58	84.33	1.06	0.33	1.34	
2	0	0	9	21.5	867.33	74.18	<b>1.81</b>	1.07	0.68	
2	0	1	9	21.5	445.25	65.63	1.15	-0.01	0.10	
2	0	2	9	21.5	735.42	63.86	<b>1.61</b>	0.86	0.50	
2	1	0	9	14.5	636.92	112.89	0.89	-0.43	-0.25	
2	1	1	9	14.5	798.42	189.00	0.58	-1.07	-0.52	
2	2	0	0	18	820.25	154.50	0.92	-0.42	-0.19	
2	2	1	0	18	1,098.75	211.00	0.73	-0.58	-0.28	
2	2	2	0	18	403.58	112.00	1.35	0.57	0.27	
3	0	0	17.45	24	867.33	137.44	<b>16.57</b>	2.38	-0.28	<b>14.05</b>
3	0	1	4.12	24	445.25	71.25	<b>2.33</b>	-0.21	-1.03	<b>1.83</b>
3	0	2	24	24	735.42	49.33	<b>2.84</b>	-0.05	-1.01	<b>2.12</b>
3	1	0	14.06	24	636.92	80.51	<b>2.00</b>	-0.30	-1.05	1.47
3	1	1	8.48	17.18	798.42	62.78	1.80	-0.34	-1.06	1.34
3	1	2	10.67	22.5	832.58	88.83	<b>2.09</b>	-0.28	-1.04	<b>1.55</b>
3	2	0	0	18	820.25	92.89	0.80	-0.66	-1.11	0.60
3	2	1	0	18	1,098.75	126.86	1.20	-0.53	-1.09	0.94

3	2	2	0	18	403.58					1.47
---	---	---	---	----	--------	--	--	--	--	------

317 “Bay” specifies the position in the vertical stack of three hydroponic trays, with “0” closest to the floor.

318 One tray in each bay contained a control condition, which had zero hours UV photoperiod and 18 hours

319 PAR photoperiod. R score > 1.5 is denoted in bold. The photoperiod hours range between 0 and 24. PAR

320 values indicate  $\mu\text{mole}/\text{m}^2\text{s}$  photosynthetic photon flux density. The R-Scores were calculated with missing

321 control imputed.

322

323 The table includes results both with and without imputed data for the control condition whose data was

324 lost in Round 3 of the experiment (denoted by dark grey boxes in table 1). Assuming control results are

325 consistent within each round, they make the results easier to compare across rounds. Imputed values for

326 each chemical for the missing control treatment in Round 3 were computed by regression, i.e., by solving

327 a fully-determined linear system that predicts the value of the third control from the other two, based the

328 values of the controls in the previous two Rounds.

329

330 Table 2 gives the correlations between input variables and metrics (Spearman, to account for nonlinearity

331 in the metrics). Correlations >0.45 are in bold. Note in particular that the R-scores are negatively

332 correlated with weight: Optimizing for flavor thus results in smaller plants, and larger plans have less

333 flavor, thus illustrating the “Dilution effect.”

334

335 **Table 2: Spearman correlations between selected input variables and metrics**

	R Score	Weight	Chemscore	Z Score
UV	0.355	-0.336	0.199	0.058
Photoperiod	<b>0.763<sup>a</sup></b>	-0.355	<b>0.477<sup>a</sup></b>	-0.149

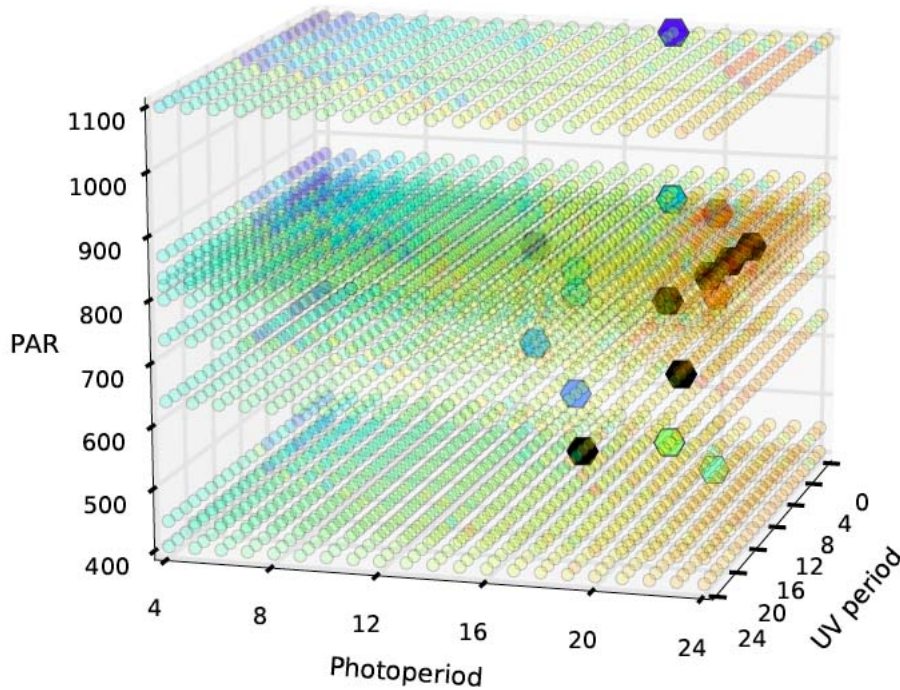
PAR	-0.131	<b>0.541<sup>a</sup></b>	-0.142	-0.070
R Score		<b>-0.471<sup>a</sup></b>	<b>0.637<sup>a</sup></b>	-0.226
R Score Imputed	<b>0.967<sup>a</sup></b>	<b>-0.502<sup>a</sup></b>	<b>0.764<sup>a</sup></b>	-0.055

336

337

338 In the first round, where an 18-hour PAR photoperiod and an 18-hour UV photoperiod were selected by  
339 hand, R-score and Chemscore did indicate that UV light or photoperiod increases volatiles. In the second  
340 round, two R-scores (both with UV light and extended PAR photoperiod of 21 hours) were above 1.5,  
341 meaning that volatiles holistically increased 50% over control. In the third round, several conditions  
342 resulted in an R-score that met or exceeded this threshold, with many conditions (all with PAR  
343 photoperiods of 22.5-24 hours and UV periods of 4-17 hours) doubling the volatile profile compared to  
344 control. The discovery of the recipes in Round 3 from the model is illustrated in Figure 2.



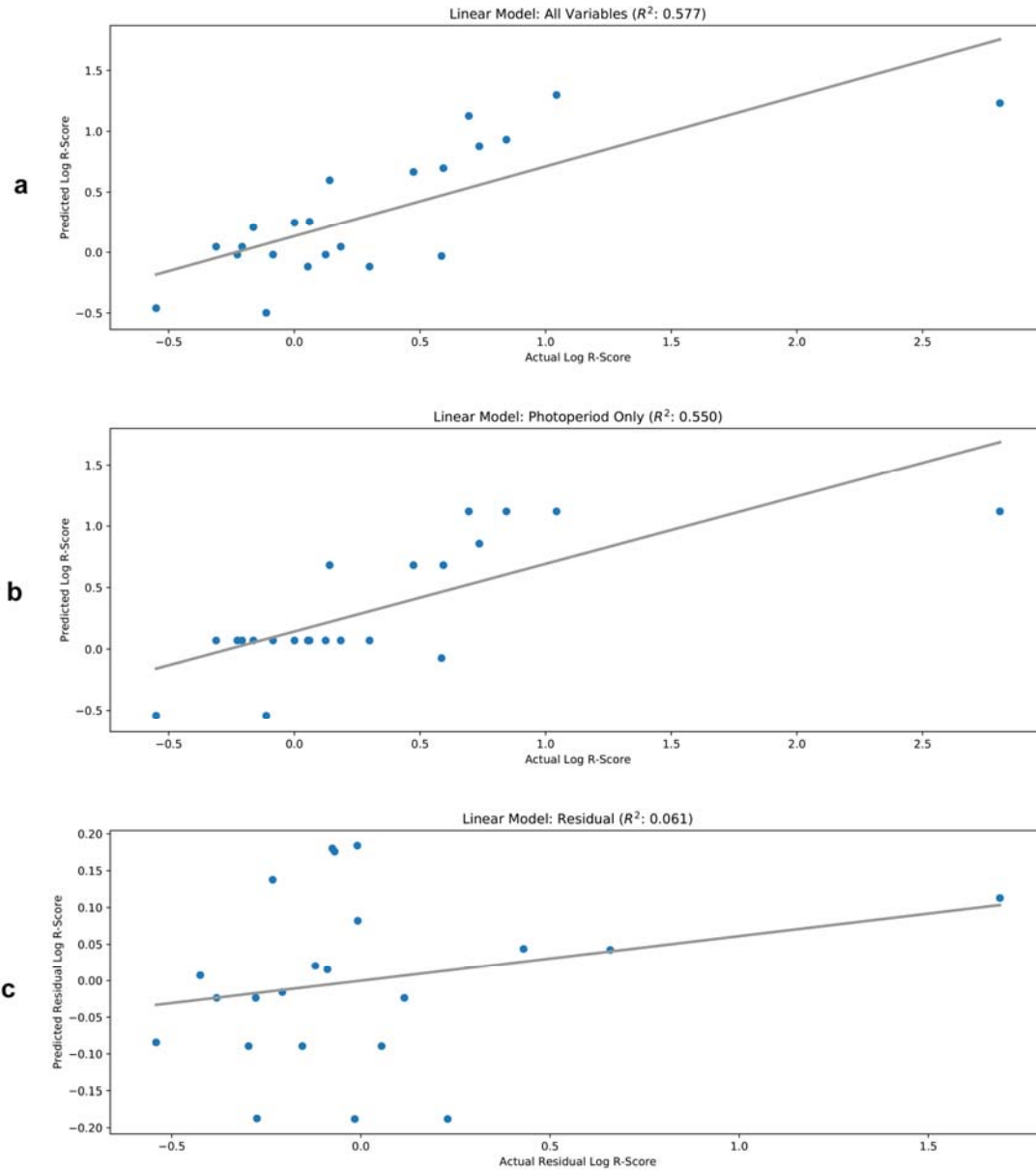


345

346 **Figure 2. An illustration of the surrogate model and the recipes suggested by the optimization.** The  
347 three axes correspond to the three actuators and the color of the small dots indicates their value predicted  
348 by the model (i.e. flavor; red > yellow > green > blue). The large dots are suggestions, and the darker dots  
349 are the most recent ones. They suggest utilizing long photoperiods and UV periods, the success of which  
350 was confirmed in growth experiments in the Food Computer.  
351

352 The most striking discovery in this experiment was the positive effect of a 24-hour photoperiod, i.e.,  
353 constant daylight. This result replicated evidence on the volatile profile effects of a 24-hour photoperiod  
354 described by Skrubis et al. [39], who found that basil plants grown with a 24-hour photoperiod weighed,  
355 upon maturity, approximately 25% more than plants grown with a nine-hour photoperiod (although they  
356 took three days longer to reach maturity) and 27% more than plants grown outdoors in natural light with  
357 an approximately 15-hour photoperiod . That study also characterized changes in the relative volatile  
358 profiles of those basil plants, but not absolute volatile content, so comparisons to chemscore in the current  
359 work are not possible. The 24-hour photoperiod discovery is notable because the hand-designed  
360 experimental conditions in Round 1 had a photoperiod of 18 hours, and the experimenters and the model

361 were blind to the Skrubis et al. study. The surrogate optimization approach nevertheless iterated the  
362 recipes into the 24-hour photoperiod, where it had a strong positive effect.  
363  
364 Aside from the high R-score in Table 1 , further evidence for the importance of photoperiod can be seen  
365 in the high correlation between R-score and Photoperiod in Table 2, and in the regression process itself:  
366 For each run of symbolic regression, the most parsimonious nontrivial model had the form  $y = cp$ , for  
367 some constant  $c$ , where  $p$  is the photoperiod. Also, Figure 3A shows the a linear model trained on all  
368 three light variables to fit the log R-score. Fig 3B shows a linear model of R-score based on photoperiod  
369 alone. Fig 3C shows the predictions of a linear model trained on all three variables, but with the effect of  
370 photoperiod removed, i.e., it is trained to fit the residuals. These modeling results are similar with  
371 imputed and outlier-handled data. The low performance of the residual model suggests that photoperiod  
372 had such a dominating effect that the effects of other variables were effectively noise. However, since  
373 significant effects of UV have been reported in previous work [26,27] and are not seen here, it is also  
374 possible that there are significant nonlinear dynamics that require further trials and nonlinear modeling to  
375 uncover and exploit.



376

377 **Figure 3. Linear regression analysis of actual vs. calculated log R-score for three different models.**

378 A: A linear model trained on UV, photoperiod, and PAR. B: A linear model trained on photoperiod only.

379 C: A linear model trained on residuals after removing photoperiod effect. Photoperiod dominates the other

380 variables (or possible there are significant nonlinear effects between these variables).

381

## 382 Discussion and Future Work

383 The experiment described in this paper confirmed that the design of climate recipes impact the

384 accumulation of volatile flavor molecules in basil, and it is possible to discover good recipes iteratively

385 through machine learning. The recipes discovered in this case replicated known principles (such as the  
386 weight/flavor tradeoff), and also demonstrated the possibility for discovering previously unknown,  
387 surprising principles (like the 24hr photoperiod). The 24-hour photoperiod in particular is impossible in  
388 nature (except around the summer solstice within the Arctic and Antarctic circles) and therefore unlikely  
389 to be discovered, except in controlled environments for cyber-physical agriculture.

390

391 The most immediate direction of future work is to expand the current experiment to a larger search space.

392 A facility with four containers, making it possible to evaluate an order of magnitude more recipes at once

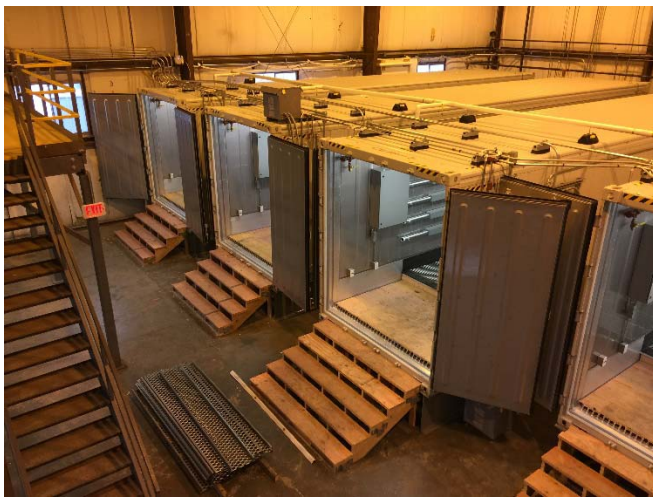
393 is in development at MIT and illustrated in Fig 4. This facility will make it possible to control a number of

394 other actuators besides light, including temperature, pH, nutrient concentration, microbial and other

395 additives, and different cultivars. It will also be possible to measure the energy and other costs associated

396 with the recipes, as well as objectives such as nutrient components, density, and yield, and more elements

397 of flavor (single compounds, and ratios of compounds).



398

399 **Figure 4. Images of MIT expansion facility under development.**

400

401 In terms of surrogate optimization, more iterations will be run to build more accurate models, and to

402 determine the proper stopping point of the method, i.e. run it until it has likely converged. The approach

403 will be extended to cover the larger search space as well as multiple objectives. Most likely, different  
404 models and optimizers will be necessary. In low-dimensional settings with unknown nonlinearities and a  
405 relatively small number of samples, Kriging [34], Gaussian processes [36,49], and symbolic regression  
406 [44] are suitable choices for building a regression model of natural phenomena. When the dimensionality  
407 and number of samples increases, deep neural networks may be a better model of the solution landscape  
408 [47,50,51], and evolutionary optimization a better way to determine most promising samples [43–45].

409  
410

411 The next step will be to extend the experiment to other plants, such as cotton, where the goal is not to  
412 optimize flavor but physical properties such as strength and length of the fibers. It will be important to  
413 verify that such plants are viable to grow artificially, and that such properties can be optimized with  
414 available actuators, in isolation and in combination with other properties. Future extensions to other areas  
415 may include biofuels and plants with specific medicinal value.

416

417 The third future step is to extend the optimization from static recipes to time-varying recipes, i.e.  
418 optimizing the actuators during the entire growth period of the plant. Of particular interest are different  
419 stress periods when the plant is exposed to, for example, drought or signals of predators (e.g. through  
420 chitosan added to the growth medium). Such periods may produce a response in the plant that results in  
421 more flavor or more rapid growth, for example. Such recipes should be reactive, i.e. conditional to real-  
422 time measurements of the growth status. One possibility is to use machine learning to establish a mapping  
423 from visual images of the plant to more destructive measurements such as chemical concentrations. Such  
424 optimization spaces are very high-dimensional, most likely making it necessary to use evolutionary  
425 optimization, and perhaps neuroevolution to construct a mapping from sensory time series to optimal  
426 actions [55,56].

427

## 428 **Conclusion**

429 Computer-controlled growth environments are a promising approach for the future of agriculture,  
430 potentially maximizing production and quality and minimizing waste and cost. Initial experiments with  
431 basil (*O. basilicum*) suggest that the cyber-physical approach to agriculture is indeed viable: such  
432 environments can be built, the plants thrive in them, the climate recipes make a difference in growth  
433 outcomes, and machine learning can be used to discover good recipes automatically. Future steps should  
434 verify these results on other plants, expand to larger search spaces with more actuators, and to optimizing  
435 entire growth periods. Higher-volume food computers need to be built and more powerful optimization  
436 methods employed, but the results suggest that such extensions are worthwhile.

437  
438

## 439 **Acknowledgements**

440 The authors would like to thank Christina Agapakis, Nate Tedford, and Scott Marr for access to  
441 and support on analysis instrumentation and Babak Hodjat and Hormoz Shahrzad for modeling  
442 and optimization insights and comments on the manuscript.

443

## 444 **References**

- 445 1. Jarrell WM, Beverly RB. The dilution effect in plant nutrition studies. *Adv Agron.*  
446 Elsevier; 1981;34: 197–224.
- 447 2. Davis DR, Epp MD, Riordan HD, Davis DR. Changes in USDA Food Composition Data  
448 for 43 Garden Crops, 1950 to 1999. *J Am Coll Nutr.* 2004;23: 669–682.  
449 doi:10.1080/07315724.2004.10719409
- 450 3. Davis DR. Declining fruit and vegetable nutrient composition: What is the evidence?  
451 *HortScience.* 2009;44: 15–19.
- 452 4. Farnham MW, Grusak MA, Wang M. Calcium and magnesium concentration of inbred  
453 and hybrid broccoli heads. *J Am Soc Hortic Sci.* 2000;125: 344–349. Available:  
454 <http://www.scopus.com/inward/record.url?eid=2-s2.0-0033998604&partnerID=tZOtx3y1>



- 455 5. Hughes M, Chaplin MH, Martin LW. Influence of mycorrhiza on the nutrition of red  
456 raspberries. *HortScience*. 1979;14: 521-523.
- 457 6. Loladze I. Rising atmospheric CO<sub>2</sub> and human nutrition: Toward globally imbalanced  
458 plant stoichiometry? *Trends Ecol Evol*. 2002;17: 457–461. doi:10.1016/S0169-  
459 5347(02)02587-9
- 460 7. Cotrufo FM, Ineson P, Scott A. Elevated CO<sub>2</sub> reduces the nitrogen concentration of plant  
461 tissues. *Glob Chang Biol*. 1998;4: 43–54. doi:10.1046/j.1365-2486.1998.00101.x
- 462 8. Fraenkel GS. The Raison d’Etre of Secondary Plant Substances. *Science* (80- ). 1959;129:  
463 1466–1470. doi:10.1126/science.129.3361.1466
- 464 9. Goff SA, Klee HJ. Plant Volatile Compounds: Sensory Cues for Health and Nutritional  
465 Value? *Science* (80- ). 2006;311: 815–819. doi:10.1126/science.1112614
- 466 10. Folta KM, Klee HJ. Sensory sacrifices when we mass-produce mass produce. *Hortic Res*.  
467 Nature Publishing Group; 2016;3. doi:10.1038/hortres.2016.32
- 468 11. Schatzker M. *The Dorito Effect: The Surprising New Truth about Food and Flavor*. Simon  
469 and Schuster; 2015.
- 470 12. Harper C, Siller M. OpenAG: A Globally Distributed Network of Food Computing. *IEEE*  
471 *Pervasive Comput*. 2015;14: 24–27. doi:10.1109/MPRV.2015.72
- 472 13. Harper C. Open-Source Agriculture Initiative—Food for the Future? In: Kozai T, editor.  
473 *LED Lighting For Urban Agriculture*. Singapore: Springer Science+Business Media;  
474 2016. pp. 37–46. doi:10.1007/978-981-10-1848-0
- 475 14. Ferrer EC, Rye J, Brander G, Savas T, Chambers D, England H, et al. Personal Food  
476 Computer: A new device for controlled-environment agriculture. *arXiv Prepr*  
477 *arXiv170605104*. 2017;
- 478 15. Chernoff H. Sequential design of experiments. *Ann Math Stat*. JSTOR; 1959;30: 755–770.
- 479 16. O’Reilly U-M, Wagy M, Hodjat B. EC-star: A massive-scale, hub and spoke, distributed  
480 genetic programming system. *Genetic Programming Theory and Practice X*. Springer;  
481 2013. pp. 73–85.
- 482 17. Meyerson E, Miikkulainen R. Discovering Evolutionary Stepping Stones through  
483 Behavior Domination. *arXiv Prepr arXiv170405554*. 2017; Available:  
484 <https://arxiv.org/pdf/1704.05554.pdf>
- 485 18. Miikkulainen R, Liang J, Meyerson E, Rawal A, Fink D, Francon O, et al. Evolving Deep  
486 Neural Networks. 2017; Available: <http://arxiv.org/abs/1703.00548>

- 487 19. Small DM. Flavor is in the brain. *Physiol Behav.* Elsevier Inc.; 2012;107: 540–52.  
488 doi:10.1016/j.physbeh.2012.04.011
- 489 20. Weng J. The evolutionary paths towards complexity□: a metabolic perspective. *New*  
490 *Phytologist.* 2014; 201(4):1141–1149.
- 491 21. Moghe G, Last RL. Something old, something new: Conserved enzymes and the evolution  
492 of novelty in plant specialized metabolism. *Plant Physiol.* 2015;169: pp.00994.2015.  
493 doi:10.1104/pp.15.00994
- 494 22. Weng J-K, Philippe RN, Noel JP. The Rise of Chemodiversity in Plants. *Science* (80- ).  
495 2012;336: 1667–1670. doi:10.1126/science.1217411
- 496 23. Deschamps C, Simon JE, Wt. Terpenoid essential oil metabolism in basil (*Ocimum*  
497 *basilicum* L.) following elicitation. *J Essent Oil Res.* 2006;18: 618–621.  
498 doi:10.1080/10412905.2006.9699183
- 499 24. Lee S-J, Umamo K, Shibamoto T, Lee K-G. Identification of volatile components in basil  
500 (*Ocimum basilicum* L.) and thyme leaves (*Thymus vulgaris* L.) and their antioxidant  
501 properties. *Food Chem.* 2005;91: 131–137. doi:10.1016/j.foodchem.2004.05.056
- 502 25. Khalid KA. Influence of water stress on growth, essential oil, and chemical composition  
503 of herbs (*Ocimum* sp.). *Int Agrophysics.* 2006;20: 289–296.  
504 doi:10.1016/j.plantsci.2004.05.034
- 505 26. Nitz G, Schnitzler W. Effect of PAR and UV-B radiation on the quality and quantity of the  
506 essential oil content in sweet basil (*Ocimum basilicum* L.). *Acta Hortic.* 2004;659: 375–  
507 381.
- 508 27. Ioannidis D, Bonner L, Johnson CB. UV-B is required for normal development of oil  
509 glands in *Ocimum basilicum* L. (sweet basil). *Ann Bot.* 2002;90: 453–460.  
510 doi:10.1093/aob/mcf212
- 511 28. Chang X, Alderson PG, Wright CJ. Solar irradiance level alters the growth of basil  
512 (*Ocimum basilicum* L.) and its content of volatile oils. *Environ Exp Bot.* 2008;63: 216–  
513 223. doi:10.1016/j.envexpbot.2007.10.017
- 514 29. Chang X, Alderson P, Wright C. Effect of temperature integration on the growth and  
515 volatile oil content of basil ( *Ocimum basilicum* L.). *J Hortic Sci Biotechnol.* 2005;80:  
516 593–598. doi:10.1080/14620316.2005.11511983
- 517 30. Banchio E, Xie X, Zhang H, Pare PW. Soil Bacteria Elevate Essential Oil Accumulation  
518 and Emissions in Sweet Basil. *J Agric Food Chem.* 2009; 653–657.  
519 doi:10.1021/jf8020305

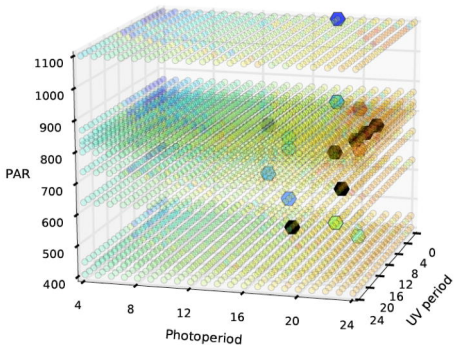


- 520 31. Kim H, Chen F, Wang X, Rajapakse N. Effect of chitosan on the biological properties of  
521 sweet basil (*Ocimum basilicum* L.). *J Agric Food Chem.* 2005;53: 3696–3701.  
522 doi:10.1021/kjf0480804
- 523 32. Said-Al Ahl HAH, Mahmoud AA. Effect of zinc and / or iron foliar application on growth  
524 and essential oil of sweet basil ( *Ocimum basilicum* L .) under salt stress. *Ozean Journal of*  
525 *Applied Sciences.* 2010;3: 97–111.
- 526 33. Johnson AJ, Hopfer H, Heymann H, Ebeler SE. Aroma Perception and Chemistry of  
527 Bitters in Whiskey Matrices: Modeling the Old-Fashioned. *Chemosens Percept.* 2017; 1–  
528 14. doi:10.1007/s12078-017-9229-3
- 529 34. Jones DR, Schonlau M, Welch WJ. Efficient Global Optimization of Expensive Black-  
530 Box Functions. *J Glob Optim.* 1998;13: 455–492. doi:10.1023/a:1008306431147
- 531 35. Koziel S, Ciaurri DE, Leifsson L. Surrogate-Based Methods. *Computational Optimization,*  
532 *Methods and Algorithms.* 2011. pp. 33–59.
- 533 36. Shahriari B, Swersky K, Wang Z, Adams RP, De Freitas N. Taking the human out of the  
534 loop: A review of Bayesian optimization. *Proc IEEE.* 2016;104: 148–175.  
535 doi:10.1109/JPROC.2015.2494218
- 536 37. Adams SR, Langton FA. Photoperiod and plant growth: A review. *J Hortic Sci Biotechnol.*  
537 2005;80: 2–10. doi:10.1080/14620316.2005.11511882
- 538 38. Yamaura T, Tanaka S, Tabata M. Light-dependent formation of glandular trichomes and  
539 monoterpenes in thyme seedlings. *Phytochemistry.* 1989;28: 741–744. doi:10.1016/0031-  
540 9422(89)80106-2
- 541 39. Skrubis B, Markakis P. The Effect of Photoperiodism on the Growth and the Essential Oil  
542 of *Ocimum basilicum* ( Sweet Basil ). *Econ Bot.* 1976;30: 389–393.
- 543 40. Koza JR. Symbolic Regression-Error-Driven Evolution. *Genetic Programming I: On the*  
544 *Programming of Computers by Means of Natural Selection.* 1992; 237–288.
- 545 41. Rodriguez Rafael GD, Solano Salinas CJ. Empirical study of surrogate models for black  
546 box optimizations obtained using symbolic regression via genetic programming.  
547 *Proceedings of the 13th annual conference companion on Genetic and evolutionary*  
548 *computation.* ACM; 2011. pp. 185–186.
- 549 42. Stijven S, Vladislavleva E, Kordon A, Willem L, Kotanchek ME. Prime-Time: Symbolic  
550 Regression Takes Its Place in the Real World. *Genetic Programming Theory and Practice*  
551 *XIII.* Springer; 2016. pp. 241–260.
- 552 43. Smits GF, Kotanchek M. Pareto-front exploitation in symbolic regression. *Genetic*  
553 *programming theory and practice II.* Springer; 2005. pp. 283–299.

- 554 44. Schmidt M, Lipson H. Distilling free-form natural laws from experimental data. *Science*  
555 (80- ). American Association for the Advancement of Science; 2009;324: 81–85.
- 556 45. Du Q, Faber V, Gunzburger M. Centroidal Voronoi tessellations: Applications and  
557 algorithms. *SIAM Rev. SIAM*; 1999;41: 637–676.
- 558 46. Bergstra JS, Bardenet R, Bengio Y, Kégl B. Algorithms for hyper-parameter optimization.  
559 *Advances in Neural Information Processing Systems*. 2011. pp. 2546–2554.
- 560 47. Snoek J, Larochelle H, Adams RP. Practical bayesian optimization of machine learning  
561 algorithms. *Advances in neural information processing systems*. 2012. pp. 2951–2959.
- 562 48. González J, Dai Z, Hennig P, Lawrence N. Batch bayesian optimization via local  
563 penalization. *Artificial Intelligence and Statistics*. 2016. pp. 648–657.
- 564 49. Srinivas N, Krause A, Kakade SM, Seeger M. Gaussian process optimization in the bandit  
565 setting: No regret and experimental design. *arXiv Prepr arXiv09123995*. 2009;
- 566 50. LeCun Y, Bengio Y, Hinton G. Deep learning. *Nature. Nature Research*; 2015;521: 436–  
567 444.
- 568 51. Snoek J, Rippel O, Adams RP. Scalable Bayesian Optimization Using Deep Neural  
569 Networks. *Int Conf Mach Learn*. 2015; 2171-2180.
- 570 52. Deb K, Myburgh C. Breaking the billion-variable barrier in real-world optimization using  
571 a customized evolutionary algorithm. *Proceedings of the 2016 on Genetic and*  
572 *Evolutionary Computation Conference. ACM*; 2016. pp. 653–660.
- 573 53. Knowles J. ParEGO: A hybrid algorithm with on-line landscape approximation for  
574 expensive multiobjective optimization problems. *IEEE Trans Evol Comput*. 2006;10: 50–  
575 66.
- 576 54. Jin Y. Surrogate-assisted evolutionary computation: Recent advances and future  
577 challenges. *Swarm Evol Comput. Elsevier B.V.*; 2011;1: 61–70.  
578 doi:10.1016/j.swevo.2011.05.001
- 579 55. Lehman J, Miikkulainen R. Neuroevolution. *Scholarpedia*. 2013;8: 30977.
- 580 56. Miikkulainen R, Iscoe N, Shagrin A, Cordell R, Nazari S, Schoolland C, et al.  
581 Conversion rate optimization through evolutionary computation. *Proceedings of the Genetic and*  
582 *Evolutionary Computation Conference. ACM*; 2017. pp. 1193–119  
583

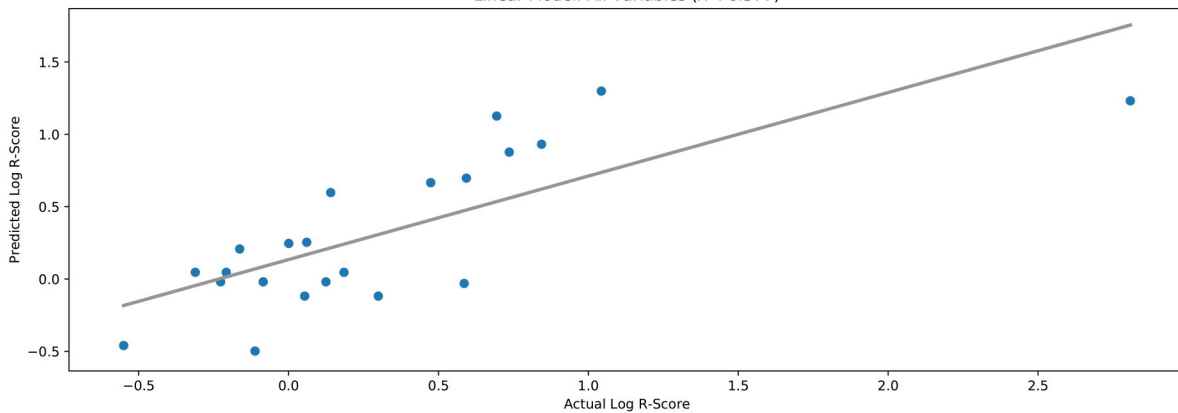






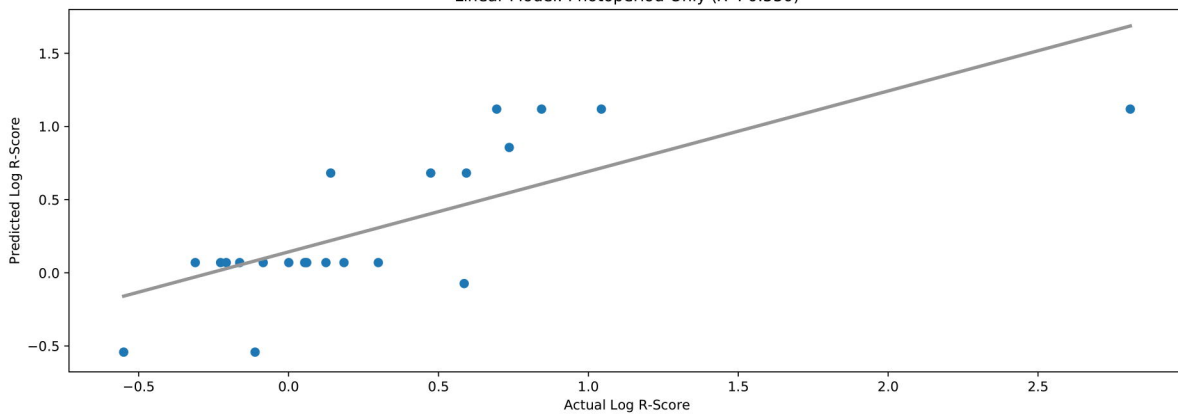
Linear Model: All Variables ( $R^2: 0.577$ )

**a**



Linear Model: Photoperiod Only ( $R^2: 0.550$ )

**b**



Linear Model: Residual ( $R^2: 0.061$ )

**c**

