

How well do crop models predict phenology, with emphasis on the effect of calibration?

1 Wallach¹, Daniel; Palosuo², Taru; Thorburn³, Peter; Seidel⁴, Sabine J.; Gourdain⁵, Emmanuelle; Asseng⁶,
2 Senthold; Basso⁷, Bruno; Buis⁸, Samuel; Crout⁹, Neil, Dibari¹⁰, Camilla; Dumont¹¹, Benjamin; Ferrise¹⁰,
3 Roberto; Gaiser⁴, Thomas; Garcia⁷, Cécile; Gayler¹², Sebastian; Ghahramani¹³, Afshin; Hochman³, Zvi;
4 Hoek¹⁴, Steven; Horan³, Heidi; Hoogenboom^{6,15}, Gerrit; Huang¹⁶, Mingxia; Jabloun⁹, Mohamed; Jing¹⁷,
5 Qi; Justes¹⁸, Eric; Kersebaum¹⁹, Kurt Christian; Klosterhalfen²⁰, Anne; Launay²¹, Marie; Luo²², Qunying;
6 Maestrini⁷, Bernardo; Moriondo²³, Marco; Nariman Zadeh²⁴, Hasti; Olesen²⁵, Jørgen Eivind; Poyda²⁶,
7 Arne; Priesack²⁷, Eckart; Pullens²⁵, Johannes Wilhelmus Maria; Qian¹⁷, Budong; Schütze²⁸, Niels;
8 Shelia^{6,15}, Vakhtang; Souissi^{29,30}, Amir; Specka¹⁹, Xenia; Srivastava⁴, Amit Kumar; Stella¹⁹, Tommaso;

¹ INRA, UMR AGIR, Castanet Tolosan, France. ORCID 0000-0003-3500-8179

² Natural Resources Institute Finland (Luke), Helsinki, Finland

³ CSIRO Agriculture and Food, Brisbane, Queensland, Australia

⁴ Institute of Crop Science and Resource Conservation, University of Bonn, Germany

⁵ ARVALIS - Institut du végétal Paris, France

⁶ Agricultural and Biological Engineering Department, University of Florida, Gainesville, Florida

⁷ Department of Earth and Environmental Sciences, Michigan State University, East Lansing, Michigan

⁸ INRA, UMR 1114 EMMAH, Avignon, France

⁹ School of Biosciences, University of Nottingham, Loughborough, UK

¹⁰ Department of Agriculture, Food, Environment and Forestry (DAGRI), University of Florence, Italy

¹¹ Department Terra & AgroBioChem, Gembloux Agro-Bio Tech, University of Liege, Gembloux, Belgium

¹² Institute of Soil Science and Land Evaluation, Biogeophysics, University of Hohenheim, Stuttgart, Germany

¹³ Centre for Sustainable Agricultural Systems, Institute for Life Sciences and the Environment, University of Southern Queensland, Toowoomba, Queensland, Australia

¹⁴ Environmental Sciences Group, Wageningen University & Research, Wageningen, The Netherlands

¹⁵ Institute for Sustainable Food Systems, University of Florida, Gainesville, Florida

¹⁶ College of Resources and Environmental Sciences, China Agricultural University, Beijing, China

¹⁷ Ottawa Research and Development Centre, Agriculture and Agri-Food Canada, Ottawa, Canada

¹⁸ CIRAD, UMR SYSTEM, Montpellier, France

¹⁹ Leibniz Centre for Agricultural Landscape Research, Müncheberg, Germany

²⁰ Institute for Bio- and Geosciences - IBG-3, Agrosphere, Forschungszentrum Jülich GmbH, Jülich, Germany

²¹ INRA, US 1116 AgroClim, Avignon, France

²² Hillridge Technology Pty Ltd, Sydney, Australia

²³ CNR-IBIMET, Firenze, Italy

²⁴ Aalto University School of Science, Espoo, Finland

²⁵ Department of Agroecology, Aarhus University, Tjele, Denmark

²⁶ Grass and Forage Science / Organic Agriculture, Institute of Crop Science and Plant Breeding, Kiel University, Kiel, Germany

²⁷ Institute of Biochemical Plant Pathology, Helmholtz Zentrum München-German Research Center for Environmental Health, Neuherberg, Germany

²⁸ Institute of Hydrology and Meteorology, Chair of Hydrology, Technische Universität Dresden, Dresden, Germany

12 Streck¹², Thilo; Trombi¹⁰, Giacomo; Wallor¹⁹, Evelyn; Wang¹⁶, Jing; Weber¹², Tobias, K.D.;

13 Weihermüller²⁰, Lutz; de Wit¹⁴, Allard; Wöhling^{28,31}, Thomas; Xiao^{32,6}, LiuJun; Zhao⁶, Chuang; Zhu³²²,

14 Yan

²⁹ National Institute of Agronomic Research of Tunisia (INRAT), Agronomy Laboratory, University of Carthage, Tunis, Tunisia

³⁰ National Agronomy Institute of Tunisia (INAT), University of Carthage, Tunis, Tunisia

³¹ Lincoln Agritech Ltd., Hamilton, New Zealand

³² National Engineering and Technology Center for Information Agriculture, Jiangsu Key Laboratory for Information Agriculture, Jiangsu Collaborative Innovation Center for Modern Crop Production, Nanjing Agricultural University, Nanjing, Jiangsu, China

15

16 ABSTRACT

17 Plant phenology, which describes the timing of plant development, is a major aspect of
18 plant response to environment and for crops, a major determinant of yield. Many studies have
19 focused on comparing model equations for describing how phenology responds to climate but
20 the effect of crop model calibration, also important for determining model performance, has
21 received much less attention. The objectives here were to obtain a rigorous evaluation of
22 prediction capability of wheat phenology models, to analyze the role of calibration and to
23 document the various calibration approaches. The 27 participants in this multi-model study
24 were provided experimental data for calibration and asked to submit predictions for sites and
25 years not represented in those data. Participants were instructed to use and document their
26 “usual” calibration approach. Overall, the models provided quite good predictions of
27 phenology (median of mean absolute error of 6.1 days) and did much better than simply using
28 the average of observed values as predictor. The results suggest that calibration can
29 compensate to some extent for different model formulations, specifically for differences in
30 simulated time to emergence and differences in the choice of input variables. Conversely,
31 different calibration approaches were associated with major differences in prediction error
32 between the same models used by different groups. Given the large diversity of calibration
33 approaches and the importance of calibration, there is a clear need for guidelines and tools to
34 aid with calibration. Arguably the most important and difficult choice for calibration is the
35 choice of parameters to estimate. Several recommendations for calibration practices are
36 proposed. Model applications, including model studies of climate change impact, should
37 focus more on the data used for calibration and on the calibration methods employed.

38 Introduction

39 Crop models are widely used to simulate the effect of weather, soil and management
40 on crops (Rauff & Bello, 2015; van Ittersum et al., 2003). Here we focus specifically on the
41 use of crop models to simulate crop phenology i.e. the cycle of biological events in plants.
42 Matching the phenology of crop varieties to the climate in which they grow is a critical crop
43 production strategy (Hunt et al., 2019; Rezaei, Siebert, & Ewert, 2015; Rezaei, Siebert,
44 Hüging, & Ewert, 2018). Thus, understanding and improving our ability to simulate
45 phenology with crop models is important step in using models for improving crop
46 management, for designing better adapted genotypes and for preparing for and adapting to
47 global change. Process-based models similar to those for crops can be used for natural
48 vegetation, so crop models can serve as examples for studies of phenology in ecosystems
49 (Piao et al., 2019).

50 Crop model evaluation is an essential aspect of modeling, assessing whether model
51 performance is acceptable for the intended use of the model. For studies of phenology two
52 major questions are a) how accurate are current models for the prediction of crop
53 development stages? and b) what determines model accuracy and what does that imply about
54 how accuracy can be improved? We use here prediction in the sense of determining outputs
55 (dates of development stages) from known inputs (weather, soil, management). The problem
56 of predicting future events, with unknown weather, is not considered.

57 There have been numerous evaluation studies of crop model simulations, including but
58 not restricted to phenology, both of individual models and of multi-model ensembles. The
59 typical procedure is to first calibrate each model using a part of the available field data and
60 then to evaluate it using the remaining data.

61 Most crop model evaluation studies focusing on crop phenology have had relatively
62 little data for calibration or evaluation. (Andarzian, Hoogenboom, Bannayan, Shirali, &
63 Andarzian, 2015) for example, used data from one location covering five growing seasons and
64 two or three sowing dates per year. Out of these data, one year was used for calibration and
65 the other two years of data to evaluate the model. (Yuan, Peng, & Li, 2017) used one year of
66 data for calibration and the second year of data from the same location for evaluation of the
67 rice crop model ORYZA. Hussain, Khaliq, Ahmad, & Akhtar (2018) tested four models using
68 data from two locations with two years of data, 11 crop planting dates, and three varieties.
69 Paucity of data means that model parameters are estimated with relatively large uncertainty
70 and model evaluation is quite uncertain.

71 Another common feature of crop model evaluation is that the data are often such that
72 model error for the evaluation data cannot be assumed to be independent of model error for
73 the calibration data. That holds for the examples listed above since the evaluation and
74 calibration data come from the same sites. In such cases, the evaluation does not give an
75 unbiased estimate of how well the model will predict for other sites not included in the
76 calibration data. Since usually the model is meant for use over a range of sites, this clearly
77 reduces the usefulness of the evaluation information.

78 A third feature often found in crop model evaluation is that the range of situations
79 from which the calibration data are drawn (the “calibration population”) is often different than
80 the range of conditions from which the evaluation data are drawn (the “evaluation
81 population”). For example, Hussain et al. (2018) used data from an experiment that included a
82 range of crop stresses to calibrate their model. They used data from the least stressed
83 treatment in the calibration process and evaluated the resultant model on the remaining
84 planting dates at the same location. The evaluation data thus represented a different range of

85 conditions than the calibration data. In a multi-model ensemble study of the effect of high
86 temperatures on wheat growth (Asseng et al., 2015) detailed crop measurements were
87 provided for one planting date and the models were evaluated using other planting dates,
88 some with additional artificial heating. Again, the evaluation data represented a much larger
89 range of temperatures than represented in the calibration data. While the capacity of crop
90 models to extrapolate to conditions quite different than those of the calibration data is
91 obviously of interest, it is a rather different type of evaluation than the case where the
92 calibration and evaluation populations are similar.

93 Thus, evaluation of crop phenology models to date has mainly concerned situations
94 that would tend to make prediction difficult, because of small amounts of data for calibration
95 and differences between the calibration and target populations. Furthermore, the quality of the
96 evaluation is often questionable, because of the relatively small amounts of data and the non-
97 independence of the errors for the calibration and evaluation situations. There is thus a need
98 for more rigorous assessments of simulation capability of crop phenology models. The first
99 objective of this paper is, therefore, to provide a rigorous evaluation of how well crop models
100 predict wheat phenology, in the situation where there is substantial data for calibration and
101 where the calibration and evaluation data can be assumed to come from the same underlying
102 population. To ensure the rigor of the evaluation, we create a situation where the model errors
103 for the calibration and simulation data can be assumed independent.

104 The emphasis in model evaluation studies is often on the role of model structure, i.e.
105 model equations (Maiorano et al., 2017; Svystun, Bhalerao, & Jönsson, 2019; Wang et al.,
106 2017). There has been relatively little work on the diversity and importance of calibration
107 approaches in crop modeling. Clearly however the simulated values also depend on the
108 parameter values estimated by calibration and therefore on the calibration approach.

109 Confalonieri et al., (2016) found that the model user, responsible for calibration, had a very
110 large effect on predictive quality. The second objective of this study then was to investigate
111 the role of calibration in determining prediction quality.

112 In a wide-ranging survey, (Seidel, Palosuo, Thorburn, & Wallach, 2018) found that
113 there is a wide diversity of calibration strategies used for crop models, but for that survey
114 each response was for a different prediction problem. This did not address the problem of the
115 diversity of calibration approaches by different groups given the same data and the same
116 prediction objectives. The third objective of this study was therefore to obtain detailed
117 information about the diversity of calibration strategies adopted by different modeling groups
118 for the same prediction problem. This is useful as a step toward developing guidelines for
119 calibration of phenology models, in that it helps define the range of possible approaches. This
120 is of practical interest not only for stand-alone phenology models, but also for crop models
121 more generally, since crop models are often calibrated first just for phenology, and then
122 separately for biomass increase and partitioning and soil processes.

123 **Materials and Methods**

124 **Experimental data**

125 The data were provided by ARVALIS – Institut du végétal, a French agricultural
126 technical institute. They run multi-year multi-purpose trials at multiple locations across
127 France, which include variety trials. The data here are from the two winter wheat check
128 varieties, *Apache* which is a common variety grown throughout France and Central Europe
129 and *Bermude*, mainly grown in Northern and Central France. The trials have three repetitions
130 and follow standard agricultural practices, with N fertilization calculated to be non-limiting.
131 Thus, both the calibration and evaluation data are drawn from sites in France where winter
132 wheat is grown, subject to standard management.

133 The observed data used in model calibration and evaluation are the dates of two
 134 development stages, namely beginning of stem elongation (growth stage 30 on the BBCH and
 135 Zadoks scales (Zadoks, Chzang, & Konzak, 1974) and middle of heading (growth stage 55 on
 136 the BBCH and Zadoks scales). These stages are of practical importance because they can
 137 easily be determined visually and are closely related to the recommended dates for the second
 138 and third nitrogen fertilizer applications.

139 The data were divided into two categories (table 1). One part, the calibration data (six
 140 sites, five years for a total of 14 environments i.e. site-year combinations), was provided to
 141 participants for calibration. A second part, the evaluation data (five sites, two years for a total
 142 of eight environments), was not given to participants. The division of the data was such that
 143 the calibration and evaluation data had no sites or years in common. To achieve this some
 144 data (denoted “other” data) were used neither for calibration nor evaluation (but were used in
 145 the calculation where overall variability in simulated values was evaluated).

146 **Table 1.**

147 **Environments (site-year combinations) that provided the data. C= calibration data. E =**
 148 **evaluation data. O = data not used for calibration or evaluation (only used for**
 149 **evaluating variability between models). Blank cells indicate no data.**

Site (longitude,latitude)	Harvest year						
	2010	2011	2012	2013	2014	2015	2016
FORESTE (3.20,49.82)			E	E	OO*	O	
MERY	C	C		O	C	C	

4.02,48.33)							
ROUVRES 5.09,47.28)			E	E	O	O	
CESSEVILLE ¹ (0.90,49.15)		C					
IVILLE ¹ (0.90,49.15)			E				
VILLETES ¹ (0.90,49.15)				E			
EPREVILLE ¹ (0.90,49.15)					C		
CRESTOT ¹ (0.90,49.15)						C	
OUZOUER (1.52,47.90)		O	E	E	O	O	
BIGNAN (-2.73,47.88)	C	C	O	O	C	C	
BOIGNEVILLE (2.38,48.33)			O	O	C	C	C

151 *There were two sowing dates at FORESTE in 2014.¹ These are separate sites that are
152 geographically close to one another and share a single weather station.

153

154 The background and input information provided to the modelers for all environments
155 included information about the sites (location, soil texture, field capacity, wilting point),
156 management (sowing dates, sowing density, irrigation and fertilization dates and amounts),
157 and daily weather data (precipitation, minimum and maximum air temperature, global
158 radiation and potential evapotranspiration). Initial soil water and N content were not measured
159 in these experiments, but best estimates were provided by the experimental scientist. If any
160 models required other input data, modeling groups were asked to derive those values in
161 whatever way that seemed appropriate.

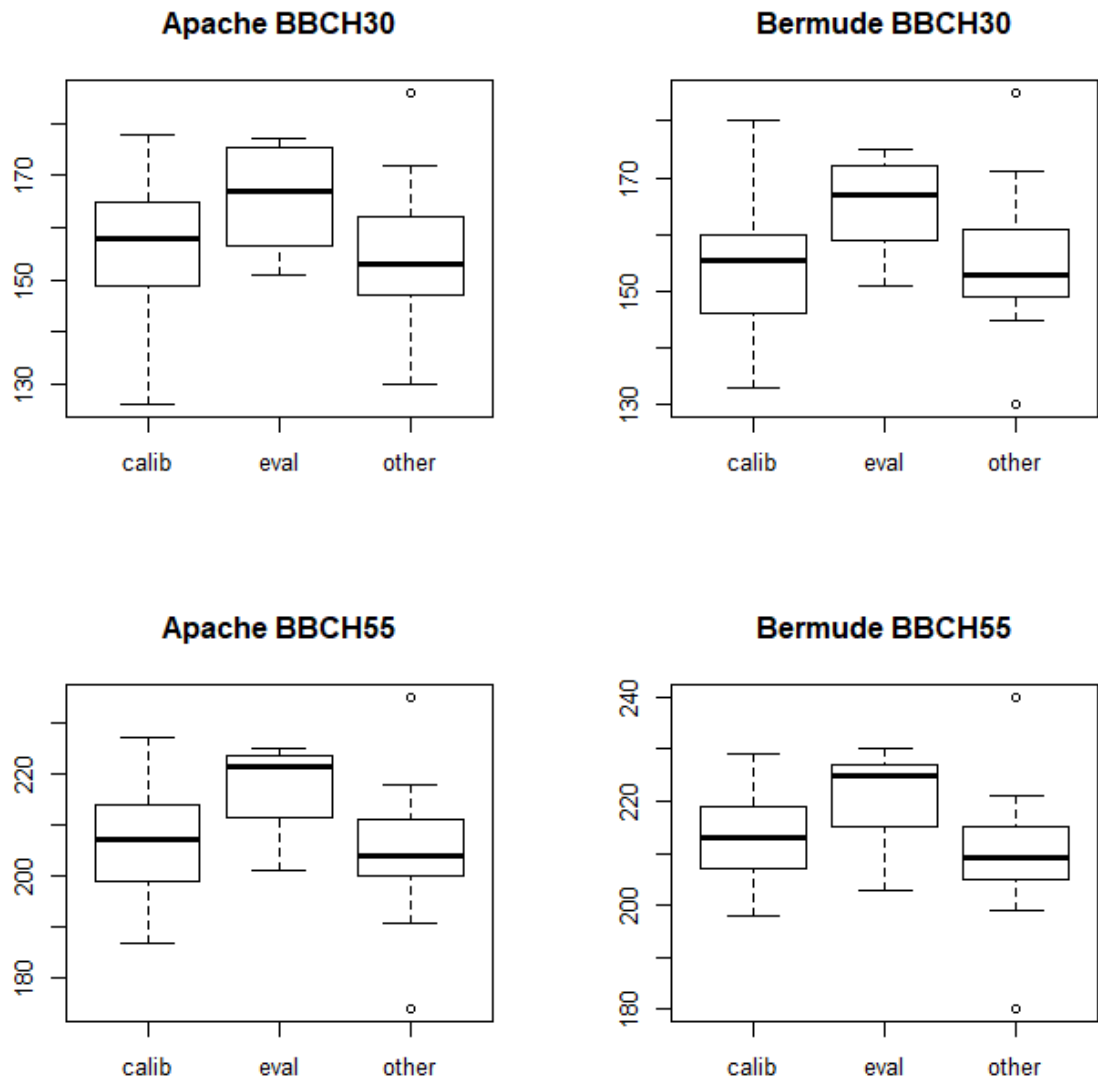
162 The range of observed days from sowing to development stages BBCH30 and
163 BBCH55 for the two varieties for each category of data (calibration, evaluation, other hidden
164 data) is shown in figure 1. The spread from minimum to maximum in the evaluation data is
165 between 24 and 27 days depending on stage and variety. The spread is larger for the
166 calibration data, and in fact, the calibration data cover the range of the evaluation data and the
167 range of other hidden data. Thus, the models are not being used to extrapolate outside the
168 observed values of the calibration data.

169

170 **Figure 1**

171 **Boxplots of calibration, evaluation and other data for development stages**
172 **BBCH30 and BBCH50 and varieties Apache (left) and Bermude (right). The y-axis**
173 **shows days from sowing to the indicated development stage. Boxes indicate the lower**
174 **and upper quartiles. The solid line within the box is the median. Whiskers indicate the**

175 **most extreme data point, which is no more than 1.5 times the interquartile range from**
176 **the box, and the outlier dots are those observations that go beyond that range.**



177
178

179

180 **Crop models**

181 Twenty-seven modeling groups participated in this study, noted M1-M27. Information
182 about the underlying model structures is given in Supplementary table S1. The four groups
183 M2, M3, M4, and M5 all used the same model structure (i.e. models with the same name),

184 denoted as model structure S1. The four models M7, M12, M13, and M25 also shared a
185 common model structure, denoted as S2. As will be seen, different groups using the same
186 model structure had different results. This could be due to differences in model version, but
187 the differences are not in the basic phenology equations, and therefore, should have no or only
188 a negligible effect on the simulated development stages. The differences are assumed to
189 mainly be due to differences in the values of the parameters, either those not fit by calibration
190 or those estimated by calibration. Since groups using the same model structure obtained
191 different results, we refer to the 27 contributions as different models. In the presentation of the
192 results the models are anonymized and are identified simply as M1 to M27. It would be
193 misleading to use the names of the model structures, since the results depend on both model
194 structure and on the values of the parameters.

195 Two of the models (M9, M18) only simulated days to development stage BBCH55
196 and not to stage BBCH30. Results for these two models are systematically included with the
197 results for the other models, but averages over development stages for these two models only
198 refer to BBCH55. This is not repeated explicitly every time an average over development
199 stages is discussed.

200 In addition to the individual model results, we show the results for the model ensemble
201 mean (“e-mean”) and the model ensemble median (“e-median”). We also define a very simple
202 predictor, denoted “naive”, which was calculated as the average of the observations in the
203 calibration data for prediction. The naive model thus predicts that all days from sowing to
204 stage BBCH30 (BBCH55) will correspond to the average of days from sowing to BBCH30
205 (BBCH55) in the calibration data, separately for each variety. The naive model predictions for
206 days from sowing to BBCH30 and BBCH55 are respectively 155.9 days and 206.9 days for
207 variety Apache, and 156.1 days and 213.1 days for variety Bermude.

208 Calibration and simulation experiment

209 The participants were provided with observed phenology data (dates of BBCH30 and
210 BBCH55) only for the calibration environments. The participants were asked to calibrate their
211 model using those data, and then to use the calibrated model to simulate phenology for all
212 environments (i.e. calibration, evaluation and hidden data environments). No guidelines for
213 calibration were provided. Participants were instructed to calibrate their model in their “usual
214 way” and fill out a questionnaire explaining what they did (Supplementary table S2).

215 Evaluation

216 A common metric of error is mean squared error (MSE). We calculated MSE for each
217 model, each development stage (BBCH30 and BBCH55) and for each variety, as well as
218 averaged over stages and varieties. This was done separately for the calibration and evaluation
219 data. For example, MSE for model m , for predicting BBCH30, variety Apache, based on the
220 evaluation data, is:

$$221 \quad MSE_{eval,m}^{BBCH30,Apache} = (1/8) \sum_{i \in eval} \left(y_i^{BBCH30,Apache} - \hat{y}_{i,m}^{BBCH30,Apache} \right)^2 \quad (1)$$

222

223 where the sum is over the eight environments used for evaluation and $y_i^{BBCH30,Apache}$ and

224 $\hat{y}_{i,m}^{BBCH30,Apache}$ are respectively the observed value and value simulated by model m for

225 evaluation environment i , development stage BBCH30 and variety Apache. For $MSE_{eval,m}^{all}$, the

226 average is over the eight evaluation environments, both stages and both varieties, so overall

227 32 predictions.

228 Mean squared error can be shown to be the sum of three positive terms, namely

229 squared bias, the difference in variance between the observed and simulated values and a term

230 related to the correlation between observed and simulated values (Kobayashi & Salam, 2000).
231 We specifically examined the bias, defined as the average over observed values minus the
232 average over simulated values.

233 The mean absolute error (MAE) is of interest as a more direct measure of error, that
234 does not give extra weight to large errors as MSE does. For example, MAE for model m for
235 predicting BBCH30, variety Apache, based on the evaluation data, is:

$$236 \quad MAE_{eval,m}^{BBCH30,Apache} = (1/8) \sum_{i \in eval} |y_i^{BBCH30,Apache} - \hat{y}_{i,m}^{BBCH30,Apache}|$$

237

238 We also look at modeling efficiency (EF) defined for model m as

$$239 \quad EF_m = 1 - MSE_m / MSE_{naive}$$

240 where MSE_m is MSE for model m and MSE_{naive} is MSE for the naive model defined above.

241 EF is a skill measure, which compares the predictive capability of a model to that of the naive
242 model. Since the naive model makes the same prediction for all environments, it does not
243 account for any of the variability between environments. A model with $EF \leq 0$ is a model that
244 does no better than the naive model, and so would be considered to be a very poor predictor.
245 A perfect model, with no error, has modeling efficiency of 1.

246

247 Results

248 Goodness-of-fit and prediction error

249 Summary statistics for MSE averaged over both varieties and over both development
250 stages, for the calibration and evaluation data, are shown in table 2. Summary MSE values for

251 the calibration data for each development stage and variety separately are shown in
252 Supplementary table S7, and results for each individual model are given in Supplementary
253 figure S1.

254 **Table 2**

255 **Summary statistics of MSE (days²) averaged over both varieties and over both**
256 **development stages.**

MSE (days ²)	Minimum	1st quartile	Median	Mean	3rd quartile	Maximum
Calibration data	15	28	47	77	63	426
Evaluation data	20	35	62	79	111	235

257

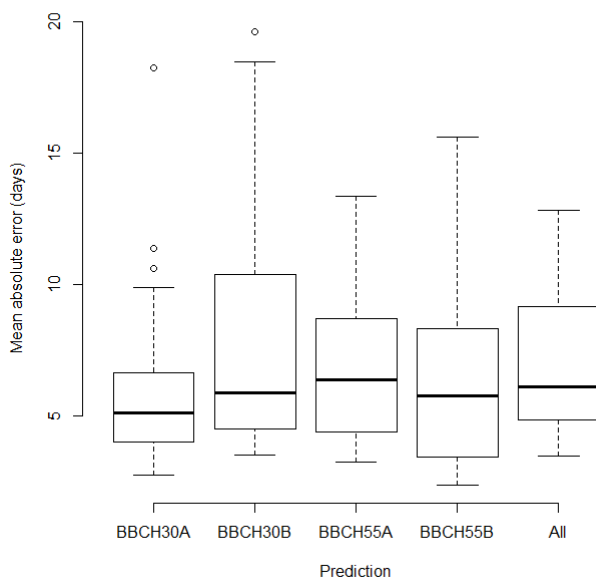
258 Figure 3 and Supplementary tables S4-S6 show results for each development stage and
259 variety and averaged over development stages and varieties for the evaluation data. Results
260 for each model are given in Supplementary table S3. The median of MAE for the evaluation
261 data is 6.1 days. The median of overall efficiency is 0.62, signifying that half of the models
262 have MSE values for the evaluation data that are at most 38% as large as that of the naive
263 predictor. Only two models have negative values of EF, indicating that one would do better
264 to predict using the average of the calibration data. For the four individual predictions (two
265 development stages, two varieties), the median of MAE ranges from 5.1 to 6.4 and the median
266 of EF ranges from 0.6 to 0.8. The ensemble models e-median and e-mean, though not the best
267 predictors, are among the best, with e-median being rated second best and e-mean fourth best.

268 The range of results among individual models is appreciable. The mean absolute errors for the
269 evaluation data averaged over all predictions (MAE_{eval}^{all}) go from 3.5 to 13 days. The MSE_{eval}^{all}
270 values vary by over a factor of 10, from a minimum of 20 days² to a maximum of 235 days².

271 **Figure 3**

272 **Box and whisker diagrams of absolute errors for evaluation data for each**
273 **prediction and on average (top panel) and modeling efficiency for each prediction and**
274 **on average (bottom panel). BBCH30A and BBCH30B refer respectively to prediction of**
275 **days to BBCH30 for variety Apache and variety Bermude. BBCH55A and BBCH55B**
276 **refer respectively to prediction of days to BBCH55 for variety Apache and variety**
277 **Bermude. The variability comes from differences between models.**

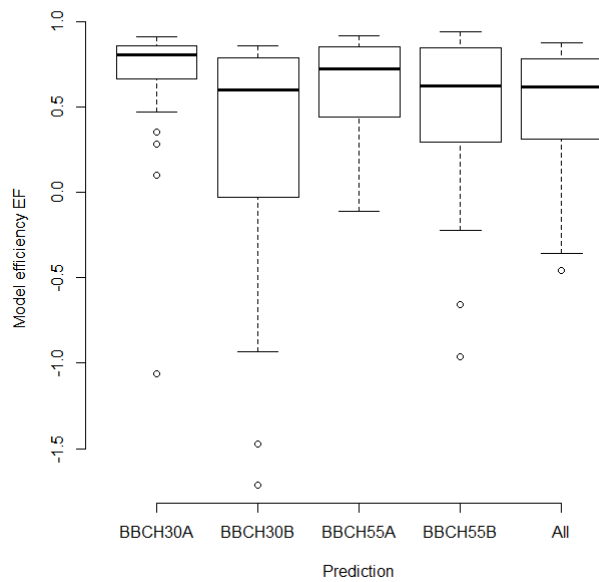
278



279

280

281



282

283

284

285

286 **Role of calibration**

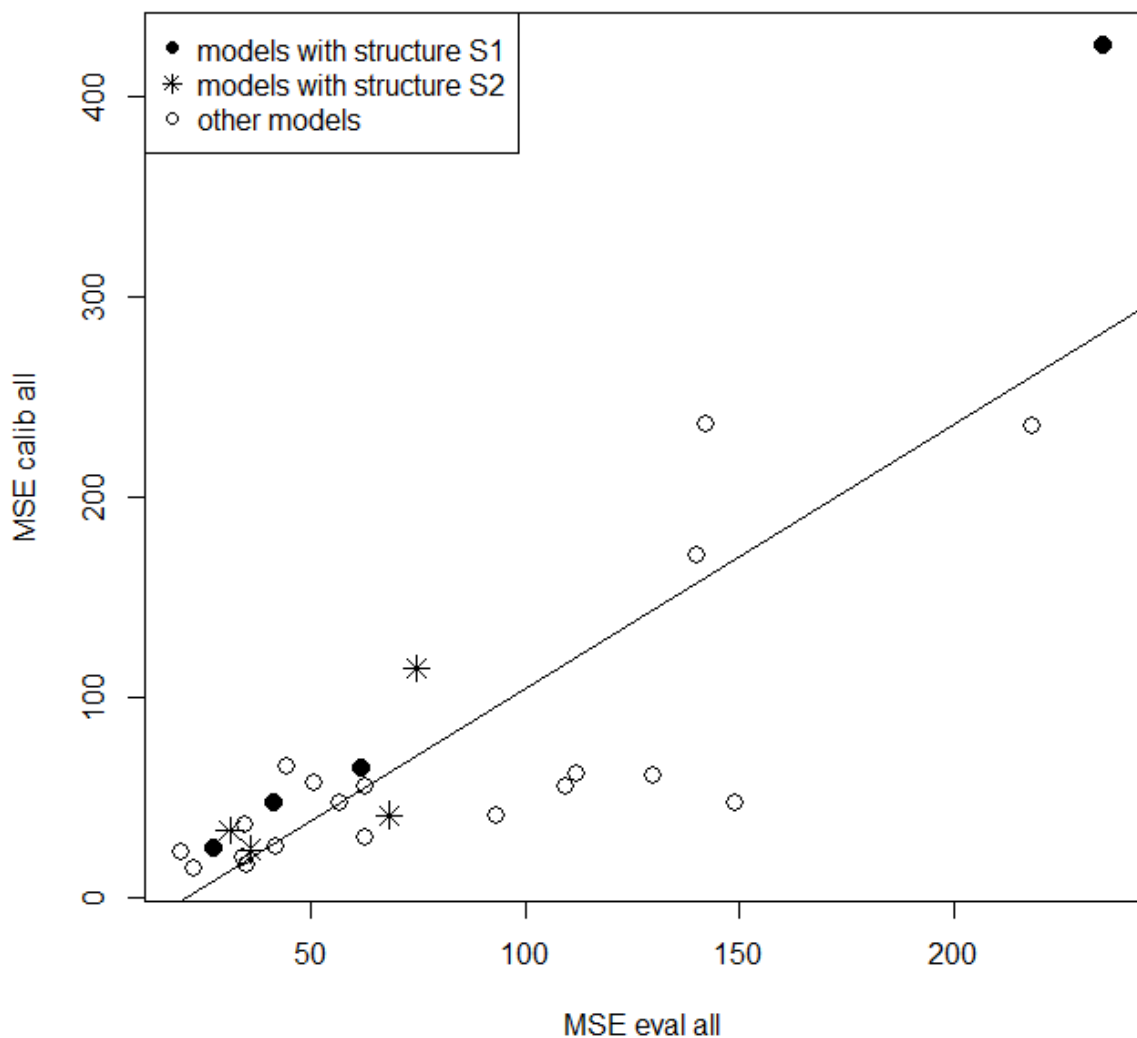
287 The relationship between overall MSE for the evaluation data and overall MSE for the
288 calibration data is quite close (adjusted $R^2=0.70$, figure 4). That is, much of the variability
289 between models in MSE for the evaluation data can be explained by the variability in MSE
290 values for the calibration data, which emphasizes the importance of obtaining a good fit to the
291 calibration data, which in turn depends to a large extent on the method of calibration.

292 The four models that have model structure S1 and the four models that have model
293 structure S2 are identified in figure 4. Models with the same structure have different MSE
294 values; the differences are particularly large for S1. The models with structure S1 are ranked
295 3rd, 9th, 14th and 27th best for overall evaluation MSE among the 27 individual models. The
296 models with structure S2 are ranked 4th, 8th, 17th and 18th best.

297

298 **Figure 4**

299 **Mean squared error (MSE) for the calibration data, averaged over environments,**
300 **development stages and varieties (MSE_{calib}^{all} days²), as related to MSE for the evaluation**
301 **data (MSE_{eval}^{all} , days²). The regression line $MSE_{calib}^{all} = -27.6 + 1.32 MSE_{eval}^{all}$ is shown**
302 **($R^2=0.70$). ● indicates models with structure S1. * indicates models with structure S2. ○**
303 **indicates other models.**



304

305

306 Twenty-one models simulated and reported time from sowing to emergence. For these
307 models, we can separate simulated time from sowing to BBCH30 (sow_30) into two
308 contributions, the simulated time from sowing to emergence (sow_em) plus the simulated
309 time from emergence to BBCH30 (em_30), so that $sow_30 = sow_em + em_30$. Figure 5 shows
310 results from two environments, typical of essentially all environments and both varieties, for
311 the relation between em_30 or sow_30 and sow_em. The average slope of the regression
312 $em_30 = a + b * sow_em$ over all environments (including calibration, evaluation and other
313 environments) and both varieties is $b = -1.04$, so that each day increase in simulated days to
314 emergence is on average associated with a 1.04 day decrease in simulated time from
315 emergence to BBCH30. The negative correlation between sow_em and em_30 leads to a
316 between-model variance for sow_30 (average variance 92 days²) that is smaller than the sum
317 of the variances of sow_em (average variance 20 days²) and em_30 (average variance 101
318 days²). The right panels of figure 5 show that different models could simulate almost exactly
319 the observed value of sow_30 with quite different values of sow_em.

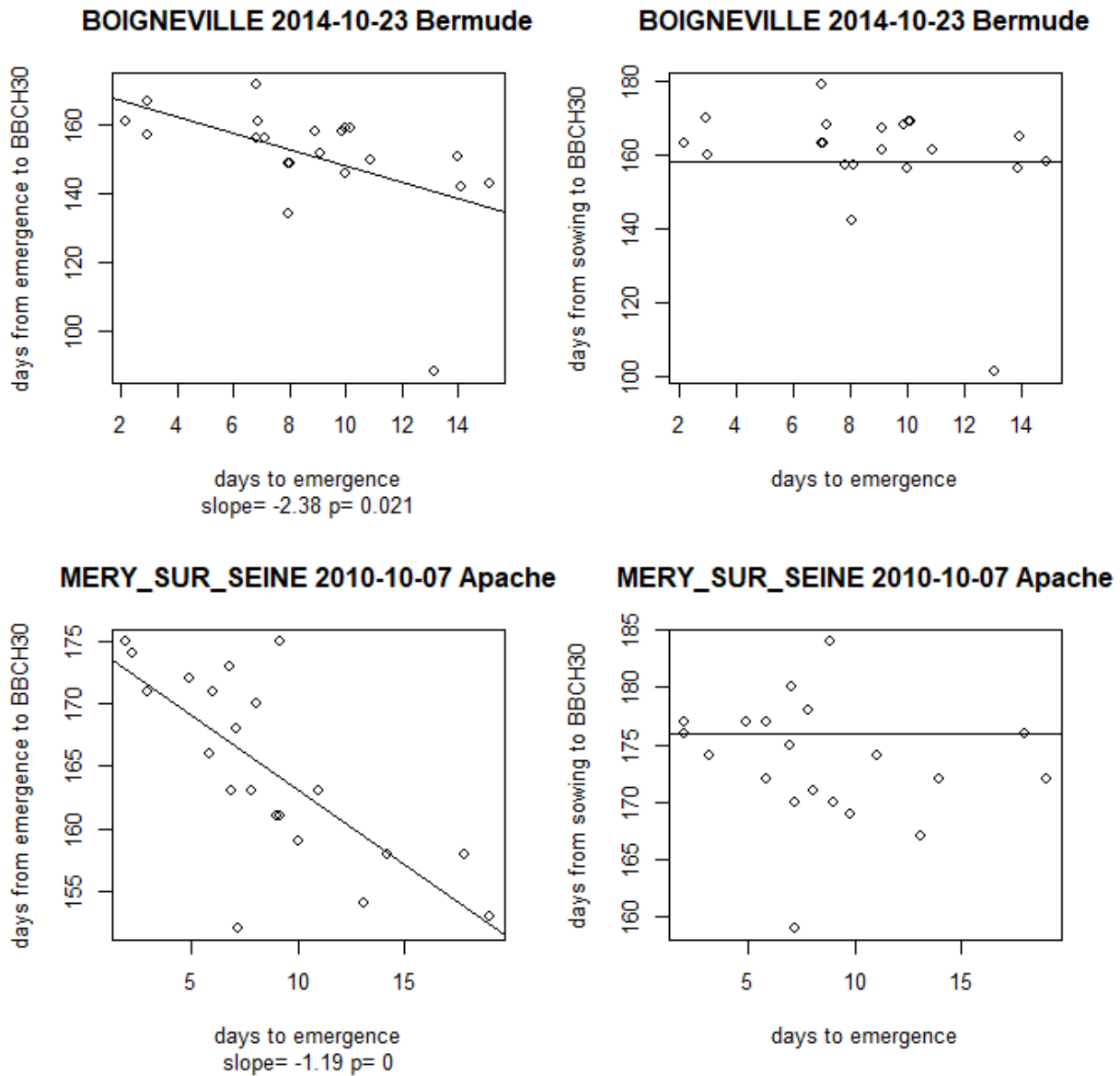
320

Figure 5

321 **Relation between simulated days from emergence to BBCH30 and simulated days**
322 **from sowing to emergence as reported by 21 crop models for two environments (left**
323 **panels). Relation between simulated days from sowing to BBCH30 and simulated days**
324 **from sowing to emergence for the same environments (right panels). A small amount of**
325 **noise has been applied to avoid overlap. The slope of the linear regression line and the p-**
326 **value for testing slope=0 are shown for the left panels. The observed days from sowing to**
327 **BBCH30 is shown as a horizontal line in the right panels.**

328

329



330

331 Bias (average over environments of observed values – average of simulated values) is
332 one aspect of goodness-of-fit. For most models, the bias for the calibration data is quite small.
333 Considering absolute bias for both development stages and both varieties, the median value
334 over models was 2 days (Supplementary table S8). In cases where the bias is relatively large,
335 it is often of opposite sign for BBCH30 and BBCH55, as in the examples of figure 2.

336

337 **Figure 2**

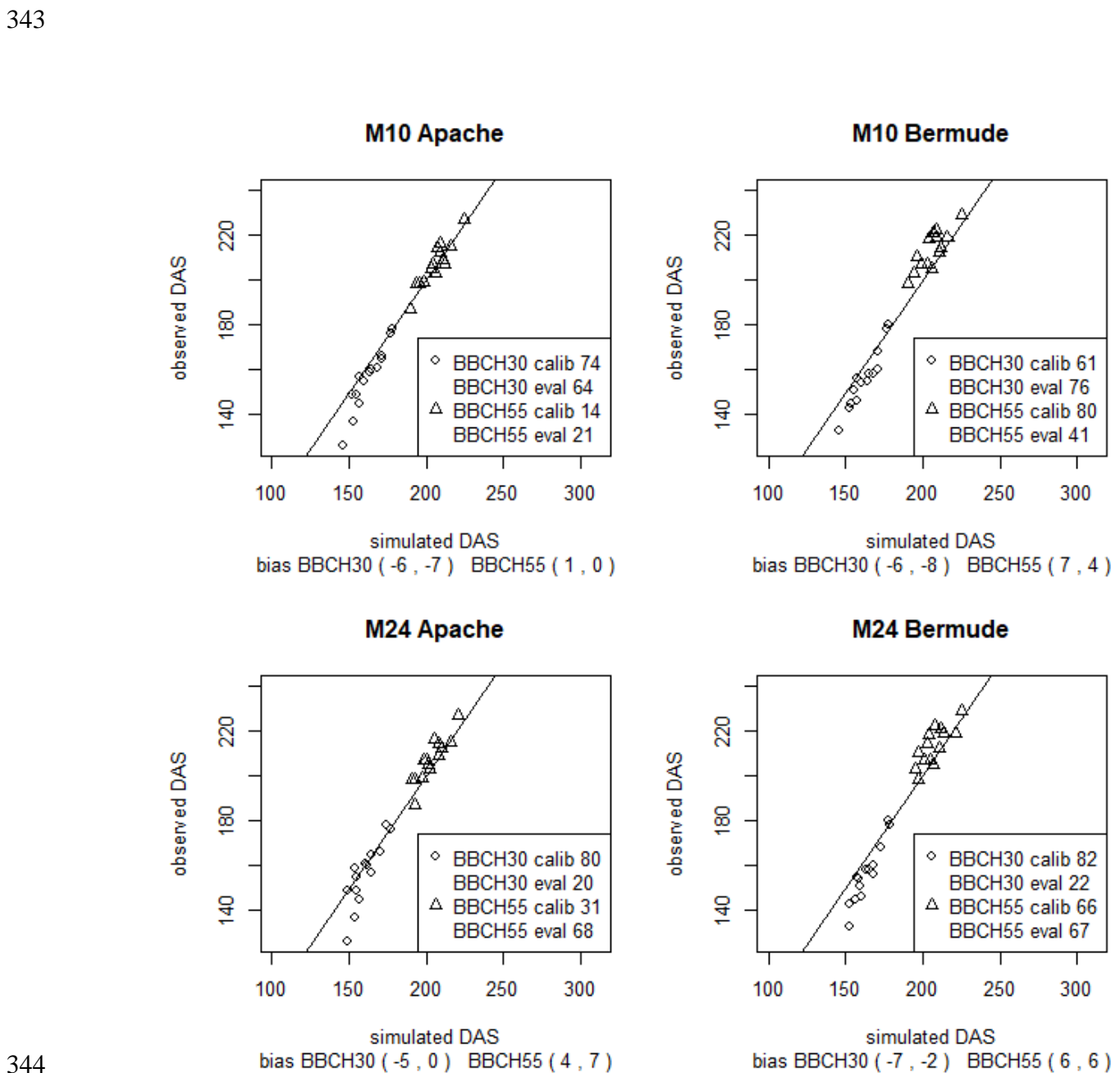
338 **Observed vs. simulated days after sowing (DAS) for calibration data for models**

339 **M10 and M24. The legend shows MSE (days²) for each stage and for calibration and**

340 **evaluation data. (The individual evaluation results are not displayed). In the subtitles,**

341 **bias values (days) for each stage are shown. The first number in parentheses is for the**

342 **calibration data, the second number is for the evaluation data.**



344

19

345

346 **Calibration approach**

347 Each participant was asked to calibrate the model in the “usual” way, using the
 348 calibration data provided. The questionnaire about calibration focused on three aspects of
 349 calibration; the criterion of error to be minimized, the software used and the choice of
 350 parameters to estimate. The choices of the participants are summarized in table 3 and choices
 351 for each model are shown in Supplementary table S9.

352 **Table 3**

353 **Summary of calibration approaches. Numbers are number of models with**
 354 **indicated choice. The specific models associated with each choice are shown in**
 355 **Supplementary tables S3 and S9. More information about the software is presented in**
 356 **Supplementary table S10.**

357

Number of parameters ¹	<table border="1"> <thead> <tr> <th>Minimum</th> <th>1st Quartile</th> <th>Median</th> <th>Mean</th> <th>3rd Quartile</th> <th>Maximum</th> </tr> </thead> <tbody> <tr> <td>1.00</td> <td>2.00</td> <td>3.00</td> <td>3.63</td> <td>4.50</td> <td>9.00</td> </tr> </tbody> </table>	Minimum	1st Quartile	Median	Mean	3rd Quartile	Maximum	1.00	2.00	3.00	3.63	4.50	9.00
Minimum	1st Quartile	Median	Mean	3rd Quartile	Maximum								
1.00	2.00	3.00	3.63	4.50	9.00								
Which parameters	Thermal time to a single development stage 16 Thermal time to two or more development stages 6 Related to vernalization 11 Related to photoperiod 11 Related to effect of temperature (e.g. base temperature) 6 Related to phyllochron 6 Related to tiller appearance 2												

	<p>Related to time to emergence ³</p> <p>Parameters unrelated to calibration data² ⁶</p>
Objective function	<p>Sum of squared errors or of root mean squared errors ² ¹</p> <p>Sum of absolute errors ²</p> <p>Concentrated likelihood ¹</p> <p>No single explicit objective function ³</p>
Software ³	<p>Trial and error ¹⁰</p> <p>DIRECT-L (Gablonsky & Kelley, 2001; Johnson, n.d.) ²</p> <p>Ucode (E. P. Poeter, Hill, Banta, Mehl, & Christensen, 2005; Eileen P. Poeter & Hill, 1999) ³</p> <p>DE Optim (Mullen, Ardia, Gil, Windover, & Cline, 2011) ³</p> <p>PEST (Doherty, Hunt, & Tonkin, 2010) ²</p> <p>SCE (Duan, Gupta, & Sorooshian, 1993; Houska, Kraft, Chamorro-Chavez, & Breuer, 2015) ²</p> <p>GLUE (Beven & Binley, 2014; J. He, Jones, Graham, & Dukes, 2010) ¹</p> <p>DREAM (J. A. Vrugt et al., 2009; Jasper A. Vrugt, 2016) ²</p> <p>Wrote code ⁴ ⁴</p>

358 ¹ Summary of number of estimated parameters for models M1-M27. ² These are
 359 parameters that do not affect simulated days to BBCH30 or BBCH55. ³ Some modeling
 360 groups used more than one software package. ⁴ Modeling groups that wrote their own
 361 software.

362 **Objective function**

363 Most modeling groups defined the sum of squared errors or the sum of root mean
 364 squared errors as the objective function to be minimized, where the sum is over the two

365 stages. (In all cases, the calibration was done separately for the two varieties). Two groups
366 minimized the sum of absolute errors. Calibration for model M21 was based on maximizing
367 the concentrated likelihood (Seber & Wild, 1989) assuming a normal distribution of errors
368 with possibly different error variances for the two development stages. In this case, the
369 objective function involves a product of errors for the two outputs, rather than a sum. Four of
370 the participants (M12, M16, M18) did not define an explicit objective function to be
371 minimized. In these cases, the parameter values were chosen to obtain a “good fit” to the data
372 by visual inspection. Finally, two of the models (M7, M8) divided the calibration into two
373 steps. In these cases three of the parameters were used to fit the BBCH30 data, and then in
374 another step another parameter was used to fit the BBCH55 data.

375 Minimizing the sum of squared errors is a standard statistical approach to model
376 calibration, which has highly desirable properties if certain assumptions about model error are
377 satisfied, including equal variance of model error for all data points and non-correlation of
378 model errors. Only model M21 took into account the possibility that the error variances are
379 different for BBCH30 and BBCH55, and none of the modeling groups took into account
380 possible correlations between errors for BBCH30 and BBCH55 in the same field. Based on
381 the errors for all the data and all the models, it was found that there is a highly significant
382 difference in variance between errors for BBCH30 (variance of error 100.7 days²) and
383 BBCH55 (variance of error 67.3 days²). Also, the correlation between the error for BBCH30
384 and the error for BBCH55 in the same field is 0.53 and highly significant. However, if only
385 results for a single model are considered, then for most models the differences in variance and
386 the correlation are not significant.

387 Two models defined a posterior probability of the parameters equal to the likelihood
388 times the prior probability, as usually assumed in a Bayesian approach. The parameters used

389 for prediction were those that maximized the posterior probability (i.e., the estimated mode of
390 the posterior distribution). In both cases, the likelihood was assumed Gaussian with
391 independent errors, and the prior distribution was assumed uniform between some minimum
392 and maximum value. This approach is equivalent to minimizing the mean squared error, with
393 constraints on the parameter values.

394 **Software**

395 Seven participants simply used trial and error to search for the optimal parameters.
396 The other participants used software specifically adapted to minimizing the objective
397 function, either written specifically for their model or, in most cases, available from other
398 sources (Supplementary table S10).

399 **Choice of parameters to estimate**

400 The choice of parameters to estimate was based on expert judgement in most cases.
401 The participants declared that they chose parameters known to affect phenology in the model,
402 or more specifically parameters expected to have a major effect on time to BBCH30 and
403 BBCH55 and expected to differ between varieties. Five participants did a sensitivity analysis
404 to aid in the choice of parameters to estimate. The number of estimated parameters ranged
405 from 1 to 9. In almost all cases, the number of parameters to estimate was decided a priori. In
406 three cases, the number was the result of testing the fit with different numbers of parameters.
407 In one of those cases the Akaike Information Criteria (AIC, Akaike, 1973) and adjusted R^2
408 were used to test whether additional parameters should be estimated.

409 Almost all modeling groups estimated one or more parameters that represent thermal
410 time between development stages (table 6). Some adjustments were necessary for models that
411 did not explicitly calculate time of BBCH30 or BBCH55. In model M2, for example, a new
412 parameter was added to the model, and estimated, representing the fraction of thermal time

413 from double ridge to heading at which BBCH30 occurs. Thirteen groups estimated a
414 parameter related to the effect of photoperiod. Ten groups estimated a parameter related to
415 vernalization. Six groups modified one or more parameters related to the temperature
416 response (for example model M6 estimated *Tbase*, the temperature below which there is no
417 development). Only three models modified parameters related to the time from sowing to
418 emergence, and only one model modified a parameter related to the effect of water stress. Six
419 models included among the parameters to estimate, parameters that have no effect on the
420 variables furnished as calibration data. Such parameters included thermal times for
421 development stages after BBCH55, potential kernel growth rate, kernel number per stem
422 weight and the temperature below which there is 50% death due to cold (Supplementary table
423 S9).

424

425 Discussion

426 Prediction error

427 The challenge in this study was to predict the time from sowing to beginning of stem
428 elongation and to heading in winter wheat field trials performed across France. This is a
429 problem of practical importance, since these two development stages are important for wheat
430 management (e.g. fertilization). The evaluation concerned years and sites not included in the
431 calibration data, making this one of the most rigorous evaluations to date of how well crop
432 models simulate phenology.

433 Twenty-seven modeling groups participated in the exercise. Most models predicted
434 times to stem elongation and heading quite well (median MAE of 6 days). Half the models
435 had MSE values of prediction that were 36% or less than MSE of a naive predictor. It must be

436 kept in mind that this study is a rather favorable situation for prediction, with a substantial
437 amount of calibration data and predictions for environments similar to those of the calibration
438 data.

439 **Role of calibration**

440 What is the role of calibration in determining prediction accuracy? We cannot answer
441 this exactly, because differences between models result not only from differences in
442 calibration approach, but also from differences in structure and from differences in the values
443 of parameters not estimated by calibration. However, several aspects of the results indicate
444 that calibration is important.

445 Consider first the comparison between models with the same structure. There are
446 fairly large differences in MSE between models with the same structure. This could partially
447 be due to different values for the parameters not estimated by calibration. However, since
448 there are major differences in calibration approach, and in general the parameters estimated
449 by calibration are among the most important controlling phenology, it seems reasonable to
450 conclude that the differences between models with the same structure are largely due to
451 differences in calibration.

452 Conversely, our results indicate that calibration can result in models with very
453 different structures achieving similar values of MSE. One essential aspect of model structure
454 is the choice of input variables. In fact, MSE can be expressed as a sum of two terms, the first
455 of which depends only on the choice of the model input variables, while the second measures
456 the distance between the model used and the optimal model for those inputs (Wallach,
457 Makowski, Jones, & Brun, 2019). Calibration has a major effect of the second term, and in
458 fact the objective of calibration is to minimize that term. The most important inputs that
459 determine spring wheat phenology are daily temperature and photoperiod (Aslam et al., 2017)

460 and for winter wheat it is also necessary to include the process of vernalization, i.e. the effect
461 of low winter temperatures on development (Li et al., 2013). Five of the best eight predicting
462 models here, with $MSE_{eval}^{all} < 40 \text{ days}^2$, do use all three of those variables (daily temperature,
463 photoperiod, vernalizing temperatures) as inputs. Two of those best eight models however do
464 not use vernalizing temperatures, and one of those best eight does not use photoperiod. Thus
465 there are similarly low values of MSE for prediction even for models so fundamentally
466 different in structure that they use different input variables. It seems likely that this is largely
467 due to the fact that the different models are calibrated using the same data.

468 Another indication that calibration compensates for differences in structure is the
469 result that there is less variability between models for predicting days from sowing to
470 BBCH30, which is provided as calibration data, than would be expected if the uncertainties in
471 days from sowing to emergence and days from emergence to BBCH30 simply added up.
472 Compensation is usually discussed in the context of single models. For example, equifinality,
473 which is a well-known phenomenon of complex models, means that different combinations of
474 parameter values, and thus different quantitative descriptions of processes, can lead to the
475 same results for outputs because there is compensation between the processes (Beven, 2006;
476 D. He et al., 2017). However, this phenomenon has not been described in the context of multi-
477 model studies. Here, we have an example of compensation for differences between models in
478 the way they partition days from sowing to BBCH30 into days from sowing to emergence
479 plus days from emergence to BBCH30. Models with longer simulated times from sowing to
480 emergence tend to have a shorter simulated time from emergence to development stage
481 BBCH30 and vice versa. In fact, each extra day from sowing to emergence is associated on
482 average with almost exactly one less day from emergence to BBCH30. The result is that
483 models with quite different simulated days from sowing to emergence can have nearly
484 identical times from sowing to BBCH30. This can be expressed in terms of model

485 uncertainty, as quantified by between-model variance. The variance of days from sowing to
486 BBCH30 is less than the sum of variances of days from sowing to emergence and days from
487 emergence to BBCH30. That is, calibration reduces, but does not eliminate, model uncertainty
488 for the variable provided for calibration.

489 We do not have observed time to emergence, but in any case, the models with
490 different simulated days to emergence can't all be right. This is an example of how models
491 can get the right answer (correct days to BBCH30, thanks to calibration) for the wrong
492 reasons (wrong days to emergence and compensating wrong days from emergence to
493 BBCH30), illustrating the problem pointed out for example by (Challinor, Martre, Asseng,
494 Thornton, & Ewert, 2014). The same compensation of errors between sowing to emergence
495 and emergence to BBCH30 will not be appropriate for all environments. This is one of the
496 main reasons that extrapolation to populations different than the calibration population is
497 dangerous.

498 Much previous work on improving the predictive capability of crop models has
499 focused on the model equations, for instance the way temperature is taken into account in
500 various processes (Maiorano et al., 2016; Wang et al., 2017). Here we show that models with
501 the same structure can have very different levels of prediction error, if the calibration methods
502 differ, while models with quite different structures can have very similar prediction accuracy,
503 thanks to calibration using the same data. This means that model comparison studies may
504 often be comparing calibration approaches as much or more as they are comparing model
505 equations. This is in line with the conclusions of Confalonieri et al. (2016), who argued that
506 one should not speak of evaluation of a model but rather of the combination of a model and a
507 model user, where a major role of the user is in implementing calibration.

508 **Calibration approach**

509 This study was designed to identify how different groups do calibration, given the
510 same data and prediction objectives. We focused here on three specific aspects of the
511 calibration approach; the choice of objective function, the software used and the choice of
512 parameters to estimate. The results show the diversity of approaches. Since different models
513 differ in multiple ways, the study does not allow us to define best practices for each aspect of
514 calibration. However, it is possible to point out practices which should probably be avoided.

515 **Objective function**

516 Most participants defined an objective function based on what one would use in a
517 statistical approach to non-linear regression, namely a sum of squared errors to be minimized
518 or a likelihood to be maximized. However three models (see Supplementary table S9) did
519 not have an explicit quantitative objective function. Those models all had relatively large
520 values of overall MSE for the evaluation data (MSE_{eval}^{all}), having 15th, 16th, and 18th largest
521 MSE_{eval}^{all} values out of the 25 models that predicted both BBCH30 and BBCH55. It seems
522 reasonable to suppose that the lack of a quantitative objective function can be a drawback
523 since then one does not have a clear criterion for deciding on the best parameter values.

524 Among the models that chose to minimize a sum of squared errors or to maximize a
525 likelihood, all but one implicitly or explicitly assumed that all model errors had equal variance
526 and were independent. This will in general not be the case when there are multiple
527 measurements in the same field, as is the case here (measurement of days to BBCH30 and
528 days to BBCH55 in each field). Ignoring unequal variances and correlated errors in non-linear
529 regression leads to inefficient estimators (Seber & Wild, 1989). One should at least test
530 whether heteroscedasticity and non-independence are important.

531 **Software**

532 Several different software solutions were used for calibration by the different models.
533 There does not seem to be any clear connection between the software used and predictive
534 quality. Various different software solutions were used by the best predicting models, but
535 largely the same software solutions were also found among the models with the largest
536 prediction errors.

537 A problem that may arise concerns the test for convergence to the parameter values
538 that minimize the chosen objective function. Having such a test allows the user to have
539 confidence that the best parameter values have been found. With trial and error, there is no
540 such test, which is a major drawback of this approach. Algorithms to estimate a Bayesian
541 posterior distribution normally test convergence to the posterior distribution, which may not
542 be relevant if one is using just the mode of the distribution. It would be good practice to adopt
543 a software option that includes an appropriate test of convergence.

544 **Choice of parameters to estimate**

545 There was a large diversity of choices of parameters to estimate by calibration, and
546 this had in certain cases an important effect on prediction error. One rather unexpected
547 observation was that several participants included, among the parameters to estimate,
548 parameters that have no effect on the variables furnished as calibration data among the
549 parameters to estimate. The data cannot in those cases give any information about the
550 parameter value. At best, including such parameters among the parameters to estimate is
551 useless, and those parameters will simply have final values exactly equal to their initial
552 values. However, there may also be serious disadvantages to including such parameters. It
553 gives the erroneous impression that one is estimating parameters that cannot in fact be
554 estimated, it increases computation time and it can cause problems for the parameter

555 estimation algorithm. The very poor fit of model M5 to the calibration data seems to be
556 directly related to the fact that for this model, several parameters unrelated to the calibration
557 data were chosen to be fitted. The software used here was PEST (Doherty et al., 2010), with
558 the singular value decomposition option, which allows one to deal with non-estimable
559 parameters, but at the cost of introducing bias in the estimated parameter values. This bias
560 may be at the origin of the poor performance. Obviously, one should not include non-
561 estimable parameters among the parameters to estimate.

562 The choice of parameters to estimate may be the principal cause of bias in fitting the
563 calibration data for some models. If a model includes an additive constant term, and squared
564 error is minimized, bias will be 0 for the calibration data. Even for more complex models,
565 calibration can bring bias close to 0, as illustrated here by the fact that many of the models
566 had very small biases for the calibration data. Eliminating bias is important, since squared
567 bias is one component of MSE, and therefore the bias necessarily adds on to MSE (Kobayashi
568 & Salam, 2000). If one does not have a parameter with a nearly additive effect for each of the
569 development stages BBCH30 and BBCH55, the elimination of bias for both outputs is not
570 assured. Model M24 estimated only a single parameter. In such a case, at best one can
571 estimate a parameter value that gives the best compromise between errors in BBCH30 and
572 BBCH55. This may lead to a negative bias for one of those outputs and more or less
573 corresponding positive bias for the other. This is exactly the behavior illustrated in figure 2.
574 Model M10 also had fairly large biases. Here three parameters were estimated, but one is
575 unrelated to the observed data and a second concerns time to emergence, which was only
576 allowed to vary in a limited range. Apparently in this case also there was not enough
577 flexibility to eliminate bias for both development stages. Models with large bias for the
578 calibration data tended to have large MSE values for the evaluation data (Supplementary
579 figure S2). This suggests that the parameters to estimate should include one parameter that is

580 nearly additive (i.e. that adds an amount that is nearly the same for all environments) for each
581 observed output, and that is not too limited in the allowed range of values.

582 Conclusions

583 Overall, we have shown in a rigorous evaluation of prediction for new environments
584 that most of the 27 crop models tested, given calibration data, provide good predictions of
585 phenology in winter wheat and do much better than predicting with the average of the
586 calibration data. Calibration has a major effect on predictive quality. Calibration can
587 compensate to some extent even for different choices of input variables. It reduces variability
588 between models for outputs used for calibration, but may lead to models getting the right
589 answer for the wrong reason. Poor practices of calibration can seriously degrade predictive
590 capability. Arguably the most difficult aspect of calibration, and yet the least studied, is the
591 choice of parameters to estimate. Unlike the choice of objective function and of software,
592 there is little guidance here from other fields. Furthermore, the problem is specific to each
593 model, since each model has a different set of parameters. Given the large diversity of
594 calibration approaches and the importance of calibration, there is a clear need for guidelines
595 and tools to aid model users with respect to calibration. Model applications, including model
596 studies of climate change impact, should focus more on the data used for calibration and on
597 the calibration methods employed.

598

599 Acknowledgements

600 This work was in part supported by the Collaborative Research Center 1253 CAMPOS
601 (Project 7: Stochastic Modelling Framework), funded by the German Research Foundation
602 (DFG, Grant Agreement SFB 1253/1 2017), the Academy of Finland through projects AI-
603 CropPro (316172) and DivCSA (316215) and Natural Resources Institute Finland (Luke)
604 through a strategic project BoostIA, the BonaRes project "Soil3" (BOMA 03037514) of the
605 Federal Ministry of Education and Research (BMBF), Germany, the project BiomassWeb of
606 the GlobeE programme (Grant number: FKZ031A258B) funded by the Federal Ministry of
607 Education and Research (BMBF, Germany), the INRA ACCAF meta-programme, the
608 German Federal Ministry of Education and Research (BMBF) in the framework of the
609 funding measure "Soil as a Sustainable Resource for the Bioeconomy – BonaRes", project
610 "BonaRes (Module B): BonaRes Centre for Soil Research, subproject B" (grant 031B0511B),
611 the National Key Research and Development Program of China (2017YFD0300205), the
612 National Science Foundation for Distinguished Young Scholars (31725020), the Priority
613 Academic Program Development of Jiangsu Higher Education Institutions (PAPD), the 111
614 project (B16026), and China Scholarship Council, the Agriculture and Agri-Food Canada's
615 Project 1387 under the Canadian Agricultural Partnership, the DFG Research Unit FOR 1695
616 'Agricultural Landscapes under Global Climate Change – Processes and Feedbacks on a
617 Regional Scale, the U.S. Department of Agriculture National Institute of Food and
618 Agriculture (award no. 2015-68007-23133) and USDA/NIFA HATCH grant N. MCL02368,
619 the National Key Research and Development Program of China (2016YFD0300105), The
620 Broadacre Agriculture Initiative, a research partnership between University of Southern
621 Queensland and the Queensland Department of Agriculture and Fisheries, the Academy of
622 Finland through project AI-CropPro (315896), the JPI FACCE MACSUR2 project, funded by

623 the Italian Ministry for Agricultural, Food, and Forestry Policies (D.M. 24064/7303/15 of
624 26/Nov/2015). The order in which the donors are listed is arbitrary.

625

626

627

628 References

629

630 Akaike, H. (1973). Information Theory and an Extension of the Maximum Likelihood
631 Principle. In B. N. Petrov & F. Csaki (Eds.), *In B. N. Petrov, & F. Csaki (Eds.),*
632 *Proceedings of the 2nd International Symposium on Information Theory* (pp. 267–281).
633 Budapest: Akademiai Kiado.

634 Andarzian, B., Hoogenboom, G., Bannayan, M., Shirali, M., & Andarzian, B. (2015).
635 Determining optimum sowing date of wheat using CSM-CERES-Wheat model. *Journal*
636 *of the Saudi Society of Agricultural Sciences*, 14(2), 189–199.
637 <https://doi.org/10.1016/J.JSSAS.2014.04.004>

638 Aslam, M. A., Ahmed, M., Stöckle, C. O., Higgins, S. S., Hassan, F. ul, & Hayat, R. (2017).
639 Can Growing Degree Days and Photoperiod Predict Spring Wheat Phenology? *Frontiers*
640 *in Environmental Science*, 5, 57. <https://doi.org/10.3389/fenvs.2017.00057>

641 Asseng, S., Ewert, F., Martre, P., Rötter, R. P., Lobell, D. B., Cammarano, D., ... Zhu, Y.
642 (2015). Rising temperatures reduce global wheat production. *Nature Climate Change*,
643 5(2), 143–147. <https://doi.org/10.1038/nclimate2470>

644 Beven, K. (2006). A manifesto for the equifinality thesis. *Journal of Hydrology*, 320(1–2),

645 18–36.

646 Beven, K., & Binley, A. (2014). GLUE: 20 years on. *Hydrological Processes*, 28(24), 5897–
647 5918. <https://doi.org/10.1002/hyp.10082>

648 Challinor, A., Martre, P., Asseng, S., Thornton, P., & Ewert, F. (2014). Making the most of
649 climate impacts ensembles. *Nature Climate Change*, 4(2), 77–80.
650 <https://doi.org/10.1038/nclimate2117>

651 Confalonieri, R., Orlando, F., Paleari, L., Stella, T., Gilardelli, C., Movedi, E., ... Acutis, M.
652 (2016). Uncertainty in crop model predictions: What is the role of users? *Environmental*
653 *Modelling & Software*, 81, 165–173. <https://doi.org/10.1016/j.envsoft.2016.04.009>

654 Doherty, J. E., Hunt, R. J., & Tonkin, M. J. (2010). *Approaches to highly parameterized*
655 *inversion: A guide to using PEST for model-parameter and predictive-uncertainty*
656 *analysis: U.S. Geological Survey Scientific Investigations Report 2010–5211*. Retrieved
657 from <http://pubs.usgs.gov/sir/2010/5211>

658 Duan, Q. Y., Gupta, V. K., & Sorooshian, S. (1993). Shuffled complex evolution approach for
659 effective and efficient global minimization. *Journal of Optimization Theory and*
660 *Applications*, 76(3), 501–521. <https://doi.org/10.1007/BF00939380>

661 Gablonsky, J. M., & Kelley, C. T. (2001). A Locally-Biased form of the DIRECT Algorithm.
662 *Journal of Global Optimization*, 21(1), 27–37. <https://doi.org/10.1023/A:1017930332101>

663 He, D., Wang, E., Wang, J., Lilley, J., Luo, Z., Pan, X., ... Yang, N. (2017). Uncertainty in
664 canola phenology modelling induced by cultivar parameterization and its impact on
665 simulated yield. *Agricultural and Forest Meteorology*, 232, 163–175.
666 <https://doi.org/10.1016/j.agrformet.2016.08.013>

667 He, J., Jones, J. W., Graham, W. D., & Dukes, M. D. (2010). Influence of likelihood function

- 668 choice for estimating crop model parameters using the generalized likelihood uncertainty
669 estimation method. *Agricultural Systems*, 103(5), 256–264.
670 <https://doi.org/10.1016/j.agsy.2010.01.006>
- 671 Houska, T., Kraft, P., Chamorro-Chavez, A., & Breuer, L. (2015). SPOTting Model
672 Parameters Using a Ready-Made Python Package. *PLOS ONE*, 10(12), e0145180.
673 <https://doi.org/10.1371/journal.pone.0145180>
- 674 Hunt, J. R., Lilley, J. M., Trevaskis, B., Flohr, B. M., Peake, A., Fletcher, A., ... Kirkegaard,
675 J. A. (2019). Early sowing systems can boost Australian wheat yields despite recent
676 climate change. *Nature Climate Change*, 9(3), 244–247. [https://doi.org/10.1038/s41558-](https://doi.org/10.1038/s41558-019-0417-9)
677 [019-0417-9](https://doi.org/10.1038/s41558-019-0417-9)
- 678 Hussain, J., Khaliq, T., Ahmad, A., & Akhtar, J. (2018). Performance of four crop model for
679 simulations of wheat phenology, leaf growth, biomass and yield across planting dates.
680 *PLOS ONE*, 13(6), e0197546. <https://doi.org/10.1371/journal.pone.0197546>
- 681 Johnson, S. G. (n.d.). The NLOpt nonlinear-optimization package.
- 682 Kobayashi, K., & Salam, M. U. (2000). Comparing simulated and measured values using
683 mean squared deviation and its components. *Agronomy Journal*, 92, 345–352.
- 684 Li, G., Yu, M., Fang, T., Cao, S., Carver, B. F., & Yan, L. (2013). Vernalization requirement
685 duration in winter wheat is controlled by TaVRN-A1 at the protein level. *The Plant*
686 *Journal: For Cell and Molecular Biology*, 76(5), 742–753.
687 <https://doi.org/10.1111/tpj.12326>
- 688 Maiorano, A., Martre, P., Asseng, S., Ewert, F., Müller, C., Rötter, R. P., ... Zhu, Y. (2017).
689 Crop model improvement reduces the uncertainty of the response to temperature of
690 multi-model ensembles. *Field Crops Research*, 202.

- 691 <https://doi.org/10.1016/j.fcr.2016.05.001>
- 692 Maiorano, Andrea, Martre, P., Asseng, S., Ewert, F., Müller, C., Rötter, R. P., ... Zhu, Y.
693 (2016). Crop model improvement reduces the uncertainty of the response to temperature
694 of multi-model ensembles. *Field Crops Research*.
695 <https://doi.org/10.1016/j.fcr.2016.05.001>
- 696 Mullen, K., Ardia, D., Gil, D., Windover, D., & Cline, J. (2011). “DEoptim”: An R Package
697 for Global Optimization by Differential Evolution. *Journal of Statistical Software*, 40, 1–
698 26.
- 699 Piao, S., Liu, Q., Chen, A., Janssens, I. A., Fu, Y., Dai, J., ... Zhu, X. (2019). Plant phenology
700 and global climate change: current progresses and challenges. *Global Change Biology*,
701 gcb.14619. <https://doi.org/10.1111/gcb.14619>
- 702 Poeter, E. P., Hill, M. C., Banta, E. R., Mehl, S., & Christensen, S. (2005). *UCODE_2005 and*
703 *Six Other Computer Codes for Universal Sensitivity Analysis, Calibration, and*
704 *Uncertainty Evaluation: U.S. Geological Survey Techniques and Methods 6-A11*.
- 705 Poeter, Eileen P., & Hill, M. C. (1999). UCODE, a computer code for universal inverse
706 modeling. *Computers & Geosciences*, 25(4), 457–462. <https://doi.org/10.1016/S0098->
707 [3004\(98\)00149-6](https://doi.org/10.1016/S0098-3004(98)00149-6)
- 708 Rauff, K. O., & Bello, R. (2015). A Review of Crop Growth Simulation Models as Tools for
709 Agricultural Meteorology. *Agricultural Sciences*, 06(09), 1098–1105.
710 <https://doi.org/10.4236/as.2015.69105>
- 711 Rezaei, E. E., Siebert, S., & Ewert, F. (2015). Intensity of heat stress in winter wheat—
712 phenology compensates for the adverse effect of global warming. *Environmental*
713 *Research Letters*, 10(2), 024012. <https://doi.org/10.1088/1748-9326/10/2/024012>

- 714 Rezaei, E. E., Siebert, S., Hüging, H., & Ewert, F. (2018). Climate change effect on wheat
715 phenology depends on cultivar change. *Scientific Reports*, 8(1), 4891.
716 <https://doi.org/10.1038/s41598-018-23101-2>
- 717 Seber, G. A. F., & Wild, C. J. (1989). *Nonlinear regression*. New York: Wiley .
- 718 Seidel, S. J., Palosuo, T., Thorburn, P., & Wallach, D. (2018). Towards improved calibration
719 of crop models – Where are we now and where should we go? *European Journal of*
720 *Agronomy*, 94, 25–35. <https://doi.org/10.1016/J.EJA.2018.01.006>
- 721 Svystun, T., Bhalerao, R. P., & Jönsson, A. M. (2019). Modelling Populus autumn phenology:
722 The importance of temperature and photoperiod. *Agricultural and Forest Meteorology*,
723 271, 346–354. <https://doi.org/10.1016/J.AGRFORMET.2019.03.003>
- 724 van Ittersum, M. ., Leffelaar, P. ., van Keulen, H., Kropff, M. ., Bastiaans, L., & Goudriaan, J.
725 (2003). On approaches and applications of the Wageningen crop models. *European*
726 *Journal of Agronomy*, 18(3–4), 201–234. [https://doi.org/10.1016/S1161-0301\(02\)00106-](https://doi.org/10.1016/S1161-0301(02)00106-5)
727 5
- 728 Vrugt, J. A., ter Braak, C. J. F., Diks, C. G. H., Robinson, B. A., Hyman, J. M., & Higdon, D.
729 (2009). Accelerating Markov Chain Monte Carlo Simulation by Differential Evolution
730 with Self-Adaptive Randomized Subspace Sampling. *International Journal of Nonlinear*
731 *Sciences and Numerical Simulation*, 10(3).
732 <https://doi.org/10.1515/IJNSNS.2009.10.3.273>
- 733 Vrugt, Jasper A. (2016). Markov chain Monte Carlo simulation using the DREAM software
734 package: Theory, concepts, and MATLAB implementation. *Environmental Modelling &*
735 *Software*, 75, 273–316. <https://doi.org/10.1016/J.ENVSOF.2015.08.013>
- 736 Wallach, D., Makowski, D., Jones, J. W., & Brun, F. (2019). *Working with Dynamic Crop*

- 737 *Models: Methods, Tools and examples for Agriculture and Environment*. London, U.K.:
738 Academic Press.
- 739 Wang, E., Martre, P., Zhao, Z., Ewert, F., Maiorano, A., Rötter, R. P., ... Asseng, S. (2017).
740 The uncertainty of crop yield projections is reduced by improved temperature response
741 functions. *Nature Plants*, 3. <https://doi.org/10.1038/nplants.2017.102>
- 742 Yuan, S., Peng, S., & Li, T. (2017). Evaluation and application of the ORYZA rice model
743 under different crop managements with high-yielding rice cultivars in central China.
744 *Field Crops Research*, 212, 115–125. <https://doi.org/10.1016/J.FCR.2017.07.010>
- 745 Zadoks, J. C., Chzang, T. T., & Konzak, C. F. (1974). A decimal code for the growth stages
746 of cereals. *Weed Research*, 14(6), 415–421. [https://doi.org/10.1111/j.1365-](https://doi.org/10.1111/j.1365-3180.1974.tb01084.x)
747 [3180.1974.tb01084.x](https://doi.org/10.1111/j.1365-3180.1974.tb01084.x)
- 748
- 749