

1 **Problems with Estimating Anthesis Phenology Parameters in *Zea mays*: Consequences for Combining**  
2 **Ecophysiological Models with Genetics**

3 Abhishes Lamsal<sup>a</sup>, Stephen M. Welch<sup>a</sup>, Jeffrey W. White<sup>b</sup>, Kelly R. Thorp<sup>b</sup>, and Nora Bello<sup>c</sup>

4 <sup>a</sup>Department of Agronomy, Kansas State University, 2104 Throckmorton Plant Science Center, Manhattan, Ks  
5 66502, USA

6 <sup>b</sup>USDA-ARS Arid-Land Agricultural Research Center, 21881 North Cardon Ln, Maricopa, Az 85138, USA

7 <sup>c</sup>Department of Statistics, Kansas State University, 002 Dickens hall, Manhattan, Ks, 66502, USA

8 **Corresponding Author:** Abhishes Lamsal, [abhilam@ksu.edu](mailto:abhilam@ksu.edu)

9

## 10 **Abstract**

11 Ecophysiological crop models encode intra-species behaviors using constant parameters that are presumed to  
12 summarize genotypic properties. Accurate estimation of these parameters is crucial because much recent work  
13 has sought to link them to genotypes. The original goal of this study was to fit the anthesis date component of  
14 the CERES-Maize model to 5266 genetic lines grown at 11 site-years and genetically map the resulting  
15 parameter estimates. Although the resulting estimates had high predictive quality, numerous artifacts emerged  
16 during estimation. The first arose in situations where the model was unable to express the observed data for  
17 many lines, which ended up sharing the same parameter value. In the second (2254 lines), the model  
18 reproduced the data but there were often many parameter sets that did so equally well (equifinality). These  
19 artifacts made genetic mapping impossible, thus, revealing cautionary insights regarding a major current  
20 paradigm for linking process based models to genetics.

## 21 **Highlights**

- 22 • CSM-CERES-Maize v. 4.5 was used to fit the anthesis date parameter for 5266 genetic lines grown at 11  
23 site-years.
- 24 • Despite the high predictive value of the model outputs, numerous artifacts emerged in the estimation  
25 process.
- 26 • The model was unable to express the observed variation in anthesis date data for many lines.
- 27 • More than one parameter set (equifinality) were found for 2254 lines that equally reproduce the data.
- 28 • These results revealed cautionary insights regarding a major current paradigm for linking process based  
29 models to genetics.

## 30 **Keywords**

31 CERES-Maize; Genotype-Specific-Parameters; Parameter estimations; Equifinality; Expressivity; Nested  
32 Association Mapping.

## 33 **1. Introduction**

34 In the opening sentences of the 1968 book, *The Population Bomb*, Paul Ehrlich (and his wife Anne,  
35 uncredited at publisher behest) wrote, “The battle to feed all of humanity is over. In the 1970s hundreds of  
36 millions of people will starve to death in spite of any crash programs embarked upon now” and, in a subsequent  
37 chapter, “I don't see how India could possibly feed two hundred million more people by 1980.” Fortunately,  
38 research started in Mexico, India and elsewhere by Norman Borlaug before 1968 created high yielding dwarf  
39 wheat varieties that, worldwide, are credited with averting one billion deaths from famine. India also  
40 introduced IR8, the so-called “miracle rice” developed at the International Rice Research Institute in the  
41 Philippines and the predicted human catastrophe was averted.

42 Nearly 50 years later, the specter of global disruption is again upon us. The challenges today are not  
43 only increasing human population (which has doubled since 1970) but emerging concerns like climate change  
44 and declining water resources. The confluence of these manifold trends makes finding ways to feed nine billion  
45 people by 2050 one of the most pressing issues of our time (Stone, 2011). However, the annual percentage  
46 increase rates for crop yields are only half those required to meet that goal (Godfray et al., 2010).

47 Beginning over 20 years ago, a paradigm has emerged offering the promise of dramatically accelerating  
48 breeding programs via improved phenotype prediction of prospective crop genotypes in novel, time-varying  
49 environments subject to sophisticated management practices (Cooper et al., 2016; Hammer et al., 2006; Welch  
50 et al., 2005a; White and Hoogenboom, 1996; Yin et al., 2003). The basic notion has two parts. The first is to  
51 exploit ecophysiological crop models (ECM's) to describe the intricate, dynamic, and environmentally  
52 responsive biological mechanisms that determine crop growth and development on daily or even hourly time  
53 scales. The aim is to use highly detailed, nonlinear simulation models to predict the phenotypes of interest  
54 within a subsample of possible environments and in-field management options. ECMs, whose origin is often  
55 credited to Wit. (1965), encode intra-species behavioral differences in terms of parameters that are intended to

56 summarize genotypic properties. On the strength of that presumption, the constants are termed *genotype-*  
57 *specific parameters* (GSP's).

58         The second part of the paradigm is to use quantitative genetic methods such as genomic prediction  
59 (Meuwissen et al., 2001) to relate the GSP's to genotypic markers (Cooper et al., 2016). Next, the outcomes of  
60 crosses are estimated by (1) calculating the GSP values that would arise from possible offspring genotypes.  
61 These values are then (2) used in ecophysiological model runs to predict the phenotypes in the target  
62 population of environments (for which detailed descriptive data must be available). In simplified instances, this  
63 approach has seen remarkable success (e.g., (Reymond et al., 2003).

64         Composed of large coupled sets of continuous-time differential equations, ecophysiological models  
65 simulate many interacting processes (Jones et al., 2003; White and Hoogenboom, 2010) operating in the soil-  
66 plant-atmosphere continuum. These processes include physiology (e.g., photosynthesis, respiration, resource  
67 partitioning to various plant parts, and growth), phenology (leaf emergent timing, the date of vegetative-to-  
68 reproductive development, etc.), as well as chemistry and physics (soil water flows, chemical transformations,  
69 energy fluxes, gas exchange, etc.). During simulation runs, model formulas compute instantaneous process  
70 rates based on plant status and environmental conditions at each time point. These rates are integrated (*sensu*  
71 calculus) to output time series of dozens of plant variables. The models typically have 10 to 20 GSP's whose  
72 estimates are read from input files at the start of model execution. Numerous other inputs (e.g. soil water  
73 holding capacities by layer; measured daily solar radiation, rainfall, maximum and minimum temperatures; etc.)  
74 further quantify the physical environment.

75         The lynchpin of the two-step paradigm is the accurate estimation of the GSP's so that these can be  
76 related to allelic states of the individual lines. Unfortunately, the direct measurement of GSP's is so time- and  
77 resource-demanding as to be infeasible for large numbers of lines. Indirect GSP estimation via model inversion  
78 is also challenging because easily-measured plant phenotypes exhibit strong interactions with the environment  
79 (Chenu et al., 2009) thus increasing data requirements by necessitating trait measurement in multiple settings

80 (Hammer et al., 1987). Even so, ecophysiological crop models enjoy extensive global use in areas ranging from  
81 global climate change, policy analysis, crop management, etc. Indeed, a Google search on the abbreviations of  
82 just two major model systems [namely “DSSAT” (Hoogenboom et al., 2015) and “APSIM” (Keating et al., 2003)]  
83 returned 134,000 hits. Not surprisingly, there is an extensive literature (reviewed briefly below) on  
84 ecophysiological model parameter estimation.

85 Initially, the authors’ intent was to apply the two-step method to anthesis date using data from over  
86 5000 lines comprising the maize nested association mapping population (NAM) (McMullen et al., 2009), which  
87 was developed specifically to enable high-resolution studies of trait genetic architectures. Not only is anthesis  
88 date a phenotype of major biological significance, but it was also studied in this same panel using conventional  
89 statistical genetic methods (Buckler et al., 2009; Hung et al., 2012). Our hypothesis was that applying the  
90 proposed 2-step paradigm would demonstrate its merit in the specific context of the large data sets increasingly  
91 used in crop breeding programs to interrelate genotypes and phenotypes. Contrasting the results of the  
92 standard and ecophysiological approaches was expected to be interesting and informative. Granted, the  
93 model fitting methods to be used were not novel, but we expected that a further demonstration of their value  
94 with data sets much larger than ever used before would have utility.

95 However, something quite different happened. We discovered modeling issues and estimation artifacts  
96 that are of sufficient severity and generality that, if not addressed, are likely to imperil the breeding  
97 acceleration paradigm. Therefore, the objectives of this paper were 1) to describe these problems and the  
98 methods that revealed them (which can be applied as detection tools in studies of other traits) and 2) to discuss  
99 research directions that might ameliorate the problems.

## 100 **2. Background**

101 Numerous optimization methods have been used to estimate parameters for ECM’s. Surprisingly,  
102 perhaps the most common approach has been that of trial and error (Wallach et al., 2001), wherein different  
103 parameters values are manually tested until an acceptable match between simulated and observed data is

104 found. This approach, of course, becomes highly inefficient as the number of model parameter increases. Thus,  
105 numerous off-the-shelf, automated optimization techniques have been developed. Examples include the  
106 simplex method (Grimm et al., 1993), simulated annealing (Mavromatis et al., 2002; Thorp et al., 2008),  
107 sequential search software (GENCALC) (Hunt et al., 2001), Uniform Covering by Probabilistic Region (UCPR)  
108 (Román-Paoli et al., 2000), particle swarm optimization (PSO) (Koduru et al., 2007), and generalized likelihood  
109 uncertainty estimation (GLUE) (He et al., 2010). While these traditional optimization techniques have  
110 advantages, they can be inefficient in terms of runtime and are highly dependent on optimization settings when  
111 thousands of combinations of line  $\times$  planting site-years are involved – a situation that is becoming common in  
112 the era of massive genetic mapping populations. The fundamental issue is that, as the number of lines and  
113 environments increases, estimating GSP's for each line independently usually involves highly redundant  
114 simulation. To this end, we adapted an algorithm pioneered by Welch et al. (2000) and Irmak et al. (2000), as  
115 described in methods section. The approach exhibits particular efficiencies when individual plantings  
116 incorporate large numbers of lines and, serendipitously, supports a close examination of the estimation  
117 process, itself.

118 The vast majority of prior ECM parameter estimation studies have been conducted in non-genetic  
119 contexts. Against these backgrounds, the sole merit criterion has been the predictive skill demonstrated by the  
120 GSP estimates obtained. However, the current setting, however, is markedly different. GSP's are not just inputs  
121 to ecophysiological crop models; GSP's simultaneously function as the outputs (i.e. dependent) variables of  
122 genetic prediction models. As such, GSP's are at least as closely related to tangible biochemical processes at  
123 the molecular level as they are summative of physiological properties (e.g. maximum photosynthetic rates) in  
124 higher organizational realms. Therefore, a deeper inspection of their estimation is warranted and two concepts  
125 are helpful in achieving the enhanced discernment now required.

126 We employ the term “expressivity” (and the adjective “expressive”) to describe a model's innate ability  
127 to reproduce a set of observations independent of particular parameter values. An expressive model may fail

128 to replicate data because an unskilled optimizer cannot find a meritorious combination of parameter values. In  
129 contrast, a model with low expressivity will fail to fully mimic actual data irrespective of what (biologically or  
130 physically reasonable) values are assigned to its parameters. In cases where the latter behavior is detected,  
131 remedies will be vigorously sought. However, as shown below, however, systematic gaps in expressivity can  
132 coexist even within an overall framework of predictively skilled model performance.

133 Another model property that has received little attention in previous estimation studies is equifinality.  
134 Equifinality describes a situation in which multiple sets of parameter values generate identical model  
135 predictions. In statistics, a synonym for “equifinality” is “parameter non-identifiability” (Luo et al., 2009). When  
136 the only concern is prediction quality and that seems “good enough”, it is easy to consider equifinality a non-  
137 problem. However, when parameters are intermediaries rather than just inputs and equifinality exists, it begs  
138 the question as to what relationship, if any, putative GSP estimates might bear to allelic states across the  
139 genotype? A moment’s reflection shows that equifinality and expressivity are different model properties. The  
140 former relates to how many different estimates yield identical predictions; the latter refers to the possible  
141 existence of systematic failures of those predictions to mimic observed data.

142 In this paper, we explore these issues in modeling and estimation using the anthesis phenology  
143 component of the CERES-Maize ECM (Jones et al., 1986; Kiniry and Bonhomme, 1991; Major and Kiniry, 1991)  
144 and observed dates from multiple plantings of three maize genetics panels totaling nearly 5300 lines. Anthesis  
145 initiates the period of grain development and is therefore a critical milestone toward grain yield. As such, it  
146 mediates the adaptation of the crop to its environment by determining the relative length of the vegetative and  
147 reproductive growth phases and is a key target of breeding programs (Buckler et al., 2009). (Although at the  
148 apical meristem, floral initiation precedes the visible morphological change of anthesis, the linkage between the  
149 two is tight enough that we follow common modeling practice and consider them as effectively synonymous.)  
150 The genetics of flowering time has been intensively studied in the model plant *Arabidopsis thaliana* where well  
151 over 100 influential genes are now known (Bratzel and Turck, 2015). Indeed, gene expression models of

152 flowering time of *A. thaliana* based on differential equations have been developed (Valentim et al., 2015), and  
153 genetically-informed approaches have established the relationships between network-level function and  
154 common ecophysiological time formulations (Wilczek et al., 2009). In maize, our understanding of the genetic  
155 control on flowering time is more limited but has been advancing in recent years. More than 30 genes have  
156 been described and conservation of key features from *A. thaliana* seems apparent (Table 1 in (Dong et al.,  
157 2012)). A quantitative gene network model based on a number of these loci has been published (Dong et al.,  
158 2012).

159         The general desire within applied quantitative genetics to probe genetic architectures has led to the  
160 construction of ever-larger and/or special purpose mapping populations (Buckler et al., 2009). The maize NAM  
161 panel (McMullen et al., 2009) was constructed by making bi-parental crosses between one common parent,  
162 B73, and each of a set of 25 other inbreds that collectively encompassed a wide range of maize diversity.  
163 Approximately 200 offspring from each of these 25 crosses were then inbred for a number of generations to  
164 ensure, to the greatest degree feasible, that the influence of each locus on any trait of interest reflected the  
165 contribution of one parent only. Individual plant genotypes produced in this fashion are called “recombinant  
166 inbred lines” (RIL’s). Buckler et al. (2009) reported a seminal study of maize anthesis dates using this NAM  
167 panel. Demonstrating the power of these lines to finely dissect genetic contributions to traits of interest, they  
168 identified 36-39 QTL, where the exact number depended on the analysis method used. Most of loci had small  
169 effects but collectively, they explained 89% of total variation in anthesis date.

170         For the reasons outlined above, accurate prediction of anthesis date is a major target for  
171 ecophysiological crop models (Román-Paoli et al., 2000). However, few studies exist have used large data sets  
172 for ECM calibration. Mavromatis et al. (2002) reported 5,109 site-year-line-parameter combinations and Welch  
173 et al. (2002) estimated 4,620 site-year-line-parameters. The effort presented herein encompassed 197,964 site-  
174 year-line-parameter combinations – to our knowledge, the largest such study ever reported. As the following  
175 sections document, it was the sheer scale of this data set and the resulting scatterplots depicting thousands of



176 lines that revealed worrisome issues of equifinality and expressivity that might be overlooked in studies of  
177 smaller scale.

### 178 **3. Materials and Methods**

#### 179 *3.1 Experimental data*

180 Observations collected on anthesis date for a total of 5266 maize lines were obtained from the Panzea  
181 data repository (<http://www.panzea.org>). The lines used were members of three genetic panels. In particular,  
182 4785 lines were from the 25 RIL panels comprising the maize NAM set described above. Also included were an  
183 additional 200 RIL lines commonly referred to as the IBM panel because they originated by **I**ntermating **B**73 ×  
184 **Mo**17 (Lee et al., 2002). Finally, a maize diversity panel (Flint-Garcia et al., 2005) contributed data on 281  
185 additional lines. Various combinations of these lines were grown at six US sites: New York (NY), North Carolina  
186 (NC), Illinois (IL), Missouri (MO), Florida (FL) and Puerto Rico (PR), during 2006 and 2007 for a total of eleven  
187 site-years. In what follows “NY6” denotes the 2006 planting in New York, respectively by state abbreviation and  
188 year for other site-years. Table 1 gives the exact locations of the experimental sites, and the respective sowing  
189 dates. The “Total Lines” row of the table gives the number of lines from the three panels that were present in  
190 each study. The “Lines with data” row lists the number of lines with available observations on anthesis date.  
191 Data on daily maximum and minimum temperatures for each site were provided by the maize NAM  
192 collaborators (H. Hung, personal communication, 2010) and did not included metadata on position of the  
193 weather stations to the field plots, types and calibration of sensors or types of radiation shields used.

194

195 **Table 1.** Sowing dates, geographical coordinates, total number of lines planted and number of lines for which  
 196 anthesis dates were observed for all site-year combinations used in this study.

	NY6	NY7	NC6	NC7	MO6	MO7	IL6	IL7	FL6	FL7	PR6
Sowing Date (DOY)	128	135	122	120	137	138	128	137	265	280	314
Latitude (deg)	42.73	42.73	35.67	35.67	38.89	38.89	40.08	40.08	25.51	25.51	18.00
Longitude (deg)	-76.66	-76.66	-78.49	-78.49	-92.23	-92.23	-88.2	-88.2	-80.49	-80.49	-66.51
Number of total lines sown	5478	5478	5478	5478	5478	5478	5478	5478	5026	3753	5131
Number of lines with data	4743	5236	5236	5160	3261	2555	5036	5178	4943	3742	4401

### 197 3.2 CERES-Maize model

198 The Crop Estimation through Resource and Environment Synthesis (CERES)-Maize model is one of the  
 199 oldest, most widely used ecophysiological crop models for maize (Quiring and Legates, 2008). We used the  
 200 CERES-Maize version incorporated in CSM (Cropping System Model) 4.5 (Hoogenboom et al., 2015; Jones et al.,  
 201 2003). The CERES-Maize simulation of development toward anthesis is controlled by a set of GSP's and  
 202 environmental inputs (Kiniry and Bonhomme, 1991; Major and Kiniry, 1991). Specifically, the GSP's studied  
 203 herein were thermal time from emergence to juvenile phase (P1), critical photoperiod (P2O), sensitivity to  
 204 photoperiods longer than P2O (P2), and the phyllochron interval (PHINT) as measured in thermal time. The  
 205 duration of Stage 1, the interval from emergence through the end of the juvenile phase, is calculated by

206 accumulating daily thermal time until P1 is reached. Stage 2 follows immediately and lasts until tassel initiation.  
207 Stage 2 lasts a minimum of four days when the photoperiod (including civil twilight) is less than P20. P2  
208 specifies the number of extra days required for every hour by which the photoperiod exceeds P20. The model  
209 continues to accumulate thermal time through Stage 2. The model assumes that (1) there are five embryonic  
210 leaves; (2) two new leaves initiate during each phyllochron interval; and (3) that anthesis date, which  
211 terminates Stage 3, occurs when all leaves present at the end of Stage 2 (i.e., total leaf number, TOLN) are fully  
212 expanded. The date on which this happens is when the ongoing thermal time accumulation reaches  $TOLN \times$   
213 PHINT.

214 Thermal time is calculated from inputs of daily maximum and minimum temperatures. Sowing dates  
215 (Table 1) determined the time series of weather data that control simulated plant growth and development.  
216 The model calculated daily photoperiods from geographic position. Other required model inputs did not affect  
217 predicted anthesis dates and were not considered here. For example, the soil water and nutrient balance  
218 components of the model do not affect simulated anthesis date in the CERES-Maize model and therefore were  
219 not used in this study. The model also requires row spacing and planting depth, which were set to 0.5 m and  
220 2.5 cm, respectively. No tillage, pest, or disease effects were simulated.

### 221 *3.3. Parameter estimation*

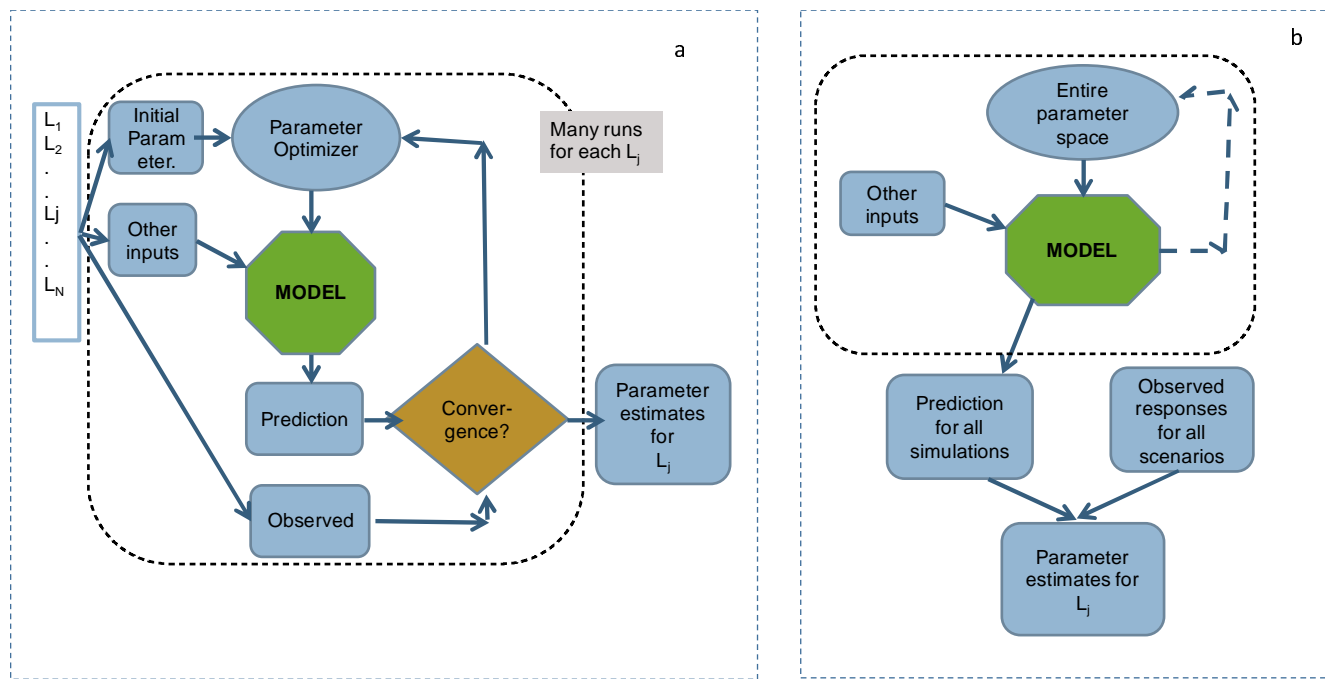
#### 222 *3.3.1 Search strategy*

223 In the conventional approach to parameter estimation (Fig. 1a), an optimizer iterates through a series  
224 of trial solutions for which model predictions are generated in each environment. The entire process is  
225 repeated for each line. This approach becomes inefficient when many lines are planted together in large  
226 experiments and are therefore exposed to identical environments. This is because estimates approaching  
227 optimal goodness-of-fit will only emerge in the latter stages of an iterative optimization run. Therefore, the  
228 majority of early iterations for each line entail the repeated evaluation of estimates with mediocre predictive  
229 ability in the same environment.

230 To overcome this problem, we adapted an approach described by Irmak et al. (2000) and Welch et al.  
231 (2002, 2000). In their scheme (Fig. 1b), model simulations were conducted for each planting across a  
232 multidimensional grid of parameter value combinations. The resulting predictions were stored in a database.  
233 As a second step, for each line the root mean square error objective function (RMSE; Gill et al., 1981) between  
234 observed and predicted anthesis day of year was evaluated with respect to all combinations of parameter  
235 values across all site-years. That is, for line  $l$ ,

$$236 \quad RMSE_l = \sqrt{\frac{1}{n} \sum_{i=1}^n (Y_p - Y_o)^2} \quad (1)$$

237 where,  $n$  is the number of observations for that line (consisting of one observation per site-year  
238 combination), and  $Y_p$  ( $Y_o$ ) is the predicted (observed) anthesis date. The optimizer goal was to minimize the  
239 RMSE for each line. If a unique minimum existed, it defined the combination of GSP values that best fit each  
240 line. Total computational time was reduced because time-consuming model simulations for each combination  
241 of GSP parameter values were only performed once, but those outputs were reused many times in the much  
242 faster RMSE calculations. Another benefit is that a combination of GSP values that yielded poor predictability  
243 for one variety might perform better for a different line. Additionally, this process ensured that identical  
244 parameter combinations were tested for each line, which can aid in comparing the results achieved. Finally,  
245 simply by retabulating the database, any number of different optimizations could be performed using different  
246 observations, alternative subsets of site-years plantings or combinations of parameter values. The use of  
247 alternative objective functions is also possible without requiring additional simulations. Because of the central  
248 role played by the database of simulation outputs, we will refer to this scheme as the *database method*.



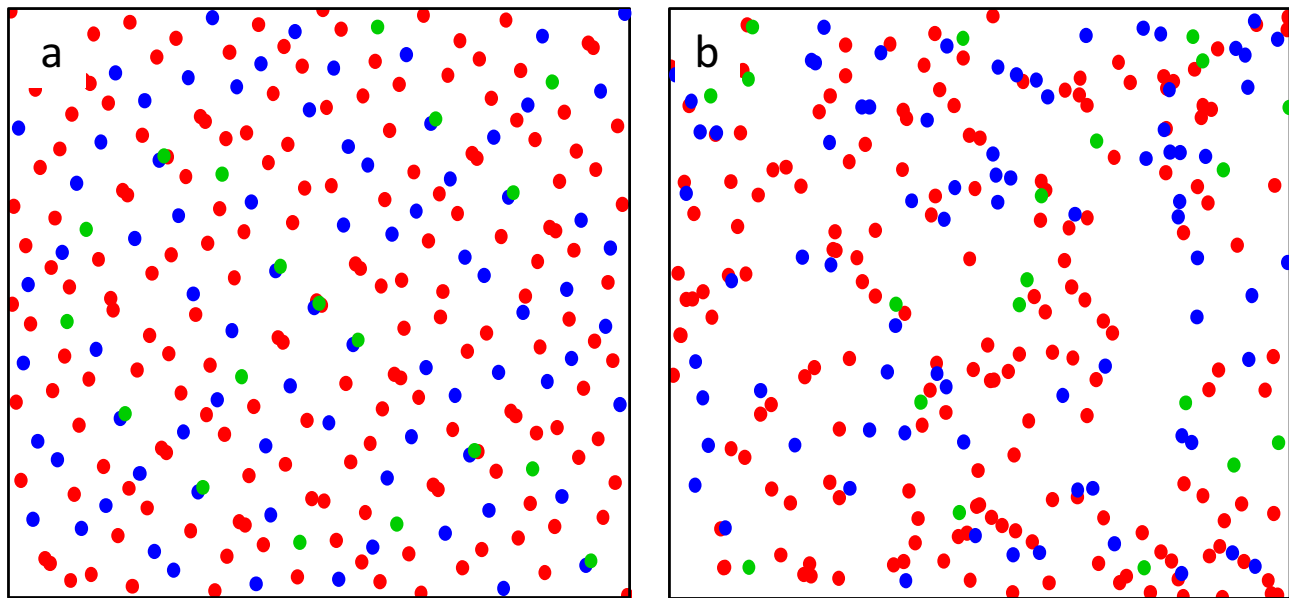
249

250 **Fig. 1.** Parameter search strategies a. Conventional method b. Database method.  $L_{1...N}$  is the number of lines.

251 *3.3.2 Sampling the model parameter space with sobol sequences*

252 Unlike Irmak et al. (2000) and Welch et al. (2002, 2000) who sampled the parameter space with a  
 253 rectilinear grid, we employed Sobol sequences so as to avoid the combinatorial explosion in computational  
 254 requirements that accompany increasing dimensionality. Sobol sequences belong to a family of quasi-random  
 255 processes designed to generate samples of multiple parameters dispersed as uniformly as possible over the  
 256 multi-dimensional parameter space (Press et al., 1992). Sobol sequences are specifically designed to generate  
 257 samples with low discrepancy – that is, a minimal deviation from equal spacing. Unlike random numbers, quasi-  
 258 random algorithms can effectively identify the position of previously sampled points and fill the gaps between  
 259 them (Saltelli et al., 2010), thus avoiding the formation of clusters. Further, Sobol sequences offer reduced  
 260 spatial variation compared to other sampling methods (e.g., random, stratified, Latin hypercube; see Fig. 2a vs.  
 261 2b), make this method more robust (Burhenne et al., 2011). We used a Python-based algorithm to generate a  
 262 Sobol sequence of quasi-random numbers for calculating 32,400,070 sets of the four CERES-Maize GSP's,  
 263 leading to a uniformly-sampled four-dimensional parameter space for P1, P2, P2O, and PHINT. To construct the

264 database, CERES-Maize calculated anthesis date for each GSP combination in each of the 11 site-years – a total  
265 of 356,400,770 model runs. Table 2 describes the upper and lower bounds and the number of distinct values  
266 obtained for each parameter.



267

268 **Fig. 2.** (a) The first 275 quasi-random points from a two-dimensional Sobol sequence. (b) The first 275  
269 points produced by the commonly used Mersenne twister pseudo-random number generator  
270 (Matsumoto and Nishimura, 1998). The Sobol sequence covers the space more evenly. The first 20  
271 points are green, the next 80 are blue, and the final 175 are red, thus demonstrating Sobol gap filling.

272

273 **Table 2.** Parameter ranges used in generating sobol sequence.

Parameter	Definition	Unit	Min	Max	No. of unique values
P1	Thermal time from seedling emergence to end of juvenile phase	GDD (°C)	150	450	30,001
P20	Critical photoperiod hour	hrs.	10	14	401
P2	Days of anthesis date delay for each hour by which the day length exceeds P20	rate	0	2	20,001
PHINT	Phylochron interval (Interval between successive leaf tip appearances)	GDD (°C)	25	70	45001

274

### 275 3.3.3 High performance computing

276 The number of model runs was too large for lab-scale computing facilities, so we used the “Stampede”  
277 supercomputer at the Texas Advanced Computing Center (TACC) (Burhenne et al., 2011). *In toto*, the CERES-  
278 Maize runs required 63,372 CPU-hours, which equates to ca. 176 simulations per second distributed across 112  
279 processors. The predicted anthesis dates were collated and transferred to the “BeoCat” computing cluster at  
280 Kansas State University ([https://support.beocat.ksu.edu/BeocatDocs/index.php/Compute\\_Nodes](https://support.beocat.ksu.edu/BeocatDocs/index.php/Compute_Nodes)). There,  
281 RMSE values were tabulated for each line × parameter value combination across all site-years in which anthesis  
282 date was observed. As combinations of GSP values were found that had progressively lower RMSE values, they  
283 were recorded by the computer. This process required ca. 15 minutes of wall clock time per line so the total  
284 estimation process was completed in ca. 7 h on 200 Xeon E5-2690 cores.

### 285 3.4 Assessing estimate properties

#### 286 3.4.1. Equifinality

287 Equifinality occurs when multiple combinations of parameter estimates generate the same minimal  
288 RMSE value, often because they generate identical model predictions (Luo et al., 2009), in this case identical  
289 integer DOY values for anthesis dates. We quantified "equifinality" by defining "number of ties" as the number  
290 of Sobol sets of parameter combinations that produced the same optimal RMSE values, minus one. No  
291 equifinality is present in a line if there is only one combination of parameter values that minimizes the RMSE.  
292 That is, there are zero ties among its estimates. To illustrate the magnitude of the problem and our motivation  
293 to study it more closely, we note that 2254 (43%) of the 5266 lines available in the data exhibited equifinality.  
294 The worst case was represented by a line that had 1,043,933 distinct combinations of GSP values that produced  
295 identical anthesis date predictions, and thus the same RMSE, thereby yielding 1,043,932 ties.

296 During the database tabulation phase, the values of the "best combination of parameter estimates seen  
297 so far" was updated only if its RMSE value was strictly better than all previously evaluated ones. So, when  
298 equifinality was present, the final GSP estimate was the first combination of parameter values encountered that  
299 had a minimal RMSE value. As a result, some of the analyses described below are sensitive to equifinality,  
300 illustrating the fact that subtle optimizer algorithm idiosyncrasies can have marked impacts on the overall  
301 results. Such cases are noted explicitly along with the procedures used to mitigate the effects.

#### 302 *3.4.2. Interrelationships between parameter estimates*

303 Correlations and other relations among parameter estimates are highly important to breeding programs  
304 and related simulation studies. When correlations between parameter estimates are present, opportunities  
305 exist to select on one plant trait by selecting on a related phenotype instead. Additionally, there have been a  
306 number of *in silico* studies where CERES models were used to design crop ideotypes (Laurila et al., 2012;  
307 Semenov and Stratonovitch, 2013). Such efforts find combinations of model parameter values that predict  
308 phenotypes well suited to the target population of environments. Once identified, lines with those values  
309 become breeding targets. However, a potential pitfall arises if realizing the desired genotype involves changing  
310 parameter values in directions contrary to the correlations that exist between them.



311 For this reason, we explored the pairwise correlation structure of the GSP parameter estimates and  
312 generated pairwise scatter plots of their line-specific values. However, the latter revealed a bizarre pattern, the  
313 diagnosis of which ultimately led us to the second problem alluded to in the introduction – the inability of the  
314 model to reproduce certain observational combinations – and to the methods presented next.

#### 315 *3.4.3. Model expressivity*

316 A common graphical method to assess the quality of model fit is to plot the predicted vs. observed  
317 values (e.g., Fig. 3). Such scatterplots can be informative in detecting areas of mismatch between observed and  
318 predicted values, thus providing specific characterization of the model’s lack of fit. By definition, each point in  
319 the scatterplot corresponds to a prediction that a model is able to make given an optimized set of parameter  
320 values. However, an entirely different question is whether there are observations that a given model cannot  
321 reproduce using *any* reasonable combination of parameter values? That is, one might seek to assess whether a  
322 given model has the requisite expressivity to reproduce the data.

323 The database approach allows such a question to be addressed using what we term *phenotype space*  
324 scatter plots. In such plots, each axis corresponds to a different site-year. The coordinates along the axes  
325 represent the observed or predicted anthesis dates for each site-year. Model expressivity is then assessed by  
326 comparing the scatter of predicted anthesis date generated from a wide range of GSP value combinations to  
327 the scatter of observed values in large data sets. Because equifinality does not affect predictions, this method  
328 of evaluating model expressivity is independent of the order in which an optimizer locates points that minimize  
329 RMSE values (see the second paragraph in section 3.4.1).

#### 330 *3.4.4 Testing for parameter stability across environments*

331 In order for the two-step paradigm outlined in the Introduction to work, the estimates of GSP’s should  
332 not vary across the set of environments used to estimate them, a property called “stability” (Hammer et al.,  
333 2006). If GSP estimates did vary across environments, there would be no way to tell what GSP values to input  
334 to the ecophysiological model to predict traits whenever daily weather time series or soils differed from those

335 used in the paradigm's first step. This might seem an insuperable barrier to readers for whom G×E interactions  
336 are virtually ubiquitous among quantitative plant phenotypes, but it is not. This is because the *raison d'être* of  
337 models like CERES-Maize is to explain crop variety × environment interactions mechanistically based on  
338 physiological principles.

339 Many GSP's, including the ones in this study, explicitly relate plant behaviors (e.g., development toward  
340 anthesis) to environmental variables (e.g., temperature and photoperiod in the current case). Modelers assert  
341 that GSP's are properties of the individual lines (i.e., stable) and, therefore, by implication, have a genetic basis  
342 because genotypes do not change with the environment. Over time, it is thus expected that research will  
343 mechanistically link at least some GSP's to molecular genetic processes. For example, both short (P20) and long  
344 day critical photoperiods are determined by the dynamics of the CONSTANS protein in a range of plants  
345 including *Arabidopsis* (Andrés and Coupland, 2012) and a number of grasses (Hammer et al., 2006), albeit not  
346 maize (Mascheretti et al., 2015). In rice (*Oryza sativa*), critical short day length has even been successfully  
347 predicted from a differential equation model of the diurnal expression patterns of the *CONSTANS* ortholog  
348 (Welch et al., 2005b).

349 Because stability is both important and reasonable to expect given the goals of ecophysiological  
350 modeling, it has been argued (Welch et al., 2005a) that finding a putative GSP to be unstable is *prima facie*  
351 evidence of a problem. Possible causes of instability include: (1) the model incompletely or incorrectly  
352 disentangles G × E; (2) a stable answer exists but the optimizer is insufficiently skilled to find it; (3) undiscovered  
353 equifinality is present, and the solutions found depend on low-level algorithmic idiosyncrasies of the optimizer  
354 (e.g. section 3.4.1); and (4) unique best GSP estimates exist that the optimizer can find, but because the model  
355 is over-parameterized, the values obtained reflect noise signals that differ between environments.

356 All sources of instability, whether these or others, are detrimental to the two-step ecophysiological  
357 genetic approach to phenotype prediction. Thus, it is critical to know when parameter instability is present, so  
358 herein we developed a statistical approach to detect and test for it. The specific question asked was "Do the

359 GSP estimates depend on the particular set of environments used to construct them?" A conceptually simple  
360 way to answer this might be to (1) obtain a combination of parameter estimates from one subset of site-years,  
361 (2) repeat the estimation with a different subset, and (3) test whether the two sets of parameter estimates  
362 differ according to an appropriate statistical test.

363 A more general and robust approach, however, might be to obtain parameter estimates from many  
364 site-year subsets chosen according to a principled method. Preliminary tabulations of the Sobol database  
365 revealed that equifinality increased dramatically when fewer than seven site-years were used for estimation  
366 (see Results). Therefore, the subset size was set to seven site-years. One method for selection of site-year  
367 subsets might be to resample site-years with replacement. However, as shown by analogy in Fig. 2b,  
368 randomization adds a source of variability to the results that could be of concern given that sampling by  
369 replacement would have  $P_7^{11} = 39,916,800$  possible site-year subsets. Therefore, analogous to Fig. 2b, we used  
370 a combinatorics-based sampling pattern leading to more uniformly-distributed site-year subsets by taking all  
371 combinations of 11 site-years 7 at a time, of which there are  $C_7^{11} = 330$  possibilities. To maximize the amount  
372 of data available for each line in any subset, we focused on the 539 lines for which observation were available in  
373 all 11 site-years.

374 We then conducted 177,870 four-dimensional optimizations to obtain GSP parameter estimates for  
375 each of the 539 line  $\times$  330 site-year set combinations. These optimizations involved only Sobol database  
376 retabulations rather than new model runs, again illustrating the computational efficiency of the database  
377 approach. When forced to generate a single result, the database search returned the combination of GSP  
378 estimates yielding a minimal RMSE that it happened to encounter first. To focus on the subset that lacked this  
379 element of optimizer arbitrariness, we first dropped the 114,314 line  $\times$  site-year combinations that had ties (i.e.  
380 more than one set of GSP estimates yielding the same RMSE). Because our primary interest was in the  
381 variability that different site-year combinations might contribute to GSP estimates, we further restricted our  
382 attention to the 297 site-year subsets that had at least 100 lines remaining after ties were removed. Each of

383 the 539 lines was present in at least 28 site-year subsets, which was deemed adequate for GSP estimation.  
384 These actions left a total of 60,834 estimates for each of the four GSP's in the study. This became our base  
385 group for analysis. We acknowledge that the estimates dropped share a common property (i.e., ties) that might  
386 have systematic effects influencing the results. So, in addition to the base group just described above, we also  
387 examined the set of (1) all 177,870 GSP sets and (2) the 114,314 results for which ties existed. In both cases we  
388 used the optimizer-selected values

389 We then specified a statistical model to test for stability in parameter estimates across environmental  
390 subsets, as follows:

$$391 \quad \rho_{l,e} = \mu_{\rho} + \alpha_l + \beta_e + \varepsilon_{l,e} \quad (2)$$

392 where  $\rho_{l,e}$  represents an estimate of the GSP  $\rho$  (i.e. either P1, P2, P2O, or PHINT) for the  $l^{\text{th}}$  line ( $l =$   
393  $1, 2, \dots, 539$ ) obtained from the  $e^{\text{th}}$  site-year set ( $e = 1, 2, \dots, 297$ ),  $\mu$  is the intercept parameter, acting as an overall  
394 mean of GSP  $\rho$  across all lines and site-year subsets;  $\alpha_l$  is the differential random effect of line  $l$ , assumed to  
395 be distributed  $\alpha_l \sim N(0, \sigma_l^2)$ ;  $\beta_e$  is the differential random effect of the  $e^{\text{th}}$  set of site-years, assumed to be  
396 distributed  $\beta_e \sim N(0, \sigma_e^2)$ ; and  $\varepsilon_{l,e}$  is the left-over residual unique to the  $l, e^{\text{th}}$  observed GSP estimate and  
397 assumed  $\varepsilon_{l,e} \sim NIID(0, \sigma_{\varepsilon}^2)$ . The differential line effects  $\alpha_l$  are considered to be random, as is common in  
398 field studies of plant population biology. Further, the differential effects of site-year sets,  $\beta_e$ , were treated as  
399 random because the corresponding environmental sets are combinations of 7 out of 11 plantings considered to  
400 be a representative, if not random, sample of the population of possible site-years to which we are interested in  
401 inferring.

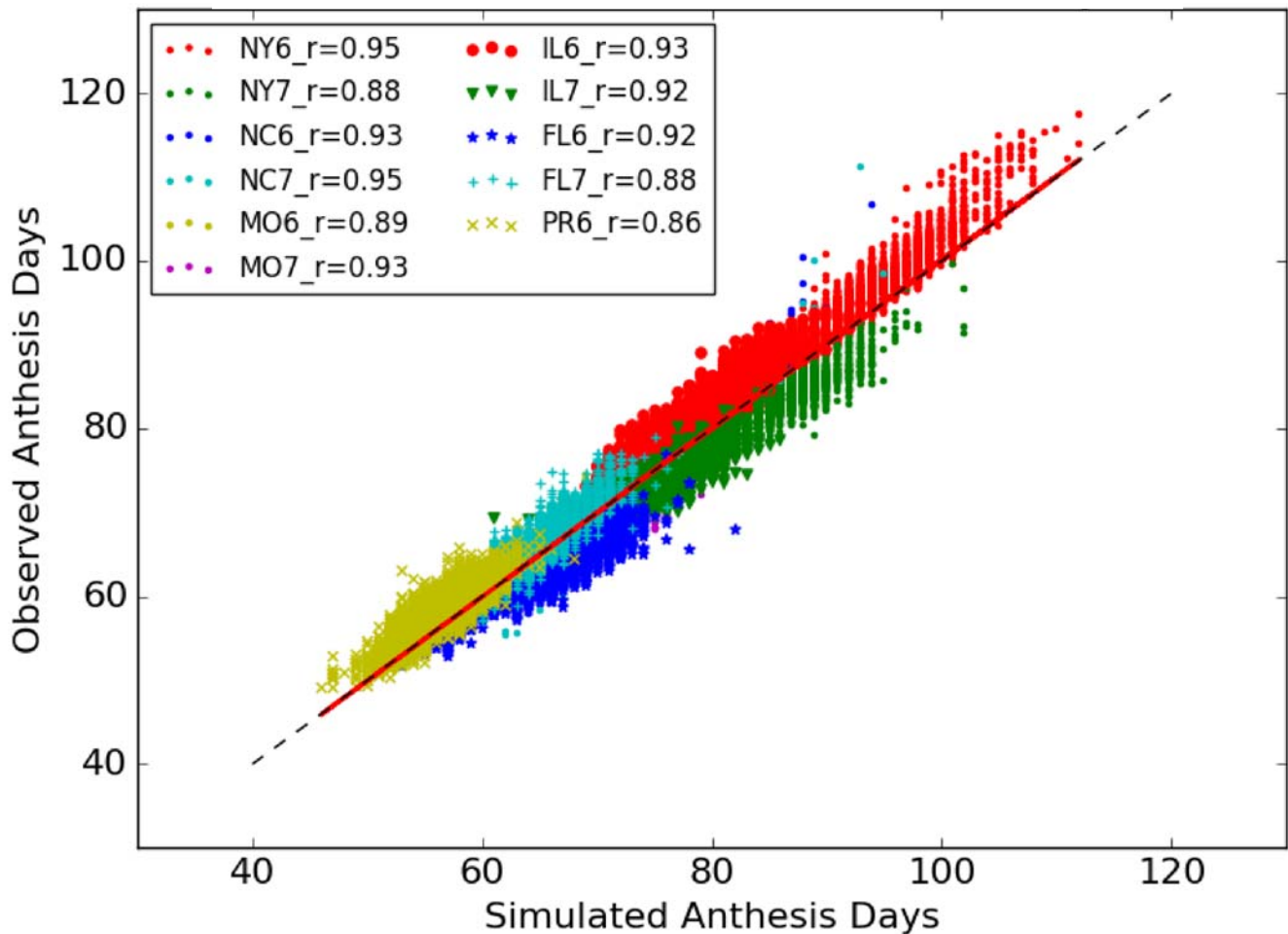
402 If the estimation of any GSP parameter  $\rho$  were stable across the site-year subsets, one would expect  
403 the variance of  $\beta_e$ , namely  $\sigma_e^2$ , to be zero; alternatively, if estimation is unstable, one would expect  $\sigma_e^2 > 0$ .

404 To test this hypothesis set, we fit two competing versions of the statistical model in equation (1), one with and  
405 one without the random effect of site-year subsets  $\beta_e$  for each of the GSP's  $\rho = P1, P2, P2O, \text{ and } PHINT$ .  
406 For each GSP, we then compared the two competing models using a likelihood ratio test statistic against a  
407 central chi-square distribution with half a degree of freedom to account for the fact that the test is being  
408 conducted on the boundary of the parameter space. Statistical models were fitted using the liner mixed-  
409 effects model package lmer in R (Bates et al., 2014) with optimization based on the log-likelihood option. The  
410 lmer package also calculated the Akaike and Bayesian Information Criteria [AIC (Akaike, 1973) and BIC (Schwarz,  
411 1978), respectively], which allowed for an additional assessment of fit for statistical models that included or  
412 excluded the random effects of site-year subsets.

## 413 **4. Results**

### 414 *4.1 Observations vs. Predictions*

415 Fig. 3 shows a color-coded scatterplot of observed vs. predicted days to anthesis for 49,491 line  $\times$  site-  
416 year combinations; the cloud of points is concentrated along the identity line, therefore suggesting accurate  
417 prediction; the overall estimated RMSE is 2.39 days. Also, there seem to be considerable differences between  
418 sites on anthesis days, whereby Florida and Puerto-Rico show very short vegetative durations (ca. 50 d), which  
419 are more than doubled in New York (120 d). Empirical correlation coefficients ( $\hat{r}$ ) were high across site-years  
420 and ranged from 0.86 to 0.95, thus indicating an overall responsiveness across lines to the range of site-year  
421 conditions on anthesis dates. The standard deviations of the predicted values and their corresponding  
422 observations are 10.336 and 10.639, respectively, which, with the overall empirical correlation coefficient of  
423 0.974, account for a close to 1-to-1 estimated regression slope of observations vs. predictions [i.e.  $1.002 =$   
424  $10.639 / 10.336) * 0.974$ ], as per the established statistical identity between these four sample quantities  
425 (Harrison and Tamaschke, 1984).



426

427 **Fig. 3.** Predicted and Observed anthesis days of all 5,266 lines from 11 site-year combinations. The graph has  
428 49,491 points and an overall RMSE of 2.39 days.

#### 429 4.2 Equifinality

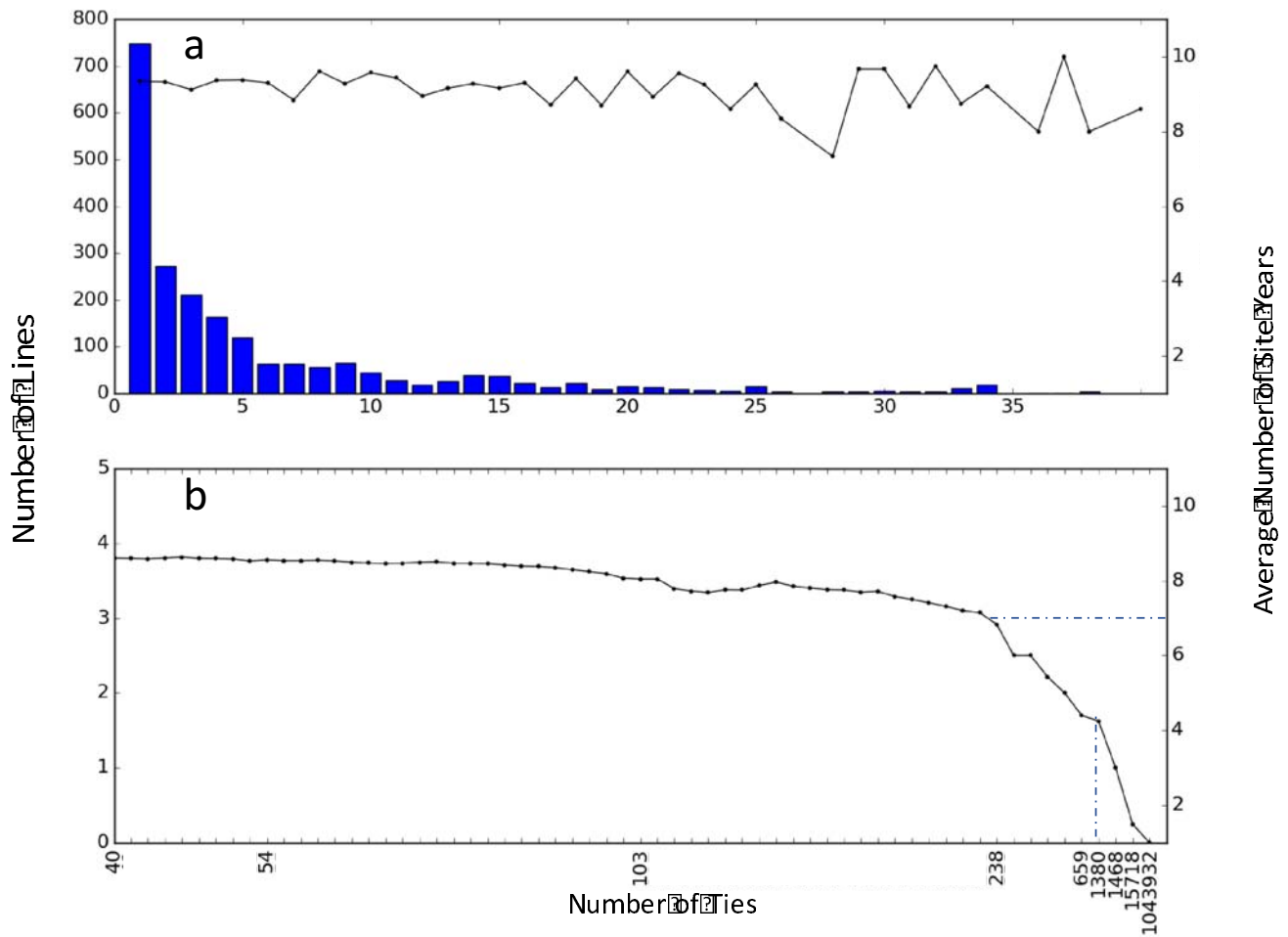
430 A more complex picture emerges when the prevalence of equifinality is considered. As noted in 3.4.1,  
431 for the 2,254 lines exhibiting equifinality, the number of ties can exceed 1M. The histogram in Fig. 4a tabulates  
432 the frequency of ties across lines. There are 2,153 lines with fewer than or equal to 40 ties. The line trace along  
433 the upper portion of the top and bottom panels shows the average number of site-years in each bin.

434 In Fig. 4a, the empirical distribution of ties was right skewed, thereby indicating that a relatively large  
435 number of maize lines had few ties and thus low levels of equifinality. This is particularly true when parameter  
436 estimates were computed using data from 7 to 11 site-years (right axis of Fig. 4b). Further, the distribution of

437 ties appears to have a very long tail to the right, whereby the number of lines with increasing amounts of  
438 equifinality declines very slowly while the number of site-year combinations used for estimation seems to  
439 plateau (Fig. 4a). This pattern continues into Fig. 4b, which shows the 101 lines with more than 40 ties. (No bars  
440 are shown in Fig. 4b due to scale of the y-axis, as each bin generally contains one to three lines.) Interestingly,  
441 the number of ties, and thus equifinality, seems to increase precipitously for the 56 out of 5,266 lines that have  
442 fewer than seven site-years of data (Fig. 4b).

443         As the number of ties increases, one can expect that the range of indistinguishable estimates for any  
444 GSP will widen. To illustrate this phenomenon, a set of GSP estimates were obtained using just two illustrative  
445 site-years (NY6 and NY7) so as to artificially inflate equifinality. Fig. 5 shows scatterplots of coordinate pairs of  
446 either predicted (a) or observed (b) values for anthesis days from NY6 (horizontal axes) and NY7 (vertical axes).  
447 Points in each scatterplot are color-coded to represent the number (on a  $\log_{10}$  scale) of tied GSP combinations.  
448 Each tied GSP combination, when simulated using the weather data for NY6 and NY7, predicts the same  
449 anthesis dates that form the point's coordinates. Dark red indicates 235,976 ties and blue indicates 1 tie. It is  
450 reasonable to expect that as the number of ties increases, the range (max minus min) of the equifinal estimates  
451 will increase. The size of each circle indicates the range of tied P1 estimates expressed as a percentage of the  
452 mean. These percentages extend from 0.36% to 65.68%. The association of redder colors with larger circles  
453 indicates that estimate ranges do, indeed, increase with the level of equifinality.

454         This is an example of a phenotype space plot that can be used to show how properties of interest (e.g.  
455 number of ties and estimate ranges in this case) are distributed across the range of predictions made by the  
456 model given the weather in a pair of site-years. Notice that (1) the cloud of observed points (Fig. 5b) is more  
457 dispersed than that of the predicted points (Fig. 5a), suggesting that model responses to the environment were  
458 less plastic than those of real plants and (2), as indicated by the red lines, the lowest numbers of ties in Fig. 5b  
459 (blue points) appear to fall in empty regions of Fig. 5a where predictions are lacking. This pattern has important  
460 consequences to be explained later in section 4.4.

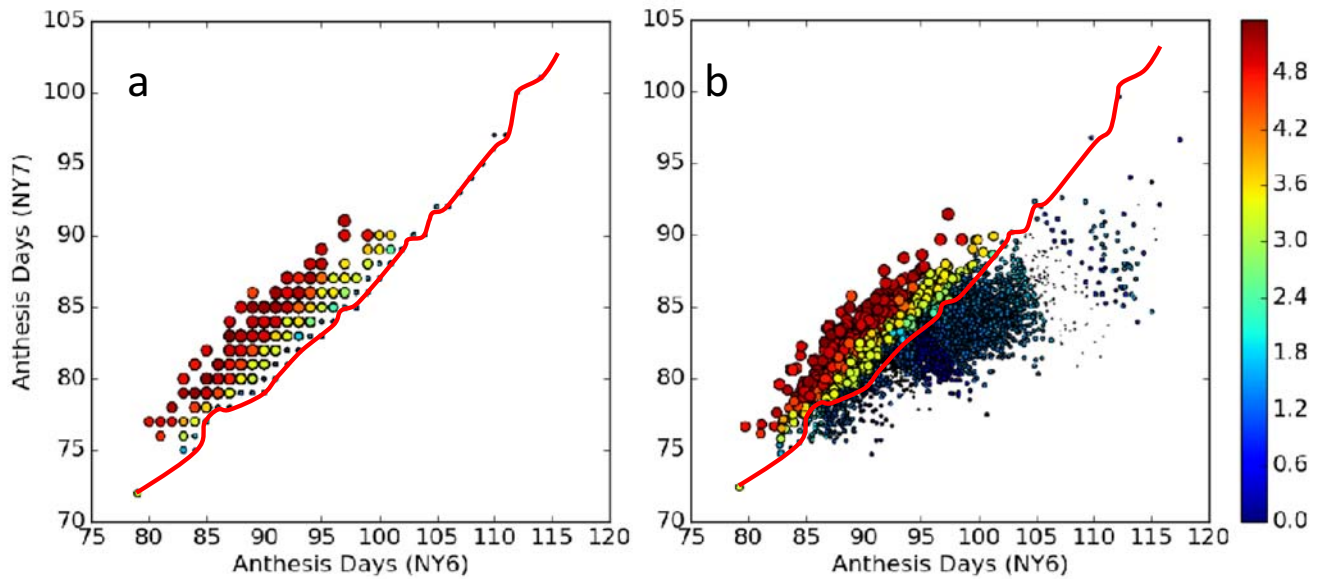


461

462 **Fig. 4.** Histogram depicting the frequency distribution of number of ties for 2,254 lines, used here to  
463 characterize equifinality. (a): Histogram of number of ties for 2153 lines with fewer than or equal to 40 ties. (b):  
464 Continuation of the histogram tail from the upper panel figure representing frequency of ties for the 101 lines  
465 with more than 40 ties. The trace at the top of each panel represents the average number of site-year  
466 combinations (right axis) used as data for parameter estimation.

467



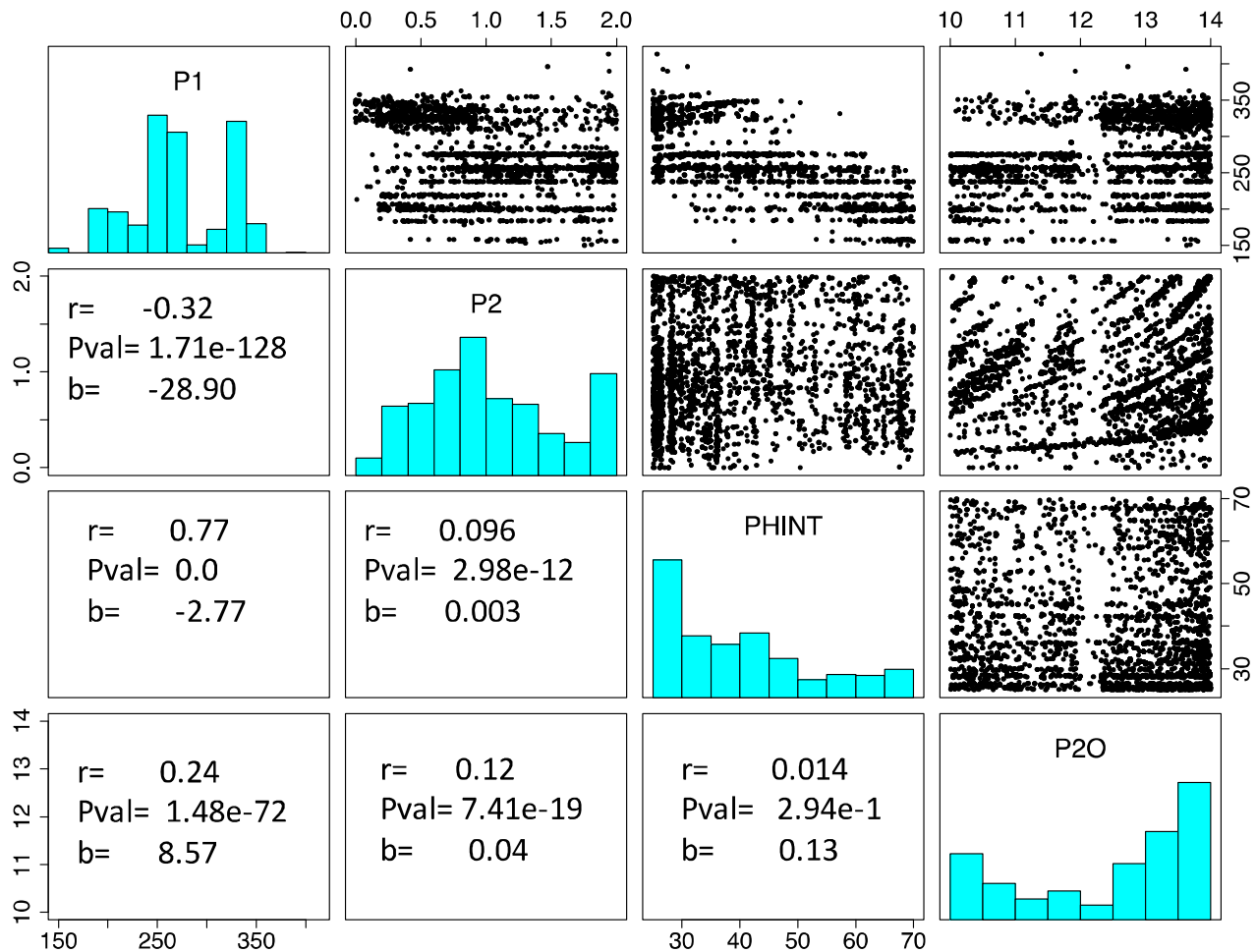


469 **Fig. 5.** Phenotype space plots of predicted (a) and observed (b) values of anthesis dates for site-years NY6 and  
470 NY7. The marker sizes and colors respectively express the levels of equifinality based on number of ties for P1  
471 ( $\log_{10}$  scale) and the relative ranges of its tied values. The red line is explained in the text.

#### 472 *4.3 Interrelationships between parameter estimates*

473 Fig. 6 presents a combined plot depicting histograms of GSP parameter estimates based on all 5,266  
474 lines along the main diagonal and corresponding pairwise GSP scatterplots in the upper right panels. The GSP  
475 estimates were obtained using all site-years. The lower left panels in Fig. 6 show the estimated Pearson  
476 correlation coefficients ( $\hat{r}$ ), estimated regression slopes ( $\hat{b}$ ), and corresponding  $p$ -values for each mirrored  
477 scatterplot. Two immediately apparent features on the scatterplots are to be noted, which might readily escape  
478 notice in data sets with fewer lines. The first is the pronounced banding pattern appearing in all plots except,  
479 perhaps, P20 vs. PHINT. Most bands seem to be linear except for those on the scatterplot of P20 and P2 plot,  
480 which exhibits curvilinearity. The second is the pronounced vertical gap in all P20 scatterplots. In an attempt to  
481 understand the reasons for such patterns, the authors explored multiple seemingly plausible hypotheses,  
482 ranging from genetics to input file coding quirks (e.g., unintended rounding of parameter values) and many

483 more, all of which were tested and discarded. Ultimately, the results presented in the following sections  
 484 provided the explanations.

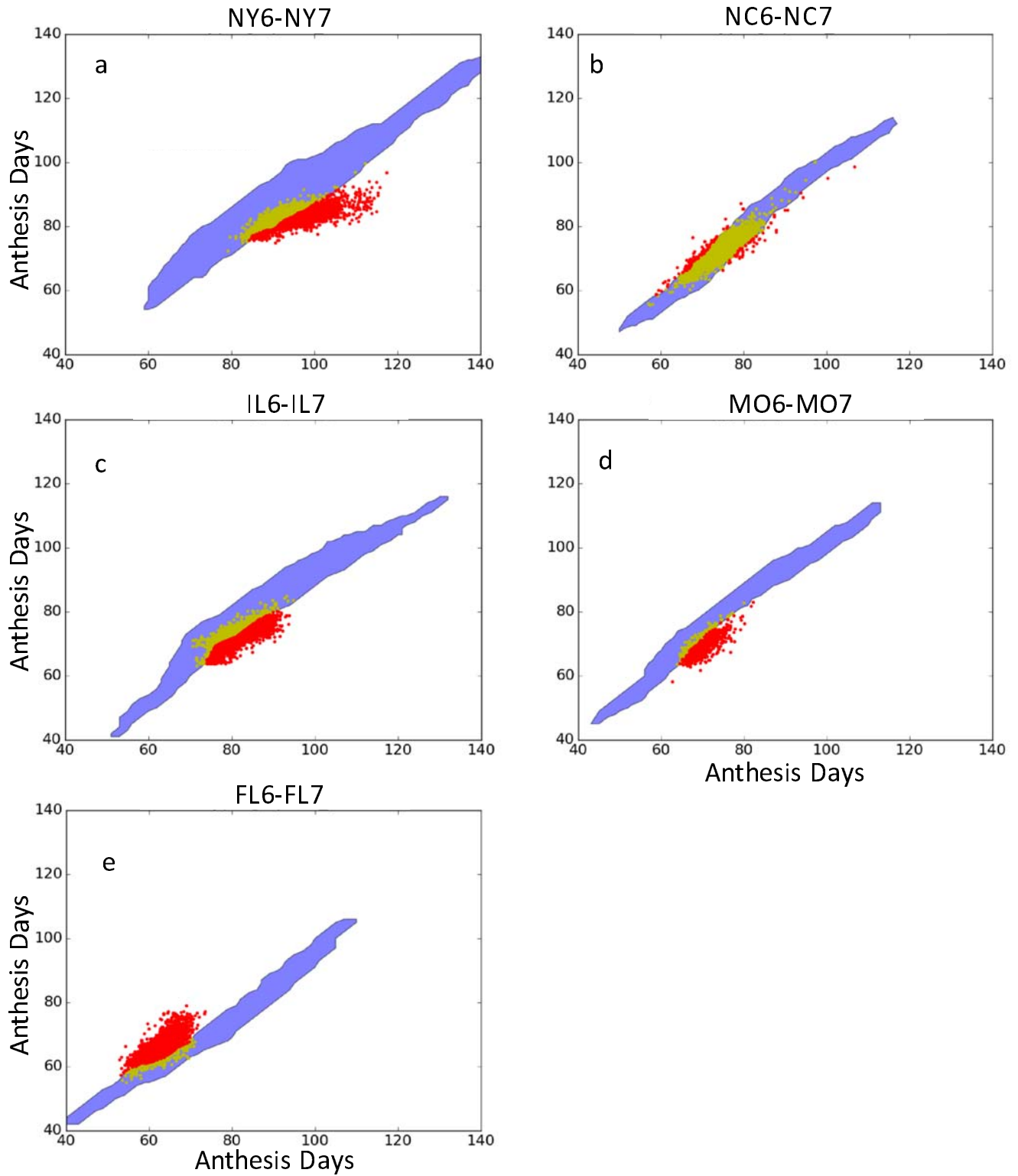


485  
 486 **Fig. 6.** Empirical distribution of selected GSP parameter estimates (main diagonal), pairwise scatterplots (upper  
 487 right triangle) and empirical estimates of Pearson correlation coefficients, regression coefficients and p-values  
 488 (Lower left triangle). Each dot in the scatter plots represents a pair of GSP estimates from a single line.

#### 489 4.4 Model expressivity

490 The first clue to the cause of the banding pattern emerges from the phenotype space plots in Fig. 7.  
 491 Each plot corresponds to an independent fit to just one particular pair of site years. The blue regions in each  
 492 panel of Fig. 7 outline predicted anthesis date pairs for two consecutive years in a given site, where model

493 prediction are constrained by the bounds imposed on the range of values allowed for each of the four GSP's  
494 (Table 2). Also, for each panel in Fig. 7, a dot depicts an observed anthesis date pair for a line present in a given  
495 site in both 2006 and 2007. Yellow (red) dots represent observed anthesis date pairs that the model was able  
496 (unable) to reproduce. We characterize each observation corresponding to a yellow (red) dot as “expressible”  
497 (“inexpressible”). Except for the two North Carolina site-years, there were many lines (Table 3) for which  
498 observations on anthesis date could not be predicted despite: (1) the seeming breadth of GSP values allowed by  
499 Table 2; and (2) the fact that the model was only being asked to match two data points, which would seem to  
500 greatly relax the constraints on GSP estimates.



501

502 **Fig. 7.** Phenotype space plots for predicted and observed anthesis dates. Each panel corresponds to a pair of  
 503 site-years for which fits were done. Regional color codes are described in the text.

504 **Table 3.** Numbers of model expressible and inexpressible observations for selected site-year pairs.

Lines that are <sup>a</sup> :	NY6/NY7	NC6/NC7	IL6/IL7	MO6/MO7	FL6/FL7
Expressible	2189	4964	2024	146	193
Inexpressible	2542	168	2946	637	3339

505 <sup>a</sup> These numbers refer to lines with data in both years of each pair and therefore do not precisely align with Table 1.

506 This begs the question as to what would happen to model expressivity if an even broader range of GSP  
 507 values were allowed. In an attempt to investigate in a computationally efficient way how the outputs of a more  
 508 conventional optimizer might appear when viewed in phenotype space, the CERES-Maize anthesis date routine  
 509 was ported to Python and fit to NY6/NY7 via Differential Evolution (DE) (Das and Suganthan, 2011). DE is a well-  
 510 established (63K Google Scholar hits on “Differential Evolution” as of October 21, 2016) and highly effective  
 511 evolutionary algorithm that embodies mechanisms reminiscent of techniques ranging from the Nelder-Mead  
 512 Simplex (Nelder and Mead, 1965) method to Particle Swarm Optimization (Kennedy, 2011). Among the  
 513 algorithm’s initiating inputs is the range of parameter values within which to search, which were set as shown in  
 514 Table 4. These ranges are greatly broadened from that used in the database search (Table 2); in fact, the values  
 515 in Table 4 are intentionally broader than biological experience would suggest as reasonable.

516 **Table 4.** Extended range of parameter values used for DE search.

Parameter	Definition	Unit	Min	Max	Percent of Sobol Range
P1	Thermal time from seedling emergence to end of juvenile phase	GDD (°C)	75	600	175%

P20	Critical photoperiod hour	hrs.	6	21	300%
P2	Days of anthesis date delay for each hour by which the day length exceeds P20	rate	0	6	375%
PHINT	Phylochron interval (Interval between successive leaf tip appearances)	GDD (°C)	20	110	200%

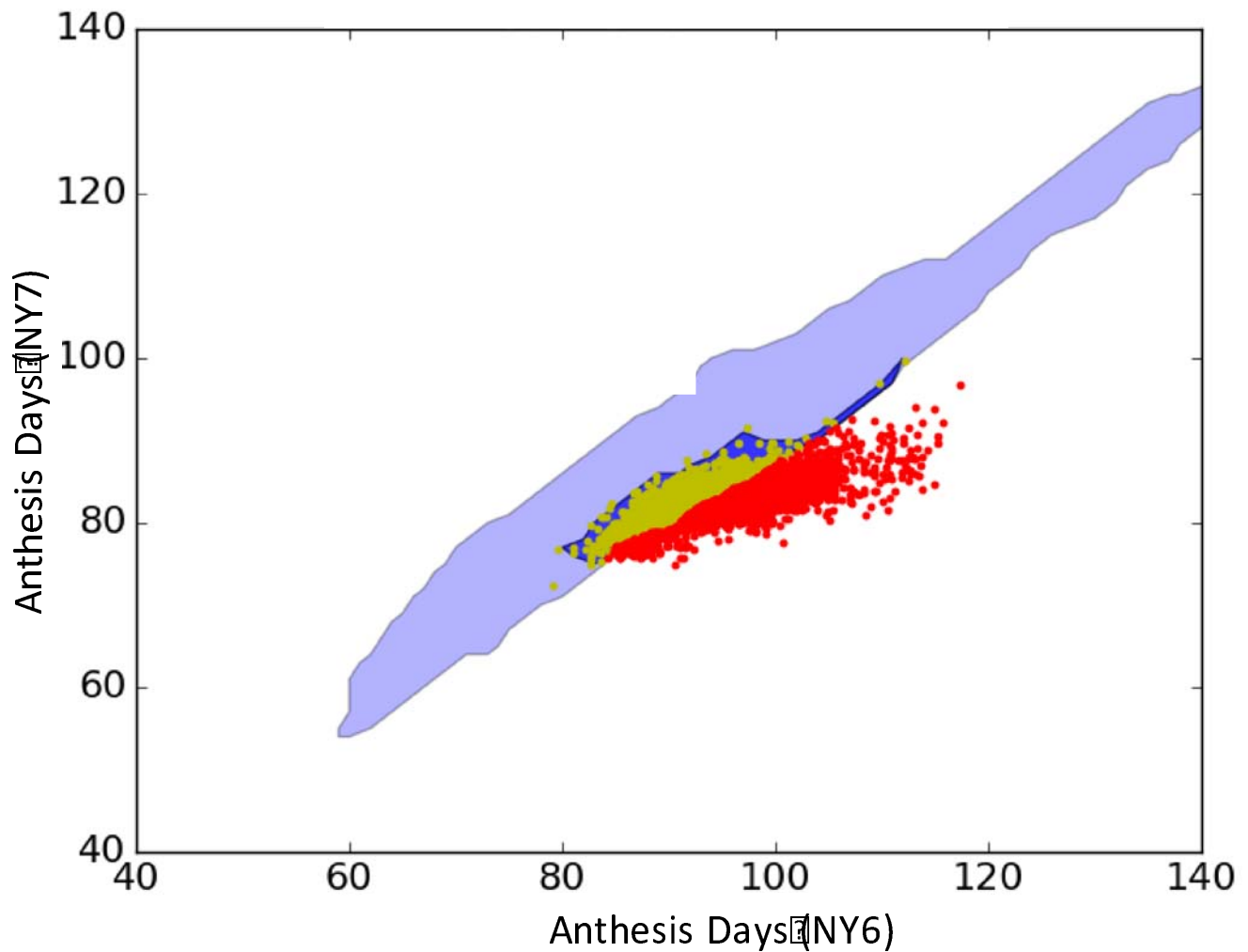
---

517

518 Fig. 8 shows overlapping predictions based on the database search under the range of parameters in  
 519 Table 2 and on the DE search under the extended range of parameter values (Table 4). Specifically, the light  
 520 blue area represents the anthesis date region that was reachable through predictions based on the database  
 521 search. In contrast, the dark blue area is the predicted anthesis date region within which the DE algorithm  
 522 converged. Note the almost perfect overlap of the lower edges of the light blue (i.e. database search) and dark  
 523 blue (i.e. DE search) areas, indicating that, despite its much larger starting parameter search space, DE did not  
 524 extend model predictions. This suggests limitations in model expressivity that go beyond the method of  
 525 parameter estimation or the initial parameter space used for the search.

526 As a corollary, it is worth noting that more site-years of data of similar quality are unlikely to improve  
 527 model expressivity, as illustrated by the following thought experiment. Suppose a community has developed  
 528 the univariate deterministic model  $y = \arctan(\theta)$ , where  $\theta$  is a parameter, with  $0 \leq \theta \leq 10$  by solid prior  
 529 knowledge and  $y$  is some dependent variable of interest. Assume that this is viewed as a very complex model  
 530 requiring simulation to solve. The community understands that no model is perfect, but no specific flaws of this  
 531 one are known. Extant data for  $y$  ranges from 1.31 to 1.61 and yields the point estimate  $\hat{\theta} = 5.79$  (RMSE =  
 532 0.12). Due to its complexity, no one has noticed that the model cannot reproduce any  $y > \arctan(10) = 1.47$   
 533 or, for any  $\theta$ , a  $y > \pi/2 \approx 1.57$ . Now suppose that: a very large set of new  $y$  data is collected. Depending on  
 534 the distribution of the new data either: (1) a new  $\hat{\theta} < 10$  will be found or (2)  $\hat{\theta}$  will rise significantly above 10,

535 leading to a rejection of the model. However, what will *not* happen is that the increase in data will enable  
536 observations >1.57 to be reproduced. The model simply lacks the expressivity to do so. Analogously, increasing  
537 the amount of anthesis date data may narrow GSP estimate confidence limits, but the reachable region of  
538 predicted phenotype space is unlikely to extend beyond the edges of the light blue regions. Therefore, any  
539 improvement in the ability to predict the large numbers of red points in Fig. 7 and 8 is unlikely.

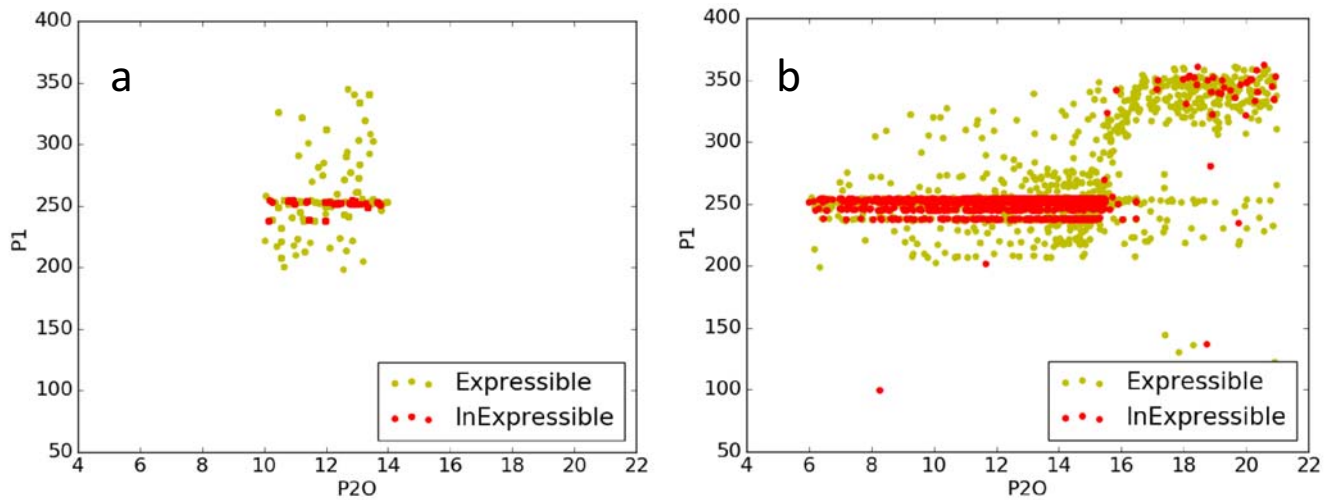


540

541 **Fig. 8.** Superimposed anthesis date results using NY6 and NY7 data illustrating that searches via database and  
542 DE optimization over a much larger parameter space are equally unable to reproduce the observations for lines  
543 shown as red dots.

544 Given these issues, a sensible follow-up question might be about what specific GSP estimates were  
545 reported for the red points? Here we report answers only for P1.

546



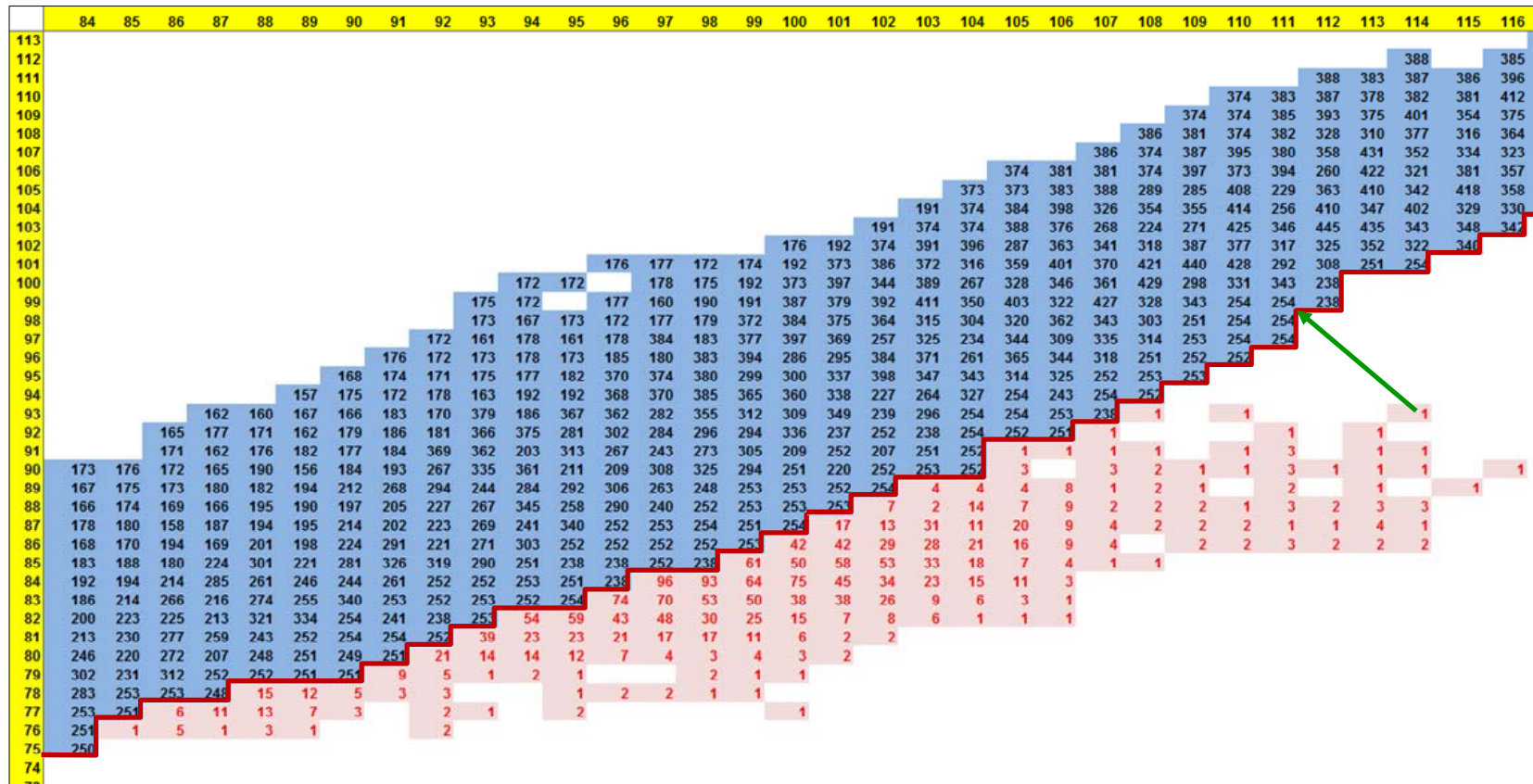
547

548 **Fig. 9.** Scatterplot of P1 vs. P2O estimates using data from NY6 and NY7 based on the database search (a) and  
549 Differential Evolution (b). Yellow and red dots are, respectively, observations characterized as expressible and  
550 inexpressible by model predictions.

551 Fig. 9 shows scatterplots of P1 and P2O estimates generated using data from NY6 and NY7 via the  
552 database search and DE. The color coding is consistent with that in Fig. 7a. The pronounced bands at ca.  
553 P1=250 in both panels are immediately striking – although the scale is small, a corresponding band is quite  
554 evident at the same position in Fig. 6. A tabulation reveals that, of all 4,731 lines represented in the Fig. 9a,  
555 3,227 (68.2%) have estimates of P1 ranging from 245 to 260. Of these, 1,493 are expressible (yellow) and 1,734  
556 (red) are not expressible. Out of the total 4,731 points in the graph 2,189 (46.2%) are expressible and 2542  
557 (53.8%) not. The Fig. 9b has similar proportions of expressible and inexpressible points (2327, 49.1%; and 2404,  
558 50.9%; respectively), reinforcing the similarity of results for parameter estimates from DE and database  
559 searches. The differences are likely due to the ability of DE to explore the parameter space continuously



560 whereas the database search is restricted to the predefined discrete Sobol points. Still, one may wonder why  
561 so many P1 estimates are near the 250 degree-days? Fig. 10 reveals the answer.



564 **Fig. 10.** P1 estimates from the database search (black) and the numbers of lines with inexpressible observations (red) arranged in a tableau organized  
 565 as a phenotype space plot corresponding to the center portion of Fig. 8. The dark red line is the expressibility frontier and the green arrow shows the  
 566 P1 value (254) from the GSP combination that minimizes the RMSE for one illustrative line. Horizontal and vertical yellow strips are the anthesis dates  
 567 for NY6 and NY7

568           The numbers in black are the “first-best-found” P1 estimated values that generate the corresponding  
569 row × column anthesis date combinations. A comparison with the corresponding dot colors and sizes in Fig. 5b  
570 indicates that, on the frontier (red borders Fig. 5a,b and 10) between expressible and inexpressible  
571 observations, there was essentially no equifinality and, concomitantly, narrow ranges of P1 values. Fig. 10  
572 shows that the P1 values along the frontier were all quite close to 250. For lines with observations falling  
573 outside the frontier, the RMSE was minimized by assigning GSP values associated with the closest achievable  
574 dates, i.e. those directly on the frontier. Therefore, all the lines counted by the red numbers were assigned P1  
575 values that are very close to 250 and have essentially no equifinality. The green arrow in Fig. 10 illustrates this  
576 phenomenon for one line. The nearest P1 estimate is 254 and the length of the arrow (ca. 5.8 days) is  
577 proportional to that line’s RMSE. Specifically, in this case the length is  $1/\sqrt{2}$  times the RMSE because there are  
578  $n = 2$  site-years.

579           Recall that the upper limit placed on P1 was 450 (and 600 in the DE search), therefore this outcome is  
580 likely not an artifact of constraints in the GSP search space but, rather, a result of poor model expressivity, that  
581 is the model inability to predict anthesis date pairs beyond those on the frontier. This mechanism accounts for  
582 the P1 band at 250 in Fig. 9a. Furthermore, as previously presented, more data cannot improve the prediction  
583 of inexpressible lines, the banding in Fig. 6 is not surprising.

#### 584 *4.5 P2O gap*

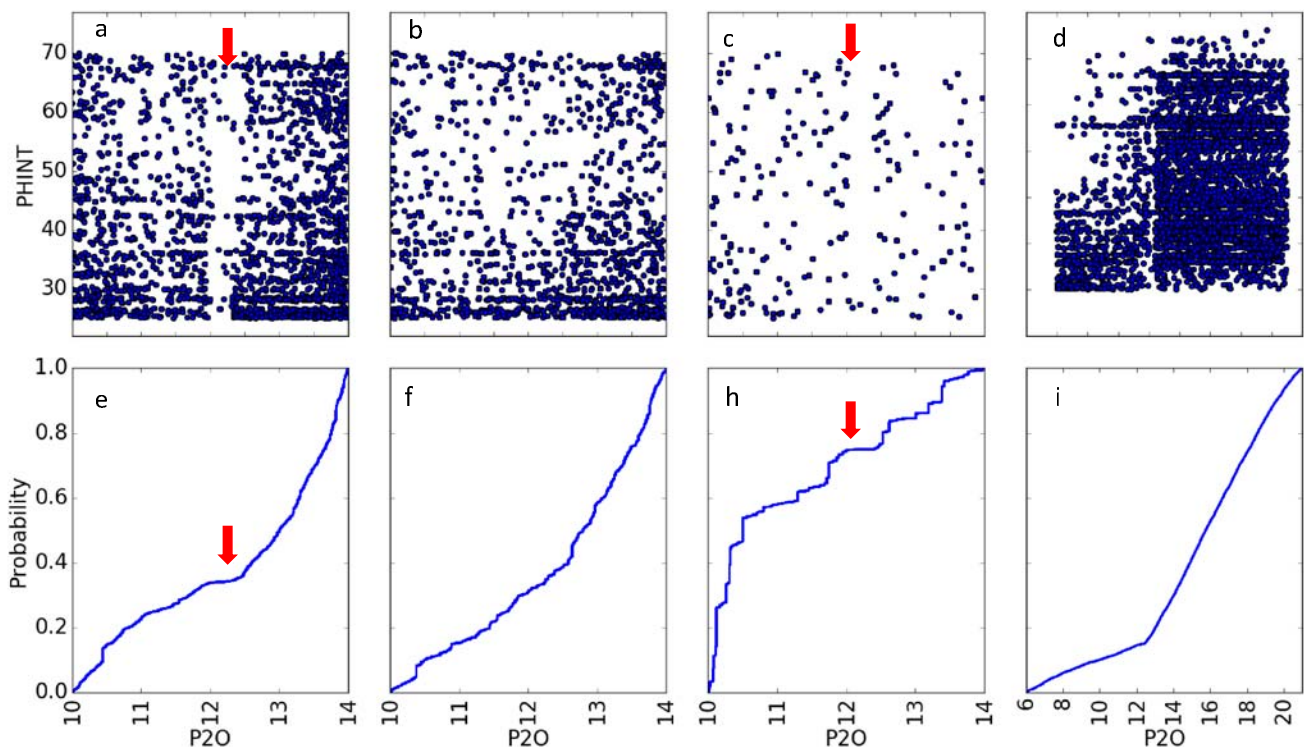
585           We now investigate the vertical gap in scatterplots involving P2O estimates (Fig. 6), which documents  
586 the intricacy of the interactions that can occur between model mechanisms, parameter ranges searched,  
587 optimization algorithms used, and environments included. Exploratory re-tabulations of the Sobol-based  
588 parameter database revealed that the P2O gap was clearly present in the three site-years having shorter day  
589 lengths (FL6, FL7, and PR6) but absent in fits obtained by only including the remaining eight site-years with  
590 longer days (Fig. 11). Fig. 12 shows that a substantial number of observations for short-day site-years are  
591 outside the predicted phenotype ranges expressible by the model under either database or DE optimization. As

592 described in section 3.2, the model operated by calculating the number of leaves initiated by the end of Stage 2  
593 and predicts anthesis only after leaves are fully emerged. For any line, leaf number was a constant across all  
594 site-years, namely  $P1/(2 \times PHINT) + 5$ . The variation of anthesis dates across plantings was such that there were  
595 few, if any, combinations of P1 and PHINT that were compatible with the data from all site-years. Therefore,  
596 the optimizer relied more heavily on the P2 and P2O parameters.

597 Specifically, the optimizer settled on very small P2O estimates, much smaller than the short southern  
598 photoperiods. Instead, the optimizer relied on P2 estimates to generate anthesis date predictions that were  
599 delayed to the greatest extent possible by lengthening Stage 2. Recall that P2O values above the day length  
600 make Stage 2 only four days long, which is not enough time for temperature differences to accumulate the  
601 needed variation. The abundance of low P2O estimates thus created the gap observed in scatterplots of P2O  
602 with other GSPs (Fig. 11a). In contrast, the photoperiods in the remaining longer-day site-years exceeded the  
603 maximum allowed P2O values in the P2O database search during (and long after) the juvenile period.  
604 Therefore, there was no empty band in the scatter plot (Fig. 11b) because the optimizer was able to exploit  
605 delays for any value of P2O.

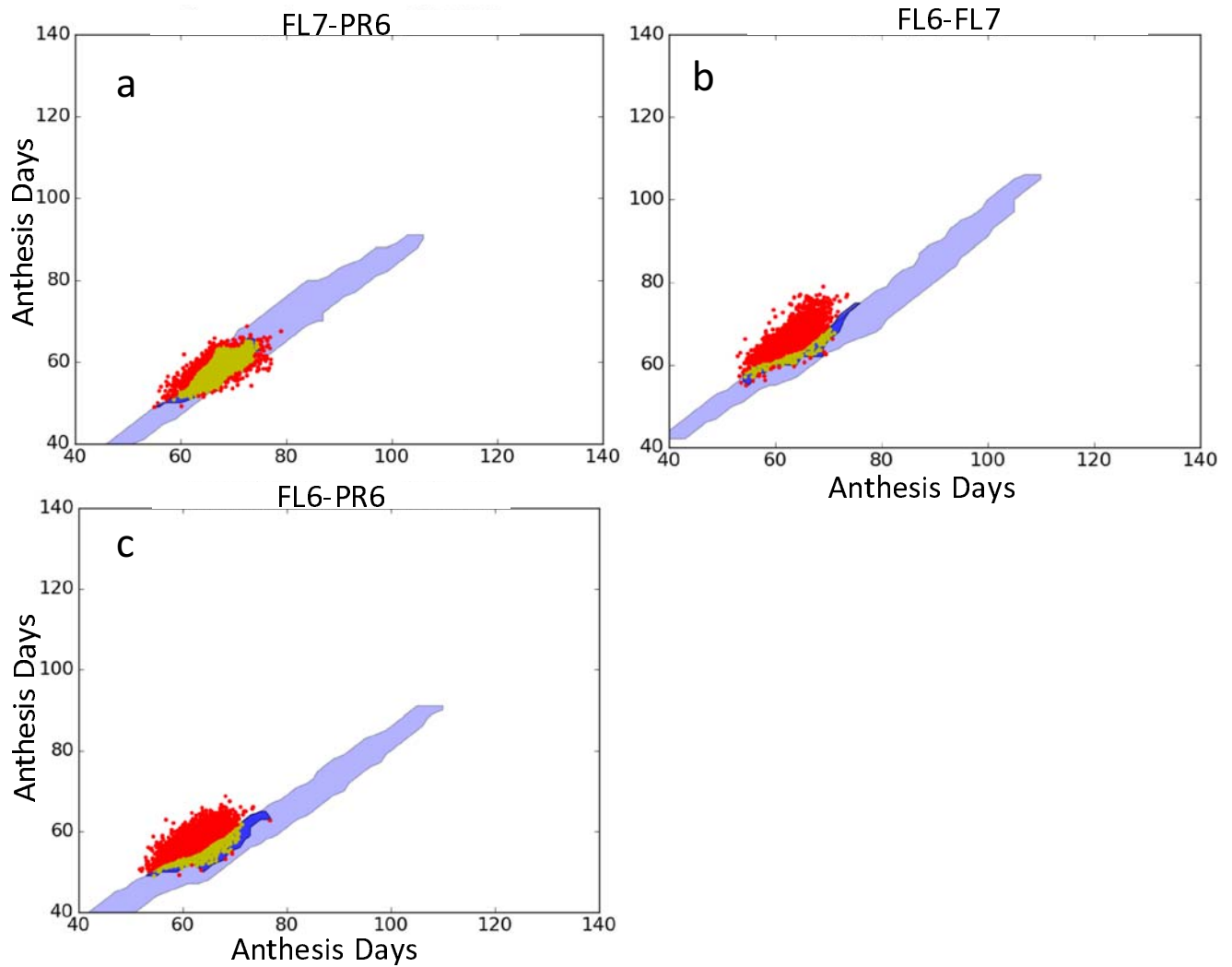
606 With the broader range of parameter values available to the DE runs and the increased flexibility  
607 available between P1 and PHINT, other options became available. In particular, in many cases DE found GSP  
608 combinations wherein P2O exceeded the southern day lengths so photoperiod had no influence on anthesis  
609 date and no gap artifact was generated (Fig. 11d,i). P1 and PHINT thus became the major explanatory  
610 parameters. This is shown in Fig. 13, whereby for each line, the parameter differences are plotted against the  
611 RMSE differences that result from changing the estimation methods from database to DE optimization. The DE  
612 estimate of P2O were larger in 4,507 out of 5,240 lines (87%; Fig. 13d), almost always by enough to put it above  
613 the local day lengths. In tandem, P1 values fell in 3,559 lines (Fig. 13a), whereas PHINT rose in 4,102 lines (Fig.  
614 13c).

615 Note, however, that for *any* (P1, PHINT) combination, *any* P2O that exceeds the local day length will  
616 give the same RMSE – a clear source of equifinality. Thus, the changes in P2O will not, in all likelihood, lead to  
617 values that can be more closely related to genetics. Moreover, because of the limits on model expressivity,  
618 none of the DE solutions gave significantly better fits than the database estimates. This is why virtually all  
619 points in Fig. 13 had DE RMSE's within 0.5 days (horizontal axes) of the database-based parameter estimates.  
620 This, too, is an illustration of equifinality because the two optimizers were finding different GSP estimates  
621 although the RMSE were of similar magnitude.



622

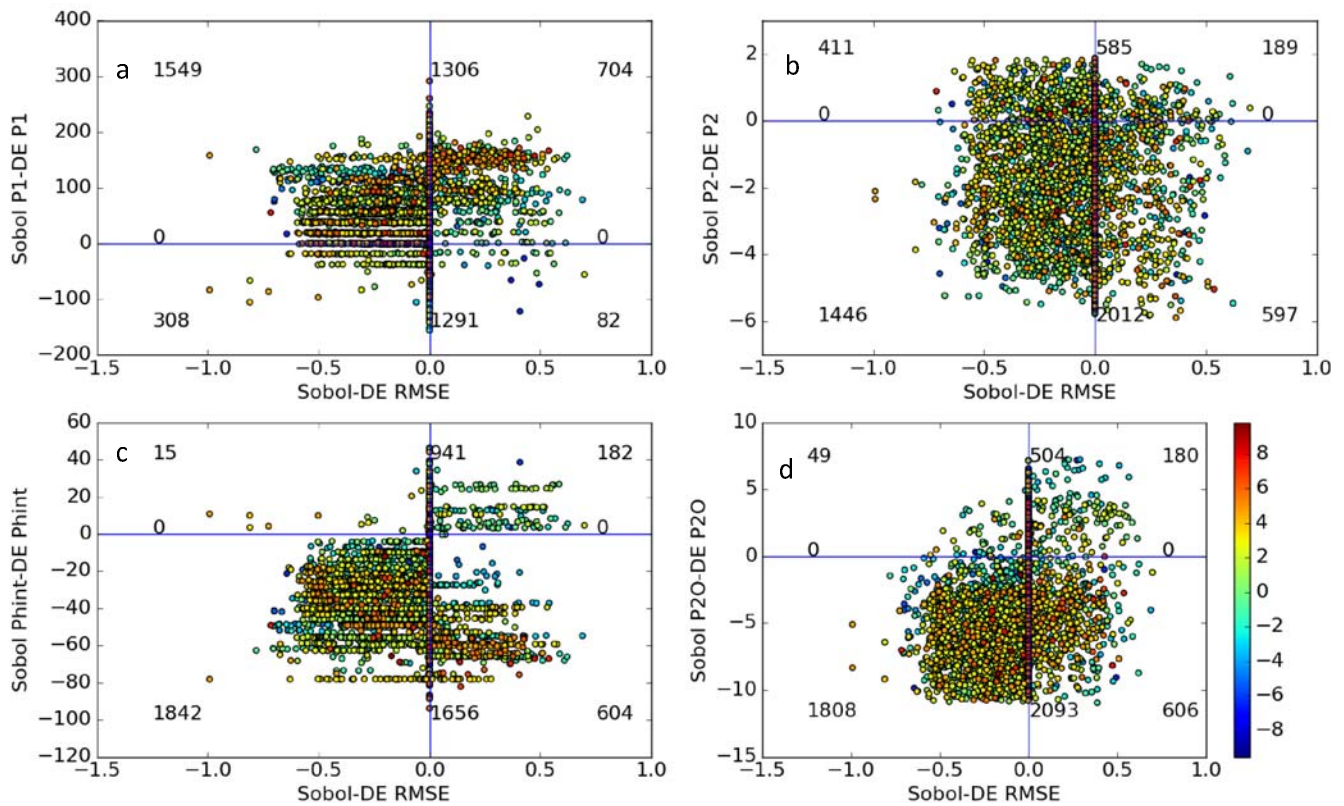
623 **Fig. 11.** P2O and PHINT scatter plots (top row) and P2O cumulative density functions (bottom row) using (a &  
624 e) all 11 site-years, ( b & f) longer day site-years, ( c & g) shorter day site-years based on the database approach,  
625 and ( d & i) shorter day site-years using the DE approach. All horizontal axes in both rows have the same scale.



626

627 **Fig. 12.** Phenotype space plots of observed and predicted values based on the three site-years with shorter  
628 days. Note the large number of points in the FL6-PR6 and FL6-FL7 plots that lie above the dark blue prediction  
629 region based on DE.

630



631

632 **Fig. 13.** The differences in parameter estimates from database search vs. DE (vertical axes) plotted against the  
 633 corresponding difference in RMSE for 5240 lines in FL6, FL7, and/or PR6. The color encodes the sum of residual  
 634 (observed minus mean) across site-years for each line.

#### 635 4.6 Tests for stability of GSP estimates

636 Table 5a shows the effect of including or excluding the effect of different subsets of site-years on the  
 637 modeling of estimates (Equation 1) for each GSP for the base set. For all GSP parameters, AIC and BIC values  
 638 were considerably smaller for models that included the random effect of site-year subsets,  $\beta_e$ , therefore  
 639 suggesting non-negligible variability across site-year subsets on the GSP estimates. The table illustrates the size  
 640 of the site-year set effects as follows. For scaling purposes, we provide the estimated intercept,  $\hat{\mu}_p$ , which also  
 641 serves also as an estimated GSP grand mean across all lines and site-year subsets. The Index of Variability  
 642 (expressed as a percent) is the standard deviation of the  $\beta_e$  effect normalized by the grand mean. The

643 percentage of the total GSP variance ( $\sigma_e^2 + \sigma_l^2 + \sigma_r^2$ ) attributable to site-year subsets is also shown. Both of  
644 these descriptors indicate substantial variability between site-year sets, with indexes of variability ranging from  
645 5.9% for P20 to 33.6% for P2 and over 20% of the total variance related to site-year sets for all GSP's.

646 The Chi square values from the likelihood ratio test and the associated  $p$ -values are presented in the  
647 last two columns of Table 5a. The extreme  $p$ -values demonstrate that the GSP values depend on the set of site-  
648 years used to estimate them. Therefore, the GSP's are not, in fact, genotype specific despite the goodness-of-fit  
649 displayed in Fig. 3. This result is completely understandable given the range of artifacts due to equifinality and  
650 model expressivity issues identified above.

651 Table 5b shows the results when only estimates having ties are tested (left) vs. an analysis that includes  
652 all estimates (right). The former corresponds to estimates for lines whose observations fall inside the  
653 expressivity frontier and the latter includes the estimates for all lines. It is clear that the grand means, index of  
654 variability, and percentages of GSP variance are highly similar between all three groupings in Table 5. Also, all  $p$ -  
655 values are extremely significant and increase with the amount of data used.

656

657



658 **Table 5a.** Estimated log likelihood, fit statistics, selected summary measures, and a likelihood ratio test for  
 659 competing statistical models fitted on GSP estimates with and without the random effect of site-year subset,  
 660 based on GSP estimates for the base group (N=60,834).

GSP	Log			GSP Grand Mean $\hat{\mu}_\rho$	Index of Variability <sup>c</sup> $\sigma_e / \hat{\mu}_\rho$	Variance pcts. for site-year sets <sup>c</sup> $\sigma_e^2 / \sigma_{tot}^2$	Chi- square test statistic	Chi- square <i>p</i> -value <sup>d</sup> (df =0.5)
	likelihood w/o (top) and w/ (bot) a site-year set effect <sup>a</sup>	AIC w/o (top) and w/ (bot) a site-year set effect <sup>b</sup>	BIC w/o (top) and w/ (bot) a site-year set effect <sup>b</sup>					
P1	-338046	676098	676125	264.625	12.30	34.38	30714	10 <sup>-13334</sup>
	-322689	645386	645422					
P2	-46154	92313	92340	1.037	33.55	33.92	41833	10 <sup>-18163</sup>
	-25237	50482	50518					
P20	-105304	210614	210642	12.2440	5.88	27.83	19894	10 <sup>8635</sup>
	-95357	190723	190759					
PHINT	-254875	509756	509783	44.167	15.44	22.62	15943	10 <sup>6919</sup>
	-246903	493815	493851					

661 <sup>a</sup> Larger is better <sup>b</sup> Smaller is better <sup>c</sup> Chernoff upper bound on Chi-squared cum. dist. function.

662

663

664 **Table 5b.** Summary measures and likelihood ratio  $p$ -values for competing statistical models fitted on GSP  
 665 estimates with and without the random effect of site-year subset from data only having ties (left) and all data  
 666 (right).

GSP	Variance				Variance			
	GSP Grand Mean $\hat{\mu}_\rho$	Index of Variability <sup>c</sup> $\sigma_e / \hat{\mu}_\rho$	pcts. for site-year sets <sup>c</sup> $\sigma_e^2 / \sigma_{tot}^2$	Chi- square $p$ -value <sup>d</sup> (df =0.5)	GSP Grand Mean $\hat{\mu}_\rho$	Index of Variability <sup>c</sup> $\sigma_e / \hat{\mu}_\rho$	pcts. for site-year sets <sup>c</sup> $\sigma_e^2 / \sigma_{tot}^2$	Chi- square $p$ -value <sup>d</sup> (df =0.5)
With Ties (N=114,314)					With all Data (N=177,870)			
P1	273.5	11.37	29.77	$10^{-23283}$	270	11.48	29.94	$10^{-34955}$
P2	0.9137	36.33	35.23	$10^{-34723}$	0.9593	35.5	33.8	$10^{-52518}$
P2O	12.49	4.43	19.70	$10^{-11883}$	12.42	4.88	21.27	$10^{-19806}$
PHINT	43.57	18.65	26.31	$10^{-17348}$	43.94	17.3	24.35	$10^{-23740}$

667 <sup>a</sup> Larger is better <sup>b</sup> Smaller is better <sup>c</sup> Chernoff upper bound on Chi-squared cum. dist. function.

## 668 5. Discussion

669 Since their inception, ecophysiological models have been evaluated in terms of predictive ability, which  
 670 are superb in many circumstances (Batchelor et al., 2002). The model parameters were considered to be *inputs*  
 671 whose genesis was secondary as long as the model outputs proved useful. However, as often happens in  
 672 science, perceived needs, desiderata, and requirements escalate as technologies evolves. In particular, we are  
 673 now demanding that the model inputs themselves be the accurate outputs of processes at the genetic level  
 674 that can be modeled by genomic prediction. It is not surprising, therefore, that modeling technologies (ranging  
 675 from data collection to estimation) that were adequate for past applications now require improvement.

676 From a fundamental but traditional perspective, there are several issues of perennial concern in crop  
677 modeling. The first is model functional structure including both its degree of expressivity and its behavior under  
678 optimization. For example, estimation procedures like DE, that primarily yield point estimates, are limited in  
679 their ability to assess equifinality. At best, one can query the flatness of the goodness-of-fit function in the  
680 neighborhood of the estimate, but this does not tell anything about the ubiquity of equifinality across the  
681 parameter space. Nor do these procedures allow one to detect observations that fall outside of the model's  
682 scope of expressivity unless the discrepancies are quite large. Doing so requires methods like the Sobol  
683 database scheme used here that can make broader assessments in both parameter and phenotype space. It  
684 may well be that the rarity with which database methods have been used has led to an underappreciation as to  
685 the prevalence of these adverse situations.

686 When expressivity issues are identified, results like those above are not likely to be solved merely by  
687 acquiring more data of the same type. In such situations, better models will often needed and modern genetic  
688 studies can help. A great many plant component subsystems are currently under study at the molecular level.  
689 Indeed, some of these (e.g., Chew et al., 2014) are even being combined into multi-scale organ and whole plant  
690 models. Even without modeling directly at the genetic level one can use the derived insights to make informed  
691 choices between alternative representations of individual ecophysiological processes. Tardieu (2003) refers to  
692 such representations as "meta-mechanisms". It would seem plausible that building models from component  
693 parts of increased biological realism should increase the ability to reproduce field variation – at the very least, it  
694 is hard to see how it can hurt. As a concrete example, the B73 parent is photoperiod insensitive. In CERES-  
695 Maize, however, the only way to express this is by setting P20 in excess of the observed photoperiods, with the  
696 consequences we have seen.

697 This is not to say, however, that both more and better data are not needed. Indeed, data quality issues  
698 can impact both expressivity and GSP stability. For example, while the date seed that are physically sown in a  
699 field is usually known and not subject to error, researchers often report a subjective notion of "effective sowing

700 date” based on their interpretation of whether low soil moisture delayed germination. If errors in sowing date  
701 push an anthesis observation across the expressivity frontier, erroneous GSP estimates will result. Such errors  
702 can also arise if different personnel are involved across locations or growing seasons, especially for visually  
703 evaluated phenotypes like most phenological traits. Providing the emergence date can provide a partial check  
704 for these problems and also for errors in simulating time from sowing to emergence. Unfortunately,  
705 emergence dates were not reported for the maize NAM dataset.

706 Another traditional modeling concern has always been the relationship between the observed  
707 environmental data and the immediate environmental conditions actually experienced by individual plants.  
708 Weather data can suffer from multiple sources of bias and error (Fall et al., 2011). For example, stations that  
709 are not located within or directly adjacent to experiments may have bias due to local variation in weather  
710 conditions. Additionally, although of limited concern for anthesis dates, the quality of soil and management  
711 data. In this study any systematic differences in protocols for collection of weather data between the sites as  
712 aggravated by small sample effects, might have contributed to some degree to the significance levels in Table 5.  
713 It would certainly be desirable to have a method by which this potential effect might be quantitatively assessed.  
714 Such a method could be instrumental in designing experimental procedures for reducing the problem. One  
715 potential example might be to eschew external measurements of some environmental variables (e.g., air  
716 temperature) and use sensors onboard UAV’s or other automated vehicles to measure plant temperatures or  
717 other critical features directly at high temporal and spatial frequencies.

718 More involved data types and structures are also needed to resolve issues of equifinality when they  
719 arise. Equifinality is fundamentally a problem of discernment. In simple terms, given an equation  $c = a + b$ , if  
720 one only has data on  $c$ , then estimates of  $a$  and  $b$  are doomed to be equifinal. If one desires otherwise, one  
721 must find a way to measure either  $a$  or  $b$ . Current technological efforts to develop high throughput  
722 phenotyping approaches might be quite helpful in this regard. For example, assuming that  
723  $TOLN = P1 / (PHINT \times 2) + 5$  is the correct way to model the number of leaves at anthesis, data on total leaf number

724 would help constrain the parameter estimates. This leads toward a range of constrained and/or multiobjective  
725 estimation procedures on which there has been significant amounts of research (Rabotyagov et al., 2012;  
726 Tatsumi, 2016). Maximum entropy methods offer another opportunity wherein one identifies a probability  
727 distribution of values that is constrained by but mathematically no more informed than is justified by a set of  
728 potentially diverse data types (Hess et al., 2002). Another alternative might be Bayesian methods with  
729 multivariate likelihood functions that combine several observational variables (Franks et al., 1999).

730 Another approach to reduce equifinality would be to use simpler models. The fewer the number of  
731 processes and GSP's in a model, the smaller the opportunity for hard-to-spot tradeoffs to exist wherein  
732 adjustments to one parameter can be offset by tweaking another one. Of course, the tradeoff may be less  
733 expressivity leading to other problems. However, Welch et al. (2005) presented 12 dichotomies comparing  
734 gene network modeling and quantitative genetics approaches, where aspects of the former might also apply to  
735 ecophysiological modeling. They opined that an optimal modeling approach should entail a synthesis of both.  
736 The key features to be contributed from the network (i.e., ecophysiological) side would be (1) the ability to  
737 handle time-varying dynamics, (2) a far more parsimonious approach to expressing biological and biology  $\times$   
738 environmental interactions, and (3) a more mechanistic explanation of how traits originate. It is at least  
739 conceivable that some way station of moderate complexity exists between statistical genetics and full crop  
740 models that can achieve this.

741 At whatever level of complexity proves appropriate, one cannot accurately estimate the parameters  
742 controlling model components without collecting data on settings wherein the relevant processes operate  
743 differentially. This is clear from the P2O gap phenomenon, which was apparent when only short day data was  
744 used and absent under long days. Both settings distorted the results, in one case compressing estimates into a  
745 restricted range, leaving a gap, and, in the other, allowing them to spread out. Furthermore, this interacted  
746 with the range of values allowed, which caused shifts between (P1, PHINT) and (P2, P2O) as to which

747 parameters appeared to be “explanatory”. The debilitating influence of such behavior on linking parameter  
748 values to genes is terribly obvious.

749           However, it also should not escape notice that the gap was evident even in a mixture of environments,  
750 suggesting that good experimental design entails more than just making sure that a suitable range of  
751 environments is included. There is some notion of balance that needs to be established and applied globally to  
752 data selection. In this context, it is worth noting that despite the fact that thousands of lines were planted in  
753 each location, there were only 539 lines where data were reported from all 11 trials. However, given the  
754 expense of such large-scale trials and the multiple purposes each one will serve, “balance” cannot mean  
755 “orthogonality” where all lines are planted at all sites. Of course, an established benefit of ecophysiological  
756 models is to serve as guides to help prioritize experimentation over time. It seems likely that as their  
757 integration with statistical genetic models expands, they might also be able to assist in the rational planning and  
758 resource allocation for large, multi-site trials.

759           Another approach entirely would be to seek to move beyond a two-step “estimate and then map”  
760 paradigm. Conventional mapping methods essentially isolate genetic markers whose pattern of assignment to  
761 lines mirrors the pattern of phenotype values of interest. A general linear model is assumed to mediate  
762 between marker states and realized phenotypes. There is no conceptual reason why that general linear model  
763 might not be replaceable by a crop model. In effect, one could conceive a hierarchical model in which a first-  
764 level model is specified on the data and higher order submodels are specified on the parameters that  
765 characterize the behavior of observed data, much like proposed by Bello et al. (2010).

766           One could conceptually implement this hierarchy in the context of crops by fitting phenotypes with an  
767 ECM whose GSP's are then specified as functions of genetic markers at another level of the hierarchical model.  
768 Indeed, this is what the current paradigm attempts, except that the two-step estimation process curtails  
769 smooth borrowing of information across hierarchical levels of the model that could potentially help resolve the  
770 equifinality problem.

771 We acknowledge that one-step hierarchical model approach might not solve the sort of expressivity  
772 problems described in the thought experiment and documented in our results (both in 4.4). Yet, it would enable  
773 the genetic structure of the population to inform the GSP estimation process. The potential utility of this  
774 hierarchical modeling approach is currently under study in one of our labs. The approach would also enable  
775 more efficient use of data. Currently, the two-step approach requires data from multiple environments (Welch  
776 et al., 2002) for each line in order to estimate the GSP's before mapping can proceed. However, consider a line  
777 that was culled very early in the selection process, perhaps even after a single round. Because the parameters  
778 estimated in putative one-step hierarchical modeling schemes would include marker effects, even just one  
779 planting becomes a usable observation if the line is genotyped. This is a sufficiently inexpensive operation now  
780 that some programs (e.g. CIMMYT; (Battenfield et al., 2016) are doing so routinely for the offspring of all  
781 crosses.

782 A one-step hierarchical modeling approach might also make it possible to utilize data taken on lines  
783 after they enter the market place. Analogously to high throughput phenotyping in breeding programs,  
784 precision agricultural management is also investing in sensor- and model-based approaches to improve  
785 productivity (Thompson et al., 2015; Thorp et al., 2015) while collecting a wealth of multivariate data. Usually,  
786 of course, hybrids are released into areas where they show low G×E interactions. For example, a line with a  
787 particular P2O is not likely to be released across a sufficient range of latitudes to have great differences in day  
788 length. This would make it difficult to directly estimate P2O for the line using the methods described in this  
789 paper.

790 However, in a one-step hierarchical model approach, one would only be looking for markers that  
791 influenced P2O. In this case, data from many lines and geographical areas could be used together. This would  
792 also make such data usable for the sorts of hypothesis testing about genes discovered by other means, thus  
793 facilitating genetically-informed ecophysiological modeling. For such approaches to be workable, however,  
794 there are many policy issues to be resolved including information property rights and fair economic returns to

795 data, not to mention the need to greatly harden cybersecurity protections (FBI, 2016). However, if this can be  
796 done then issues of environmental coverage would likely be ameliorated due to the extent of the data that  
797 would become available.

## 798 **6. Conclusions**

799 The original and seemingly simple goal of this study was to first fit the anthesis date component of the  
800 CERES-Maize model to data from over 5000 genotyped lines and then genetically map the resulting GSP values.  
801 However, we were unexpectedly detoured when we found that despite the high predictive quality of the values  
802 obtained, there were numerous artifacts that emerged in the estimation process, thereby making our  
803 immediate goal unachievable. We find it interesting that the problems we encountered would likely be  
804 invisible, though present, in smaller data sets and, unless addressed by suitable research, these problems bode  
805 ill for understanding any genetic underpinnings of ecophysiological models. This is worrisome given the recent  
806 escalating attention that has been given to this method of melding ecophysiological and statistical genetic  
807 models as a way of accelerating the crop improvement process so as to help meet global food and fiber needs  
808 by 2050.

809 The constraining issues fall into two categories. The first arises in situations where the model is unable  
810 to express the observed data for some line even by a relatively few days. In this circumstance, the line is  
811 assigned the GSP associated with the nearest point on model's expression frontier – values which can, however,  
812 change only slowly along that boundary. The result is that many and in some cases a large majority of lines are  
813 assigned the same GSP values independent of their actual genetics.

814 The second symptom arises when the model can reproduce the data. In these instances, there can be  
815 many combinations of GSP values that predict equally well. When such equifinality exists, there is no principled  
816 way to assign the line a genetically relevant value. In short, when the model can express the data there is no  
817 unique combination of GSP values and, when unique combinations do exist they are often values being given to  
818 many lines because of a deficiency in model expressivity.



819           This finding is rather remarkable because in both breeding efforts and, indeed, genetic studies as a  
820 whole, anthesis date is considered, if not a simple trait, at least one that has proved much easier to elucidate  
821 than many others. In addition, it is generally, much more readily predicted by classical phenology models for  
822 reasons that, themselves, have become generally understood (Wilczek et al., 2009). This cannot but make one  
823 wonder, what pitfalls might lie in wait for efforts to probe other, more involved traits.

824           Therefore, the next question to be asked by follow-on research is how prevalent are these phenomena.  
825 The best way to do that would seem to be to use Sobol database search methods. This is because, unlike  
826 optimizers that find single “best estimates”, the database approach will reveal the both the extent of the  
827 expressible phenotype regions as well as a direct measure of the extent of any equifinality.

828           However, despite the ability to reuse results databases for many searches, undertaking such a program  
829 in any broadly based fashion will be highly demanding computationally. For this reason, strong consideration  
830 should be given to disaggregating comprehensive models into separate modules that can be studied  
831 independently at much lower computational cost. (This is what we did for the limited DE run, although Python  
832 certainly is not a high performance language.) A better long-term strategy would be to program future models  
833 in a manner that supports single-module testing at the source code level. Doing so will facilitate the whole-  
834 model verifications needed to ensure that fragmentation into modules for testing and improvement by  
835 different labs does not compromise integration at the level of the scientific community.

836           As module testing and innovation progress, it will be of strategic value to ground improvements in  
837 advancing genetic understanding at the molecular level. While this might seem daunting to those versed in  
838 purely physiological approaches, it need not be so. One of the most venerable concepts in all of the life  
839 sciences is that of the biological hierarchy that is, a series of many functional levels extending from molecules to  
840 the biosphere. One of the perspectives emerging from molecular science is that that hierarchy might, be  
841 operationally much flatter than commonly believed. That is, simple changes at lower levels can easily create  
842 tangible responses multiple levels higher. To the extent that this is true, it greatly reduces the complexity of

843 bridging across those levels. This is the philosophy behind the meta-mechanism approach mentioned earlier  
844 (Tardieu, 2003; Tardieu and Tuberosa, 2010).

845         That approach has a proven ability to account for environmental interactions with sufficient skill to  
846 eliminate observed G×E interactions from GSP's in the data sets used (Reymond et al., 2003). However, as  
847 shown by the *p*-values in Table 5, the very large data set used herein conveyed an extraordinary power to  
848 detect site-year dependencies in GSP estimation. Indeed, so powerful as to make one wonder if an insignificant  
849 result is scientifically achievable by any even remotely feasible research effort? A better number to use for  
850 practical evaluations might be the index of variability in Table 5. This would give a clear index of the size of the  
851 effect as a percentage of the parameter values. Also, means exist for comparing such indices to see if  
852 reductions in their values (i.e. by an improved model with lowered site-year set dependency) are statistically  
853 significant (Vangel, 1996).

854         A final message from our research is that one cannot fix problems that one does not know exist.  
855 Community interest in the fitting-and-mapping paradigm has been high as shown by the heavy citation rates for  
856 the seminal papers in this area. For example, as of September, 2016, the Hammer et al. (2006) paper had been  
857 cited 257 times and those publications, *themselves*, had been cited by 6,370 others (Source: Google Scholar).  
858 There is also no doubt as to the importance of the ability to predict the behaviors of novel genotypes in novel  
859 environments while crosses are still in the planning stage. Indeed, this is precisely the genotype-to-phenotype  
860 problem, which has been declared by the National Research Council to be a top-priority goal for applied biology  
861 (NRC, 2008). So these impediments need to be overcome. However, with methods now in hand to detect  
862 adverse model behaviors under estimation, research that is probing ever more deeply into the control  
863 mechanisms of plant growth and development, and concrete tests to document model improvements, there is  
864 no reason to believe that we cannot do so.

865 **Acknowledgements**

866           The plan to use the maize NAM data to first developed through discussions with iPlant  
867 ([www.iplantcollaborative.org](http://www.iplantcollaborative.org)) on novel applications of cyberinfrastructure in plant science. The authors  
868 acknowledge the Texas Advanced Computing Center (TACC; <http://www.tacc.utexas.edu>) at The University of  
869 Texas at Austin and Beocat, Kansas State University for providing high performance computing resources that  
870 have contributed to the research results reported within this paper. Support for this effort was also supplied by  
871 the Department of Agronomy at Kansas State University. Additional support derived from a Higuchi-KU  
872 Endowment Research Achievement Award through the University of Kansas and the University of Kansas  
873 Endowment Association. This paper is contribution number 17-134-J of the Kansas Agricultural Experiment  
874 Station at Kansas State University, Manhattan, KS.

875

876 **References**

- 877 Akaike, H., 1973. Maximum likelihood identification of Gaussian autoregressive moving average models.  
878 *Biometrika* 60, 255–265.
- 879 Andrés, F., Coupland, G., 2012. The genetic basis of flowering responses to seasonal cues. *Nat. Rev. Genet.* 13,  
880 627–639.
- 881 Batchelor, W.D., Basso, B., Paz, J.O., 2002. Examples of strategies to analyze spatial and temporal yield  
882 variability using crop models. *Eur. J. Agron.* 18, 141–158.
- 883 Bates, D., Mächler, M., Bolker, B., Walker, S., 2014. Fitting linear mixed-effects models using lme4. *ArXiv Prepr.*  
884 *ArXiv14065823*.
- 885 Battenfield, S.D., Guzmán, C., Gaynor, R.C., Singh, R.P., Peña, R.J., Dreisigacker, S., Fritz, A.K., Poland, J.A., 2016.  
886 Genomic selection for processing and end-use quality traits in the CIMMYT spring bread wheat breeding  
887 program. *Plant Genome* 9.
- 888 Bello, N.M., Steibel, J.P., Tempelman, R.J., 2010. Hierarchical Bayesian modeling of random and residual  
889 variance–covariance matrices in bivariate mixed effects models. *Biom. J.* 52, 297–313.
- 890 Bratzel, F., Turck, F., 2015. Molecular memories in the regulation of seasonal flowering: from competence to  
891 cessation. *Genome Biol.* 16, 1.
- 892 Buckler, E.S., Holland, J.B., Bradbury, P.J., Acharya, C.B., Brown, P.J., Browne, C., Ersoz, E., Flint-Garcia, S.,  
893 Garcia, A., Glaubitz, J.C., others, 2009. The genetic architecture of maize flowering time. *Science* 325,  
894 714–718.
- 895 Burhenne, S., Jacob, D., Henze, G.P., 2011. Sampling based on Sobol’ sequences for Monte Carlo techniques  
896 applied to building simulations, in: *Proc. Int. Conf. Build. Simulat.* pp. 1816–1823.
- 897 Chenu, K., Chapman, S.C., Tardieu, F., McLean, G., Welcker, C., Hammer, G.L., 2009. Simulating the Yield  
898 Impacts of Organ-Level Quantitative Trait Loci Associated With Drought Response in Maize: A “Gene-to-  
899 Phenotype” Modeling Approach. *Genetics* 183, 1507–1523. doi:10.1534/genetics.109.105429
- 900 Chew, Y.H., Wenden, B., Flis, A., Mengin, V., Taylor, J., Davey, C.L., Tindal, C., Thomas, H., Ougham, H.J., de  
901 Reffye, P., others, 2014. Multiscale digital Arabidopsis predicts individual organ and whole-organism  
902 growth. *Proc. Natl. Acad. Sci.* 111, E4127–E4136.
- 903 Cooper, M., Technow, F., Messina, C., Gho, C., Totir, L.R., 2016. Use of Crop Growth Models with Whole-  
904 Genome Prediction: Application to a Maize Multienvironment Trial. *Crop Sci.*
- 905 Das, S., Suganthan, P.N., 2011. Differential evolution: a survey of the state-of-the-art. *IEEE Trans. Evol. Comput.*  
906 15, 4–31.
- 907 Dong, Z., Danilevskaya, O., Abadie, T., Messina, C., Coles, N., Cooper, M., 2012. A gene regulatory network  
908 model for floral transition of the shoot apex in maize and its dynamic modeling. *PLoS One* 7, e43450.
- 909 Fall, S., Watts, A., Nielsen-Gammon, J., Jones, E., Niyogi, D., Christy, J.R., Pielke, R.A., 2011. Analysis of the  
910 impacts of station exposure on the US Historical Climatology Network temperatures and temperature  
911 trends. *J. Geophys. Res. Atmospheres* 116.
- 912 FBI, 2016. Smart Farming May Increase Cyber Targeting Against US Food and Agriculture Sector (No. 160331-  
913 001). United States Federal Bureau of Investigation. Private Industry Notification.
- 914 Flint-Garcia, S.A., Thuillet, A.-C., Yu, J., Pressoir, G., Romero, S.M., Mitchell, S.E., Doebley, J., Kresovich, S.,  
915 Goodman, M.M., Buckler, E.S., 2005. Maize association population: a high-resolution platform for  
916 quantitative trait locus dissection. *Plant J.* 44, 1054–1064.
- 917 Franks, S.W., Beven, K.J., Gash, J.H., 1999. Multi-objective conditioning of a simple SVAT model. *Hydrol. Earth*  
918 *Syst. Sci. Discuss.* 3, 477–488.
- 919 Gill, P.E., Murray, W., Wright, M.H., 1981. *Practical optimization*. Academic Press.

- 920 Godfray, H.C.J., Beddington, J.R., Crute, I.R., Haddad, L., Lawrence, D., Muir, J.F., Pretty, J., Robinson, S.,  
921 Thomas, S.M., Toulmin, C., 2010. Food security: the challenge of feeding 9 billion people. *science* 327,  
922 812–818.
- 923 Grimm, S.S., Jones, J.W., Boote, K.J., Hesketh, J.D., 1993. Parameter estimation for predicting flowering date of  
924 soybean cultivars. *Crop Sci.* 33, 137–144.
- 925 Hammer, G., Cooper, M., Tardieu, F., Welch, S., Walsh, B., van Eeuwijk, F., Chapman, S., Podlich, D., 2006.  
926 Models for navigating biological complexity in breeding improved crop plants. *Trends Plant Sci.* 11, 587–  
927 593.
- 928 Hammer, G.L., Woodruff, D.R., Robinson, J.B., 1987. Effects of climatic variability and possible climatic change  
929 on reliability of wheat cropping—a modelling approach. *Agric. For. Meteorol.* 41, 123–142.
- 930 Harrison, S.R., Tamaschke, H.U., 1984. Applied statistical analysis. Prentice-Hall of Australia.
- 931 He, J., Jones, J.W., Graham, W.D., Dukes, M.D., 2010. Influence of likelihood function choice for estimating crop  
932 model parameters using the generalized likelihood uncertainty estimation method. *Agric. Syst.* 103,  
933 256–264. doi:10.1016/j.agsy.2010.01.006
- 934 Hess, C.P., Liang, Z.-P., Lauterbur, P.C., 2002. Maximum cross-entropy generalized series reconstruction, in: 5th  
935 IEEE EMBS International Summer School on Biomedical Imaging, 2002. Presented at the 5th IEEE EMBS  
936 International Summer School on Biomedical Imaging, 2002, p. 8 pp.–. doi:10.1109/SSBI.2002.1233979
- 937 Hoogenboom, G., Jones, J.W., Wilkens, P.W., Porter, C.H., Hunt, L.A., Singh, U., Lizaso, I., White, J., Uryasev, O.,  
938 Ogoshi, R.M., Koo, J., Shelia, V., Tsuji, G.Y., 2015. Decision Support System for Agrotechnology Transfer  
939 (DSSAT) version 4.5 (<http://dssat.net>), DSSAT Foundation. Prosser, Washington.
- 940 Hung, H.-Y., Shannon, L.M., Tian, F., Bradbury, P.J., Chen, C., Flint-Garcia, S.A., McMullen, M.D., Ware, D.,  
941 Buckler, E.S., Doebley, J.F., Holland, J.B., 2012. ZmCCT and the genetic basis of day-length adaptation  
942 underlying the postdomestication spread of maize. *Proc. Natl. Acad. Sci.* 109, E1913–E1921.  
943 doi:10.1073/pnas.1203189109
- 944 Hunt, L., White, J., Hoogenboom, G., 2001. Agronomic data: advances in documentation and protocols for  
945 exchange and use. *Agric. Syst.* 70, 477–492.
- 946 Irmak, A., Jones, J.W., Mavromatis, T., Welch, S.M., Boote, K.J., Wilkerson, G.G., 2000. Evaluating methods for  
947 simulating soybean cultivar responses using cross validation. *Agron. J.* 92, 1140–1149.
- 948 Jones, C.A., Richie, J.T., Kiniry, J.R., Godwin, D.C., 1986. Subroutine structure. CERES-Maize Simul. Model Maize  
949 Growth Dev. CA Jones JR Kiniry Contrib. PT Dyke AI.
- 950 Jones, J.W., Hoogenboom, G., Porter, C.H., Boote, K.J., Batchelor, W.D., Hunt, L., Wilkens, P.W., Singh, U.,  
951 Gijssman, A.J., Ritchie, J.T., 2003. The DSSAT cropping system model. *Eur. J. Agron.* 18, 235–265.
- 952 Keating, B.A., Carberry, P.S., Hammer, G.L., Probert, M.E., Robertson, M.J., Holzworth, D., Huth, N.I.,  
953 Hargreaves, J.N., Meinke, H., Hochman, Z., others, 2003. An overview of APSIM, a model designed for  
954 farming systems simulation. *Eur. J. Agron.* 18, 267–288.
- 955 Kennedy, J., 2011. Particle swarm optimization, in: *Encyclopedia of Machine Learning*. Springer, pp. 760–766.
- 956 Kiniry, J.R., Bonhomme, R., 1991. Predicting maize phenology. *Predict. Crop Phenol.* 11, 5–131.
- 957 Koduru, P., Welch, S.M., Das, S., 2007. A particle swarm optimization approach for estimating parameter  
958 confidence regions, in: *Proceedings of the 9th Annual Conference on Genetic and Evolutionary*  
959 *Computation*. ACM, pp. 70–77.
- 960 Laurila, H., Mäkelä, P., Kleemola, J., Peltonen, J., 2012. A comparative ideotype, yield component and cultivation  
961 value analysis for spring wheat adaptation in Finland. *Agric. Food Sci.* 21, 384–408.
- 962 Lee, M., Sharopova, N., Beavis, W.D., Grant, D., Katt, M., Blair, D., Hallauer, A., 2002. Expanding the genetic map  
963 of maize with the intermated B73  $\times$  Mo17 (IBM) population. *Plant Mol. Biol.* 48, 453–461.
- 964 Luo, Y., Weng, E., Wu, X., Gao, C., Zhou, X., Zhang, L., 2009. Parameter identifiability, constraint, and equifinality  
965 in data assimilation with ecosystem models. *Ecol. Appl.* 19, 571–574.
- 966 Major, D.J., Kiniry, J.R., 1991. Predicting daylength effects on phenological processes. *Predict. Crop Phenol.* 15–  
967 28.

- 968 Mascheretti, I., Turner, K., Brivio, R.S., Hand, A., Colasanti, J., Rossi, V., 2015. Florigen-encoding genes of day-  
969 neutral and photoperiod-sensitive maize are regulated by different chromatin modifications at the  
970 floral transition. *Plant Physiol.* 168, 1351–1363.
- 971 Matsumoto, M., Nishimura, T., 1998. Mersenne twister: a 623-dimensionally equidistributed uniform pseudo-  
972 random number generator. *ACM Trans. Model. Comput. Simul. TOMACS* 8, 3–30.
- 973 Mavromatis, T., Boote, K., Jones, J., Wilkerson, G., Hoogenboom, G., 2002. Repeatability of model genetic  
974 coefficients derived from soybean performance trials across different states. *Crop Sci.* 42, 76–89.
- 975 McMullen, M.D., Kresovich, S., Villeda, H.S., Bradbury, P., Li, H., Sun, Q., Flint-Garcia, S., Thornsberry, J.,  
976 Acharya, C., Bottoms, C., Brown, P., Browne, C., Eller, M., Guill, K., Harjes, C., Kroon, D., Lepak, N.,  
977 Mitchell, S.E., Peterson, B., Pressoir, G., Romero, S., Rosas, M.O., Salvo, S., Yates, H., Hanson, M., Jones,  
978 E., Smith, S., Glaubitz, J.C., Goodman, M., Ware, D., Holland, J.B., Buckler, E.S., 2009. Genetic Properties  
979 of the Maize Nested Association Mapping Population. *Science* 325, 737–740.  
980 doi:10.1126/science.1174320
- 981 Medlyn, B.E., Robinson, A.P., Clement, R., McMurtrie, R.E., 2005. On the validation of models of forest CO<sub>2</sub>  
982 exchange using eddy covariance data: some perils and pitfalls. *Tree Physiol.* 25, 839–857.
- 983 Meuwissen, T.H., Hayes, B.J., Goddard, M.E., 2001. Prediction of total genetic value using genome-wide dense  
984 marker maps. *Genetics* 157, 1819–1829.
- 985 Nelder, J.A., Mead, R., 1965. A simplex method for function minimization. *Comput. J.* 7, 308–313.
- 986 Press, W.H., Teukolsky, S.A., Vetterling, W.T., Flannery, B.P., 1992. *Numerical recipes in FORTRAN: the art of*  
987 *scientific computing.* Cambridge University Press Cambridge.
- 988 Quiring, S.M., Legates, D.R., 2008. Application of CERES-Maize for within-season prediction of rainfed corn  
989 yields in Delaware, USA. *Agric. For. Meteorol.* 148, 964–975.
- 990 Rabotyagov, S., Campbell, T., Valcu, A., Gassman, P., Jha, M., Schilling, K., Wolter, C., Kling, C., 2012. Spatial  
991 multiobjective optimization of agricultural conservation practices using a SWAT model and an  
992 evolutionary algorithm. *J. Vis. Exp. JoVE* e4009. doi:10.3791/4009
- 993 Reymond, M., Muller, B., Leonardi, A., Charcosset, A., Tardieu, F., 2003. Combining quantitative trait Loci  
994 analysis and an ecophysiological model to analyze the genetic variability of the responses of maize leaf  
995 growth to temperature and water deficit. *Plant Physiol.* 131, 664–675. doi:10.1104/pp.013839
- 996 Román-Paoli, E., Welch, S.M., Vanderlip, R.L., 2000. Comparing genetic coefficient estimation methods using  
997 the CERES-Maize model. *Agric. Syst.* 65, 29–41.
- 998 Saltelli, A., Annoni, P., Azzini, I., Campolongo, F., Ratto, M., Tarantola, S., 2010. Variance based sensitivity  
999 analysis of model output. Design and estimator for the total sensitivity index. *Comput. Phys. Commun.*  
1000 181, 259–270.
- 1001 Schwarz, G., 1978. Estimating the Dimension of a Model. *Ann. Stat.* 6, 461–464. doi:10.1214/aos/1176344136
- 1002 Semenov, M.A., Stratonovitch, P., 2013. Designing high-yielding wheat ideotypes for a changing climate. *Food*  
1003 *Energy Secur.* 2, 185–196.
- 1004 Stone, T., 2011. Sustainability and the needs of 2050 agriculture: Developed and developing world perspectives  
1005 (NABC No. 23), *Food Security: The Intersection of Sustainability, Safety and defense.*
- 1006 Tardieu, F., 2003. Virtual plants: modelling as a tool for the genomics of tolerance to water deficit. *Trends Plant*  
1007 *Sci.* 8, 9–14.
- 1008 Tardieu, F., Tuberosa, R., 2010. Dissection and modelling of abiotic stress tolerance in plants. *Curr. Opin. Plant*  
1009 *Biol.* 13, 206–212. doi:10.1016/j.pbi.2009.12.012
- 1010 Tatsumi, K., 2016. Effects of automatic multi-objective optimization of crop models on corn yield reproducibility  
1011 in the U.S.A. *Ecol. Model.* 322, 124–137. doi:10.1016/j.ecolmodel.2015.11.006
- 1012 Thompson, L.J., Ferguson, R.B., Kitchen, N., Frazen, D.W., Mamo, M., Yang, H., Schepers, J.S., 2015. Model and  
1013 Sensor-Based Recommendation Approaches for In-Season Nitrogen Management in Corn. *Agron. J.* 107,  
1014 2020–2030.

- 1015 Thorp, K.R., DeJonge, K.C., Kaleita, A.L., Batchelor, W.D., Paz, J.O., 2008. Methodology for the use of DSSAT  
1016 models for precision agriculture decision support. *Comput. Electron. Agric.* 64, 276–285.
- 1017 Thorp, K.R., Gore, M.A., Andrade-Sanchez, P., Carmo-Silva, A.E., Welch, S.M., White, J.W., French, A.N., 2015.  
1018 Proximal hyperspectral sensing and data analysis approaches for field-based plant phenomics. *Comput.*  
1019 *Electron. Agric.* 118, 225–236.
- 1020 Valentim, F.L., van Mourik, S., Posé, D., Kim, M.C., Schmid, M., van Ham, R.C., Busscher, M., Sanchez-Perez, G.F.,  
1021 Molenaar, J., Angenent, G.C., others, 2015. A quantitative and dynamic model of the Arabidopsis  
1022 flowering time gene regulatory network. *PLoS One* 10, e0116973.
- 1023 Wallach, D., Goffinet, B., Bergez, J.-E., Debaeke, P., Leenhardt, D., Aubertot, J.-N., 2001. Parameter Estimation  
1024 for Crop Models. *Agron. J.* 93, 757. doi:10.2134/agronj2001.934757x
- 1025 Welch, S.M., Dong, Z., Roe, J.L., Das, S., 2005a. Flowering time control: gene network modelling and the link to  
1026 quantitative genetics. *Crop Pasture Sci.* 56, 919–936.
- 1027 Welch, S.M., Roe, J.L., Das, S., Dong, Z., He, R., Kirkham, M.B., 2005b. Merging genomic control networks and  
1028 soil-plant-atmosphere-continuum models. *Agric. Syst.* 86, 243–274.
- 1029 Welch, S.M., Wilkerson, G., Whiting, K., Sun, N., Vagts, T., Buol, G., Mavromatis, T., 2002. Estimating soybean  
1030 model genetic coefficients from private–sector variety performance trial data. *Trans. ASAE* 45, 1163.
- 1031 Welch, S.M., Zhang, J., Sun, N., Mak, T.Y., 2000. Efficient estimation of genetic coefficients of crop models, in:  
1032 *The Third International Symposium on System Approaches for Agricultural Development.*
- 1033 White, J.W., Hoogenboom, G., 2010. Crop response to climate: ecophysiological models, in: *Climate Change and*  
1034 *Food Security.* Springer, pp. 59–83.
- 1035 White, J.W., Hoogenboom, G., 1996. Simulating effects of genes for physiological traits in a process-oriented  
1036 crop model. *Agron. J.* 88, 416–422.
- 1037 Wilczek, A.M., Roe, J.L., Knapp, M.C., Cooper, M.D., Lopez-Gallego, C., Martin, L.J., Muir, C.D., Sim, S., Walker,  
1038 A., Anderson, J., Egan, J.F., Moyers, B.T., Petipas, R., Giakountis, A., Charbit, E., Coupland, G., Welch,  
1039 S.M., Schmitt, J., 2009. Effects of genetic perturbation on seasonal life history plasticity. *Science* 323,  
1040 930–934. doi:10.1126/science.1165826
- 1041 Wit, C.T. de., 1965. Photosynthesis of leaf canopies, *Agricultural Research Reports*; 663; *Verslagen van*  
1042 *Landbouwkundige Onderzoekingen*; 663. Centre for Agricultural Publications and Documentation,  
1043 Wageningen.
- 1044 Yin, X., Stam, P., Kropff, M.J., Schapendonk, A.H., 2003. Crop modeling, QTL mapping, and their complementary  
1045 role in plant breeding. *Agron. J.* 95, 90–98.
- 1046