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

2

3

4 # **Running head: Crop model phenology prediction**

5 Wallach[1], Daniel; Palosuo[2], Taru; Thorburn[3], Peter; Seidel[4], Sabine J.; Gourdain[5], Emmanuelle; Asseng[6],
6 Senthold; Basso[7], Bruno; Buis[8], Samuel; Crout[9], Neil, M.J.; Dibari[10], Camilla; Dumont[11], Benjamin;
7 Ferrise[10], Roberto; Gaiser[4], Thomas; Garcia[7], Cécile; Gayler[12], Sebastian; Ghahramani[13], Afshin;
8 Hochman[3], Zvi; Hoek[14], Steven; Horan[3], Heidi; Hoogenboom[6,15], Gerrit; Huang[16], Mingxia; Jabloun[9],
9 Mohamed; Jing[17], Qi; Justes[18], Eric; Kersebaum[19], Kurt Christian; Klosterhalfen[20], Anne; Launay[21],
10 Marie; Luo[22], Qunying; Maestrini[7], Bernardo; Mielenz[23], Henrike; Moriondo[24], Marco; Nariman Zadeh[25],
11 Hasti; Olesen[26], Jørgen Eivind; Poyda[27], Arne; Priesack[28], Eckart; Pullens[26], Johannes Wilhelmus Maria;
12 Qian[17], Budong; Schütze[29], Niels; Shelia[6,15], Vakhtang; Souissi[30,31], Amir; Specka[19], Xenia; Srivastava[4],
13 Amit Kumar; Stella[19], Tommaso; Streck[12], Thilo; Trombi[10], Giacomo; Wallor[19], Evelyn; Wang[16], Jing;
14 Weber[12], Tobias, K.D.; Weihermüller[20], Lutz; de Wit[14], Allard; Wöhling[29,32], Thomas; Xiao[33,6], Liujun;
15 Zhao[6], Chuang; Zhu[33], Yan

[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] Institute for Crop and Soil Science, Federal Research Centre for cultivated Plants, Julius Kühn-Institut (JKI), Braunschweig, Germany

[24] CNR-IBIMET, Firenze, Italy

[25] Aalto University School of Science, Espoo, Finland

[26] Department of Agroecology, Aarhus University, Tjele, Denmark

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

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

[29] Institute of Hydrology and Meteorology, Chair of Hydrology, Technische Universität Dresden, Dresden, Germany

[30] National Institute of Agronomic Research of Tunisia (INRAT), Agronomy Laboratory, University of Carthage, Tunis, Tunisia

[31] National Agronomy Institute of Tunisia (INAT), University of Carthage, Tunis, Tunisia

[32] Lincoln Agritech Ltd., Hamilton, New Zealand

[33] 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

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. Since climate
19    change is projected to alter crop phenology worldwide, there is a large effort to predict
20    phenology as a function of environment. Many studies have focused on comparing model
21    equations for describing how phenology responds to weather but the effect of crop model
22    calibration, also expected to be important, has received much less attention. The objective
23    here was to obtain a rigorous evaluation of prediction capability of wheat crop phenology
24    models, and to analyze the role of calibration. The 27 participants in this multi-model study
25    were provided experimental data for calibration and asked to submit predictions for sites and
26    years not represented in those data. Participants were instructed to use and document their
27    "usual" calibration approach. Overall, the models provided quite good predictions of
28    phenology (median of mean absolute error of 6.1 days) and did much better than simply using
29    the average of observed values as predictor.  Calibration was found to compensate to some
30    extent for differences between models, specifically for differences in simulated time to
31    emergence and differences in the choice of input variables. Conversely, different calibration
32    approaches led to major differences in prediction error between models with the same
33    structure. Given the large diversity of calibration approaches and the importance of
34    calibration, there is a clear need for guidelines and tools to aid with calibration. Arguably the
35    most important and difficult choice for calibration is the choice of parameters to estimate.
36    Several recommendations for calibration practices are proposed. Model applications,
37    including model studies of climate change impact, should focus more on the data used for
38    calibration and on the calibration methods employed.

39    ## Introduction

40    Global change is placing unprecedented pressure on food security (Campbell et al.,
41    2018; Godfray et al., 2010; Wheeler & von Braun, 2013). The search for ways to increase the
42    volume and efficiency of food production needs to account for the effects of climate change
43    on agricultural production systems (Porter, Howden, & Smith, 2017). Modeling is an
44    important tool in projecting and understanding the trajectory of food production under future
45    climates   (Asseng et al., 2019; Bindi, Palosuo, Trnka, & Semenov, 2015; Carberry et al.,
46    2013; Hochman, Gobbett, & Horan, 2017; Parry, Rosenzweig, Iglesias, Livermore, & Fischer,
47    2004). One of the likely effects of climate change will be an increase in air temperature
48    (IPCC, 2014). Temperature directly affects plant growth through a number of pathways, one
49    of which is phenology (Went, 2003). Phenology describes the cycles of biological events in
50    plants. These events include seedling emergence, leaf appearance, and flowering. Matching
51    the phenology of crop varieties to the climate in which they grow is a critical crop production
52    strategy (Hunt et al., 2019; Rezaei, Siebert, & Ewert, 2015; Rezaei, Siebert, Hüging, & Ewert,
53    2018). Thus, understanding and improving our ability to simulate phenology with crop
54    models is an important activity in preparing for and adapting to global change. Process-based
55    models similar to those for crops can be used for natural vegetation, so crop models can serve
56    as examples for studies of phenology in ecosystems (Piao et al., 2019).

57    Crop model evaluation is an essential aspect of modeling, assessing whether model
58    performance is acceptable for the intended use of the model. For studies of phenology two
59    major questions are a) how accurate are current models for the prediction of crop
60    development stages? and  b) what determines model accuracy and what does that imply about
61    how accuracy can be improved?  We use here prediction in the sense of determining outputs

62    (dates of development stages) from known inputs (weather, soil, management). The problem
63    of predicting future events, with unknown weather, is not considered.

64          There have been numerous evaluation studies of crop model simulations, including but
65    not restricted to phenology, both of individual models and of multi-model ensembles. The
66    typical procedure is to first calibrate the model using a part of the available field data and then
67    to evaluate it using the remaining data. Most crop model evaluation studies focusing on crop
68    phenology have had relatively little data for calibration or evaluation. (Andarzian,
69    Hoogenboom, Bannayan, Shirali, & Andarzian, 2015) for example, used data from one
70    location covering five growing seasons and two or three sowing dates per year. Out of these
71    data, one year was used for calibration and the other two years of data to evaluate the model.
72    (Yuan, Peng, & Li, 2017) used one year of data for calibration and the second year of data
73    from the same location for evaluation of the rice crop model ORYZA. Hussain, Khaliq,
74    Ahmad, & Akhtar (2018) tested four models using data from two locations with two years of
75    data, 11 crop planting dates, and three varieties. Paucity of data means that model parameters
76    are estimated with relatively large uncertainty and model evaluation is quite uncertain.

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

84          A third feature often found in crop model evaluation is that the range of situations
85    from which the calibration data are drawn (the "training population") is often different than
86    the range of conditions from which the evaluation data are drawn (the "evaluation
87    population"). For example, Hussain et al. (2018) used data from an experiment that included a
88    range of crop stresses to calibrate their model. They used data from the least stressed
89    treatment in the calibration process and evaluated the resultant model on the remaining
90    planting dates at the same location. The evaluation data thus represented a different range of
91    conditions than the calibration data. In a multi-model ensemble study of the effect of high
92    temperatures on wheat growth (S. Asseng et al., 2015) detailed crop measurements were
93    provided for one planting date and the models were evaluated using other planting dates,
94    some with additional artificial heating. Again, the evaluation data represented a much larger
95    range of temperatures than represented in the calibration data. While the capacity of crop
96    models to extrapolate to conditions quite different than those of the calibration data is
97    obviously of interest, it is a rather different type of evaluation than the case where the training
98    and evaluation populations are similar.

99          Thus, evaluation of crop phenology models to date has mainly concerned situations
100    that would tend to make prediction difficult, because of small amounts of data for calibration
101    and differences between the training and target populations. Furthermore, the quality of the
102    evaluation is often questionable, because of the relatively small amounts of data and the non-
103    independence of the errors for the calibration and evaluation situations. There is thus a need
104    for more rigorous assessments of simulation capability of crop phenology models in the well-
105    defined situation where the calibration and evaluation data can be assumed to come from the
106    same underlying population.  The first objective of this paper is, therefore, to evaluate how
107    well crop models predict wheat phenology in such a case. To ensure the rigor of the

2

108  evaluation, we create a situation where the model errors for the calibration and simulation
109  data can be assumed independent.

110      The emphasis in model evaluation studies is often on the role of model structure, i.e.
111  model equations (Maiorano et al., 2017; Svystun, Bhalerao, & Jönsson, 2019; Wang et al.,
112  2017), and not on model calibration. Clearly however the simulated values depend on the
113  parameter values estimated by calibration and therefore on the calibration approach.
114  (Confalonieri et al., 2016) found that the model user, responsible for calibration, had a very
115  large effect on predictive quality. In a wide-ranging survey, (Seidel, Palosuo, Thorburn, &
116  Wallach, 2018) found that there  is a wide diversity of calibration strategies used for crop
117  models.  The second objective of this study was therefore to obtain detailed information about
118  the calibration strategies in use for phenology models and to better understand the effect of
119  calibration methodology in determining predictive capability for phenology. This is of
120  practical interest not only for stand-alone phenology models, but also for crop models more
121  generally, since crop models are often calibrated first just for phenology, and then separately
122  for biomass increase and partitioning and soil processes.

## Materials and Methods

**Experimental data**

125      The data were provided by ARVALIS – Institut du vegetal, 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 check varieties, *Apache*
128  which is a common variety grown throughout France and Central Europe and *Bermude*,
129  mainly grown in Northern and Central France. The trials have three repetitions and follow
130  standard agricultural practices, with N fertilization calculated to be non-limiting. Thus, both
131  the calibration and evaluation data are drawn from sites in France where winter wheat is
132  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 three categories (table 1). One part, the calibration data (14
140  environments i.e. site-year combinations), was provided to participants for calibration. A
141  second part, the evaluation data (eight environments), was not given to participants. The
142  division of the data was such that the calibration and evaluation data had no sites or years in
143  common. The only way to achieve this was have a third category, "other" (from 13
144  environments), with data that were not revealed to participants but were used neither for
145  calibration nor evaluation. These environments had either the site or the year in common with
146  the calibration data. Errors for these environments cannot be assumed to be independent of
147  errors for the calibration data, and so they were not used for evaluation. All conclusions about
148  predictive capability are based on the evaluation data. The individual observed values for the
149  evaluation and other hidden data are not presented here because they will be used again in a
150  subsequent study where all groups will be asked to use the same calibration approach.

151

152

**Table 1.**

**Environments (site-year combinations) that provided the data. C= calibration data. E = evaluation data. O = other hidden data.**

| Site (longitude,latitude) | Harvest year | | | | | | |
|---|---|---|---|---|---|---|---|
| | 2010 | 2011 | 2012 | 2013 | 2014 | 2015 | 2016 |
| FORESTE (3.20,49.82) | | | E | E | OO[*] | O | |
| MERY 4.02,48.33) | C | C | | O | C | C | |
| ROUVRES 5.09,47.28) | | | E | E | O | O | |
| CESSEVILLE[1] (0.90,49.15) | | C | | | | | |
| IVILLE[1] (0.90,49.15) | | | E | | | | |
| VILLETTES[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 |

[*]There were two sowing dates at FORESTE in 2014. [1] These are separate sites that are geographically close to one another and share a single weather station.

The background and input information provided to the modelers for all environments included information about the sites (location, soil texture, field capacity, wilting point), management (sowing dates, sowing density, irrigation and fertilization dates and amounts), and daily weather data (precipitation, minimum and maximum air temperature, global radiation and potential evapotranspiration). Initial soil water and N content were not measured

4

165   in these experiments, but best guesses were provided by the experimental scientist. If any
166   models required other input data, modeling groups were asked to derive those values in
167   whatever way that seemed appropriate.

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

175



176

177       **Figure 1**

178       **Boxplots of calibration, evaluation and other data for development stages**
179   **BBCH30 and BBCH50 and varieties Apache and Bermude. The y-axis shows days from**
180   **sowing to the indicated development stage. Boxes indicate the lower and upper quartiles.**

5

181 **The solid line within the box is the median. Whiskers indicate the most extreme data**
182 **point, which is no more than 1.5 times the interquartile range from the box, and the**
183 **outlier dots are those observations that go beyond that range.**

184
185

## Crop models

187     Twenty-seven modeling groups participated in this study, noted M1-M27. The four
188 groups M2, M3, M4, and M5 all used the same model structure (i.e. models with the same
189 name), denoted as model structure S1. The four models M7, M12, M13, and M25 also shared
190 a common model structure, denoted as S2. As will be seen, different groups using the same
191 model structure had different results. The model versions for the same model structure
192 differed in some cases, but the differences are not in the basic phenology equations
193 implemented, and therefore, should have no or only a negligible effect on the simulated
194 development stages. Some differences in results may have resulted from different values for
195 the parameters that were not fit to the calibration data, although this was not recorded. The
196 differences in calibration approach were recorded, and this certainly led to differences in
197 simulated values. Since even groups using the same model structure obtained different results,
198 we refer to the 27 contributions as different models, In the presentation of the results the
199 models are anonymized and are identified simply as M1 to M27.  It would be misleading to
200 use the names of the model structures, since the results depend on both model structure and on
201 the approach to calibration. Information about the model structures is given in Supplementary
202 table S1.

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

208     In addition to the individual model results, we show the results for the model ensemble
209 mean ("e-mean") and the model ensemble median ("e-median"). We also define a very simple
210 predictor, denoted "naive", which was calculated as the average of the observations in the
211 calibration data for prediction. The naive model thus predicts that all days from sowing to
212 stage BBCH30 (BBCH55) will correspond to the average of days from sowing to BBCH30
213 (BBCH55) in the calibration data, separately for each variety. The naive model predictions for
214 days from sowing to BBCH30 and BBCH55 are respectively 155.9 days and 206.9 days for
215 variety Apache, and 156.1 days and 213.1 days for variety Bermude.

## Calibration and simulation experiment

217     The participants were provided with observed phenology data (dates of BBCH30 and
218 BBCH55) only for the calibration environments. The participants were asked to calibrate their
219 model using those data, and then to use the calibrated model to simulate phenology for all
220 environments (i.e. calibration, evaluation and hidden data environments). No guidelines for
221 calibration were provided. Participants were instructed to calibrate their model in their "usual
222 way" and fill out a questionnaire explaining what they did (Supplementary table S2).

## Evaluation

224     A common metric of error is mean squared error (MSE). We calculated MSE for each
225 model, each development stage (BBCH30 and BBCH55) and for each variety, as well as
226 averaged over stages and varieties. This was done separately for the calibration and evaluation

6

227  data. For example, MSE for model m, for predicting BBCH30, variety Apache, based on the
228  evaluation data, is:

229
$$MSE_{eval,m}^{BBCH30,Apache} = (1/8) \sum_{i \in eval} \left( y_i^{BBCH30,Apache} - \hat{y}_{i,m}^{BBCH30,Apache} \right)^2 \qquad (1)$$

230

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

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

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

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

235  32 predictions.

236      Mean squared error can be shown to be the sum of three positive terms, namely
237  squared bias, the difference in variance between the observed and simulated values and a term
238  related to the correlation between observed and simulated values (Kobayashi & Salam, 2000).
239  We specifically examined the bias, defined as the average over observed values minus the
240  average over simulated values.

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

244
$$MAE_{eval,m}^{BBCH30,Apache} = (1/8) \sum_{i \in eval} \left| y_i^{BBCH30,Apache} - \hat{y}_{i,m}^{BBCH30,Apache} \right|$$

245

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

247
$$EF_m = 1 - MSE_m / MSE_{naive}$$

248  where $MSE_m$ is MSE for model m and $MSE_{naive}$ is MSE for the naive model defined above.
249  EF is a skill measure, which compares the predictive capability of a model to that of the naive
250  model. Since the naive model makes the same prediction for all environments, it does not
251  account for any of the variability between environments. A model with  EF≤0 is a model that
252  does no better than the naive model, and so would be considered to be a very poor predictor.
253  A perfect model, with no error, has modeling efficiency of 1.

254

## Results

**Goodness-of-fit and prediction error**

257      Summary statistics for MSE averaged over both varieties and over both development
258  stages, for the calibration and evaluation data, are shown in table 2. Summary MSE values for
259  the calibration data for each development stage and variety separately are shown in
260  Supplementary table S7, and results for each individual model are given in Supplementary
261  figure S1.

262      **Table 2**

7

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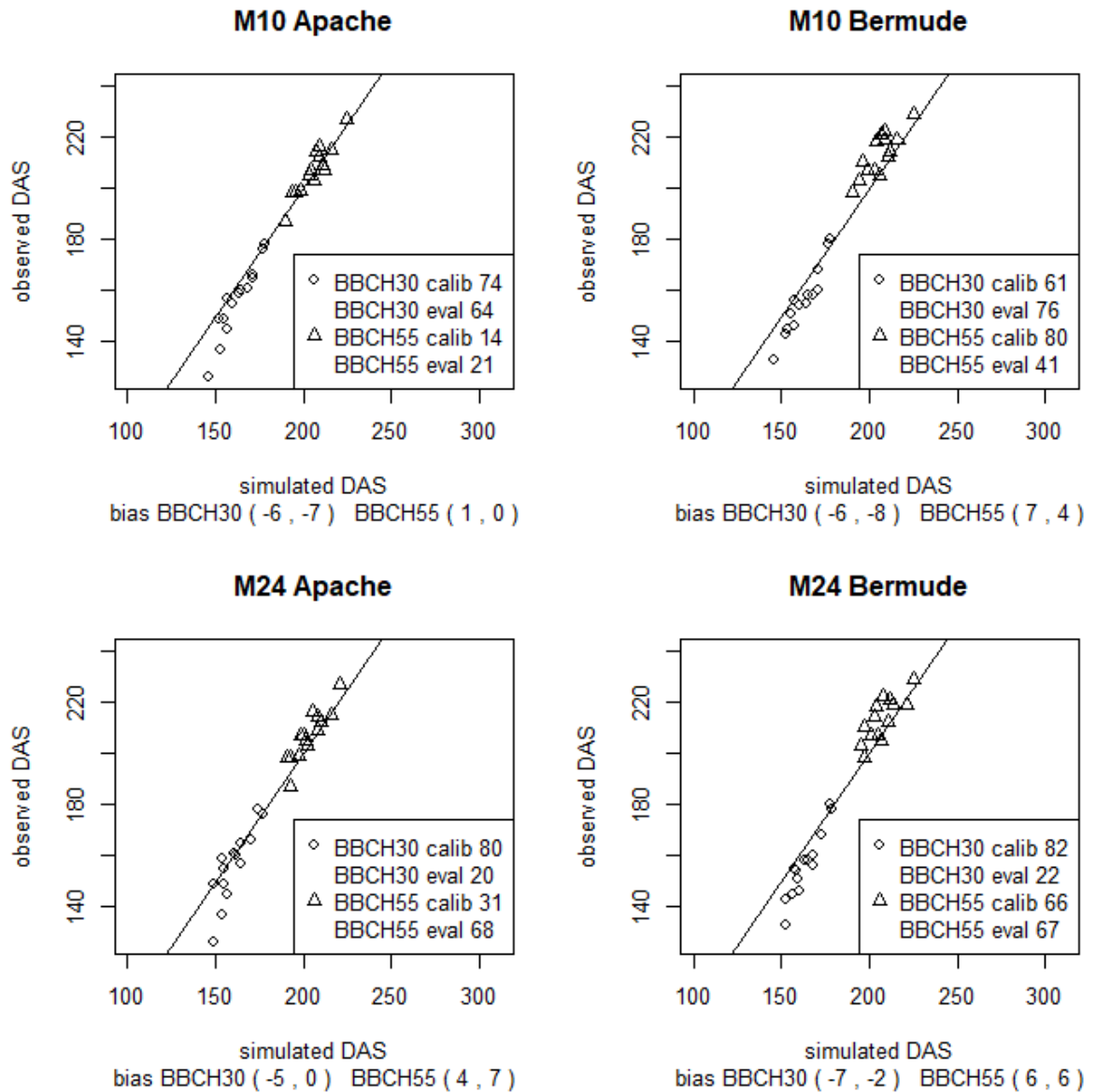264  **development stages.**

| MSE (days²) | Minimum | 1st quartile | Median | Mean | 3rd quartile | Maximum |
|---|---|---|---|---|---|---|
| Calibration data | 15 | 28 | 47 | 77 | 63 | 426 |
| Evaluation data | 20 | 35 | 62 | 79 | 111 | 235 |

265

266  In most cases, the bias for the calibration data is quite small. Considering absolute bias
267  for both development stages and both varieties, the median value was 2 days (Supplementary
268  table S8). In cases where the bias is relatively large, it is often of opposite sign for BBCH30
269  and BBCH55, as in the examples of figure 2.

270

271

**Figure 2**

**Observed vs. simulated days after sowing (DAS) for calibration data for models M10 and M24. The legend shows MSE (days²) for each stage and for calibration and evaluation data. (The individual evaluation results are not displayed). In the subtitles, bias values (days) for each stage are shown. The first number in parentheses is for the calibration data, the second number is for the evaluation data.**

Figure 3 and Supplementary tables S4-S6 show results for each development stage and variety and averaged over development stages and varieties for the evaluation data. Results for each model are given in Supplementary table S3. The median of MAE for the evaluation data is 6.1 days. The median of overall efficiency is 0.62, signifying that half of the models have MSE values for the evaluation data that are at most 38% as large as that of the naive predictor.   Only two models have negative values of EF, indicating that one would do better to predict using the average of the calibration data. For the four individual predictions (two

287     development stages, two varieties), the median of MAE ranges from 5.1 to 6.4 and the median
288     of EF ranges from 0.6 to 0.8. The ensemble models e-median and e-mean, though not the best
289     predictors, are among the best, with e-median being rated second best and e-mean fourth best.
290     The range of results among individual models is appreciable. The mean absolute errors for the
291     evaluation data averaged over all predictions ($MAE^{all}_{eval}$) go from 3.5 to 13 days. The $MSE^{all}_{eval}$
292     values vary by over a factor of 10, from a minimum of 20 days² to a maximum of 235 days².

293



294

295

296



297

298     **Figure 3**

299       **Box and whisker diagrams of absolute errors for evaluation data for each**
300       **prediction and on average (top panel) and modeling efficiency for each prediction and**
301       **on average (bottom panel). BBCH30A and BBCH30B refer respectively to prediction of**
302       **days to BBCH30 for variety Apache and variety Bermude. BBCH55A and BBCH55B**
303       **refer respectively to prediction of days to BBCH55 for variety Apache and variety**
304       **Bermude. The variability comes from differences between models.**

305

306

307       The relationship between overall MSE for the evaluation data and overall MSE for the
308       calibration data is quite close (adjusted $R^2$=0.70, figure 3). That is, much of the variability
309       between models in MSE for the evaluation data can be explained by the variability in the
310       calibration data, which further emphasizes the importance of calibration.

311       The four models that have model structure S1 and the four models that have model
312       structure S2 are identified in figure 4. Models with the same structure have different MSE
313       values; the differences are particularly large for S1. The models with structure S1 are ranked
314       3rd, 9th, 14th and 27th best for overall evaluation MSE among the 27 individual models. The
315       models with structure S2 are ranked 4th, 8th, 17th and 18[th] est.
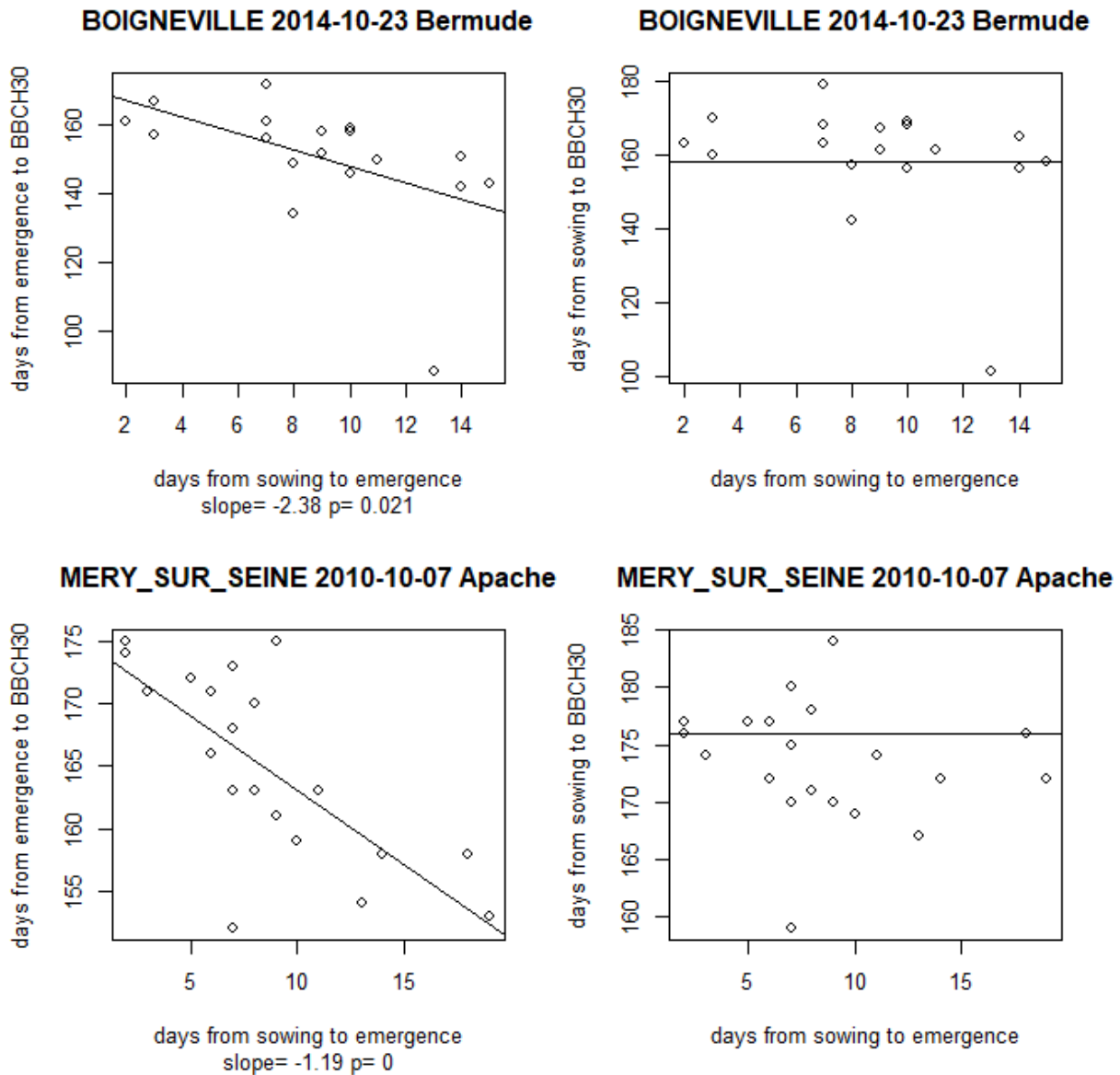
316

11

317

**Figure 4**

318

**Mean squared error (MSE) for the calibration data, averaged over environments,**
319
**development stages and varieties ($MSE_{calib}^{all}$ days²), as related to MSE for the evaluation**
320
**data ($MSE_{eval}^{all}$, days²). The regression line $MSE_{calib}^{all}$ = -27.6 + 1.32 $MSE_{eval}^{all}$ is shown**
321
**(R²=0.70). ● indicates models with structure S1. ＊ indicates models with structure S2. ○**
322
**indicates other models.**
323

324

Twenty-one models simulated and reported time from sowing to emergence. For these
325
models, we can separate simulated time from sowing to BBCH30 (sow_30) into two
326
contributions, the simulated time from sowing to emergence (sow_em) plus the simulated
327
time from emergence to BBCH30 (em_30), so that sow_30=sow_em+em_30. Figure 5 shows
328
results from two environments, typical of essentially all environments and both varieties, for
329
the relation between em_30 or sow_30 and sow_em. The average slope of the regression
330
em_30=a+b*sow_em over all environments (including calibration, evaluation and other
331

12

332 environments) and both varieties is b=-1.04, so that each day increase in simulated days to
333 emergence is on average associated with a 1.04 day decrease in simulated time from
334 emergence to BBCH30. The negative correlation between sow_em and em_30 leads to a
335 between-model variance for sow_30 (average variance 92 days²) that is smaller than the sum
336 of the variances of sow_em (average variance 20 days²) and em_30 (average variance 101
337 days²). The right panels of figure 5 show that different models could simulate almost exactly
338 the observed value of sow_30 with quite different values of sow_em.



339

340 **Figure 5**

341       **Relation between simulated days from emergence to BBCH30 and simulated days**
342 **from sowing to emergence as reported by 21 crop models for two environments (left**
343 **panels). Relation between simulated days from sowing to BBCH30 and simulated days**
344 **from sowing to emergence for the same environments (right panels). The slope of the**
345 **linear regression line and the p-value for testing slope=0 are shown for the left panels.**
346 **The observed days from sowing to BBCH30 is shown as a horizontal line in the right**
347 **panels.**

13

**Calibration approaches**

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

**Table 3**

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

| | | | | | | |
|---|---|---|---|---|---|---|
| Number of parameters[1] | 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 | | | | | |
| | Related to time to emergence 3 | | | | | |
| | Parameters unrelated to calibration data[2] 6 | | | | | |
| Objective function | Sum of squared errors or of root mean squared errors 21 | | | | | |
| | Sum of absolute errors 2 | | | | | |
| | Concentrated likelihood 1 | | | | | |
| | No single explicit objective function 3 | | | | | |
| Software[3] | Trial and error 10 | | | | | |
| | DIRECT-L (Gablonsky & Kelley, 2001; Johnson, n.d.) 2 | | | | | |
| | Ucode (E. P. Poeter, Hill, Banta, Mehl, & Christensen, 2005; Eileen P. Poeter & Hill, 1999) 3 | | | | | |
| | DE Optim (Mullen, Ardia, Gil, Windover, & Cline, 2011) 3 | | | | | |
| | PEST (Doherty, Hunt, & Tonkin, 2010) 2 | | | | | |
| | SCE (Duan, Gupta, & Sorooshian, 1993; Houska, Kraft, Chamorro-Chavez, & Breuer, 2015) 2 | | | | | |
| | GLUE (Beven & Binley, 2014; J. He, Jones, Graham, & Dukes, 2010) 1 | | | | | |
| | DREAM (J. A. Vrugt et al., 2009; Jasper A. Vrugt, 2016) 2 | | | | | |
| | Wrote code [4] 4 | | | | | |

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

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

374    Almost all modeling groups estimated one or more parameters that represent thermal
375    time between development stages (table 6). Some adjustments were necessary for models that
376    did not explicitly calculate time of BBCH30 or BBCH55. In model M2, for example, a new
377    parameter was added to the model, and estimated, representing the fraction of thermal time
378    from double ridge to heading at which BBCH30 occurs. Thirteen groups estimated a
379    parameter related to the effect of photoperiod. Ten groups estimated a parameter related to
380    vernalization. Six groups modified one or more parameters related to the temperature
381    response (for example model M6 estimated *Tbase*, the temperature below which there is no
382    development). Only three models modified parameters related to the time from sowing to
383    emergence, and only one model modified a parameter related to the effect of water stress. Six
384    models included among the parameters to estimate, parameters that have no effect on the
385    variables furnished as calibration data. Such parameters included thermal times for
386    development stages after BBCH55, potential kernel growth rate, kernel number per stem
387    weight and the temperature below which there is 50% death due to cold (Supplementary table
388    S9).

389    Most modeling groups defined the sum of squared errors or the sum of root mean
390    squared  errors  as the objective function to be minimized, where the sum is over the two
391    stages. (In all cases, the calibration was done separately for the two varieties). Two groups
392    minimized the sum of absolute errors. Calibration for model M21 was based on maximizing
393    the concentrated likelihood (Seber & Wild, 1989) assuming a normal distribution of errors
394    with possibly different error variances for the two development stages. In this case, the
395    objective function involves a product of errors for the two outputs, rather than a sum. Four of
396    the participants (M15, M16, M17, M18) did not define an explicit objective function to be
397    minimized. In these cases, the parameter values were chosen to obtain a "good fit" to the data
398    by visual inspection. Finally, two of the models (M7, M8) divided the calibration into two
399    steps. In these cases  three of the parameters were used to fit the BBCH30 data, and then in
400    another step another parameter was used to fit the BBCH55 data.

401    Minimizing the sum of squared errors is a standard statistical approach to model
402    calibration, which has highly desirable properties if certain assumptions about model error are
403    satisfied, including equal variance of model error for all data points and non-correlation of
404    model errors. Only one model took into account the possibility that the error variances are
405    different for BBCH30 and BBCH55, and none of the modeling groups took into account
406    possible correlations between errors for BBCH30 and BBCH55 in the same field. Based on
407    the errors for all the data and all the models, it was found that there is a highly significant
408    difference in variance between errors for BBCH30 (variance of error 100.7 days²) and
409    BBCH55 (variance of error 67.3 days²). Also, the correlation between the error for BBCH30
410    and the error for BBCH55 in the same field is 0.53 and highly significant. However, if only
411    results for a single model are considered, then for most models the differences in variance and
412    the correlation are not significant.

15

413        Two models defined a posterior probability of the parameters equal to the likelihood
414    times the prior probability, as usually assumed in a Bayesian approach. The parameters used
415    for prediction were those that maximized the posterior probability (i.e., the estimated mode of
416    the posterior distribution). In both cases, the likelihood was assumed Gaussian with
417    independent errors, and the prior distribution was assumed uniform between some minimum
418    and maximum value. This approach is equivalent to minimizing the mean squared error, with
419    constraints on the parameter values.

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

424

## Discussion

425

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

432        Twenty-seven modeling groups participated in the exercise. Most models predicted
433    times to stem elongation and heading quite well (median MAE of 6 days). Half the models
434    had MSE values of prediction that were 36% or less than MSE of a naive predictor. It must be
435    kept in mind that this study is a rather favorable situation for prediction, with a substantial
436    amount of calibration data and predictions for environments similar to those of the calibration
437    data.

438        The results for each individual model depend on the model equations, the values of the
439    fixed parameters and the calibration approach which determines the values of the estimated
440    parameters, and on the interactions among those elements. A full, detailed analysis of each
441    model is beyond the objectives of this study. We focused on the calibration approaches used
442    and its relation to errors for the calibration and evaluation data.

443        The results show that calibration can, in some cases, compensate for differences in
444    model equations and/or values of fixed parameters. Compensation is usually discussed in the
445    context of single models. For example, equifinality, which is a well-known phenomenon of
446    complex models, means that different combinations of parameter values, and thus different
447    quantitative descriptions of processes, can lead to the same results for outputs because there is
448    compensation between the processes (Beven, 2006; D. He et al., 2017). However, this
449    phenomenon has not been described in the context of multi-model studies. Here, we have an
450    example of compensation for differences between models in the way they partition days from
451    sowing to BBCH30 into days from sowing to emergence plus days from emergence to
452    BBCH30. Models with longer simulated times from sowing to emergence tend to have a
453    shorter simulated time from emergence to development stage BBCH30 and vice versa. In fact,
454    each extra day from sowing to emergence is associated on average with almost exactly one
455    less day from emergence to BBCH30.  The result is that models with quite different simulated
456    days from sowing to emergence can have nearly identical times from sowing to BBCH30.
457    This can be expressed in terms of model uncertainty, as quantified by between-model

16

458    variance. The variance of days from sowing to BBCH30 is less than the sum of variances of
459    days from sowing to emergence and days from emergence to BBCH30. That is, calibration
460    reduces, but does not eliminate, model uncertainty for the variable provided for calibration.

461          We don't have observed time to emergence, but in any case the models with different
462    simulated days to emergence can't all be right. This is an  example of how models can get the
463    right answer (correct days to BBCH30) for the wrong reasons (wrong days to emergence),
464    illustrating the problem pointed out for example by  (Challinor, Martre, Asseng, Thornton, &
465    Ewert, 2014). The same compensation of errors between sowing to emergence and emergence
466    to BBCH30 will not be appropriate for all environments. This is one of the main reasons that
467    extrapolation to populations different than the training population is dangerous.

468          Another indication of compensation induced by calibration is the fact that after
469    calibration, models with quite different choices for the variables that affect development can
470    have very similar levels of prediction error. The most important inputs that determine spring
471    wheat phenology are daily temperature and photoperiod (Aslam et al., 2017) and for winter
472    wheat it is also necessary to include the process of vernalization, i.e. the effect of low winter
473    temperatures on development (Li et al., 2013). Five of the best eight predicting models here,
474    with $MSE_{eval}^{all} < 40$ days², do use all three of those variables (daily temperature, photoperiod,
475    vernalizing temperatures) as inputs. Two of the best eight models however do not use
476    vernalizating temperatures, and one of those best eight does not use photoperiod. The choice
477    of input variables is a fundamental aspect of model structure. In fact, MSE can be expressed
478    as a sum of two terms, the first of which depends only on the choice of the model input
479    variables and not on any other aspects of structure, while the second measures the distance
480    between the model used and the optimal model for those inputs (Wallach, Makowski, Jones,
481    & Brun, 2019). It seems that calibration can lead to similar values of MSE for prediction even
482    for quite different choices of input variables.

483          While models with different structures can give similar results thanks to calibration,
484    our results also show that models with the same structure can provide different results, if
485    different calibration approaches are used. This is illustrated here by the results for two groups
486    of models sharing the same structure. There are major differences in prediction error between
487    models with the same structure depending on how calibration was done. Much previous work
488    on improving the predictive capability of crop models has focused on the model equations, for
489    instance the way temperature is taken into account in various processes (Maiorano et al.,
490    2016; Wang et al., 2017). Here we show that models with quite different structures can have
491    very similar prediction accuracy, thanks to calibration using the same data, while models with
492    the same model structure can have very different levels of prediction error, if the calibration
493    methods differ. This means that model comparison studies may often be comparing
494    calibration approaches as much or more as they are comparing model equations. This is in
495    line with the conclusions of Confalonieri et al. (2016), who argued that one should not speak
496    of evaluation of a model but rather of a model-user combination; a major role of  the user is in
497    determining the method of calibration.

498          The choice of objective function for calibration can have an effect on quality of
499    prediction.  While most participants defined an explicit objective function (e.g. minimizing
500    sum of squared errors, or some closely related criterion) three models (see Supplementary
501    table S9) did not have an explicit quantitative objective function. Those models all had
502    relatively large values of  overall MSE for the evaluation data ( $MSE_{eval}^{all}$ ) , having 15th, 16th,
503    and 18[th] largest $MSE_{eval}^{all}$ values out of the 25 models that predicted both BBCH30 and

17

504 BBCH55. These results suggest that the lack of a quantitative objective function can be a
505 drawback since then one does not have a clear criterion for judging the results of calibration.

506 There was a large diversity of choices of parameters to estimate by calibration, and
507 this had in certain cases an important effect on prediction error. One rather unexpected
508 observation was that several participants included parameters that have no effect on the
509 variables furnished as calibration data among the parameters to estimate. The data cannot in
510 those cases give any information about the parameter value. At best, including such
511 parameters among the parameters to estimate is useless, and those parameters will simply
512 have final values exactly equal to their initial values. However, there may also be serious
513 disadvantages to including such parameters. It gives the erroneous impression that one is
514 estimating parameters that cannot in fact be estimated, it increases computation time and it
515 can cause problems for the parameter estimation algorithm. The very poor fit of model M5 to
516 the calibration data seems to be directly related to the fact that for this model, several
517 parameters unrelated to the calibration data were chosen to be fitted. The software used here
518 was PEST (Doherty et al., 2010), with the singular value decomposition option, which allows
519 one to deal with non-estimable parameters, but at the cost of introducing bias in the estimated
520 parameter values. Obviously, one should not include non-estimable parameters among the
521 parameters to estimate.

522 The choice of parameters to estimate may be the principal cause of bias in fitting the
523 calibration data for some models. If a model includes an additive constant term, and squared
524 error is minimized, bias will be 0 for the calibration data. Even for more complex models,
525 calibration can bring bias close to 0, as illustrated here by the fact that many of the models
526 had very small biases for the calibration data. Eliminating bias is important, since squared
527 bias is one component of MSE, and therefore the bias necessarily adds on to MSE (Kobayashi
528 & Salam, 2000). If one does not have a parameter with a nearly additive effect for each of the
529 development stages BBCH30 and BBCH55, the elimination of bias for both outputs is not
530 assured. Model M24 estimated only a single parameter. In such a case, at best one can
531 estimate a parameter value that gives the best compromise between errors in BBCH30 and
532 BBCH55. This may lead to a negative bias for one of those outputs and more or less
533 corresponding positive bias for the other. This is exactly the behavior illustrated in figure 2.
534 Model M10 also had fairly large biases. Here three parameters were estimated, but one is
535 unrelated to the observed data and a second concerns time to emergence, which was only
536 allowed to vary in a limited range. Apparently in this case also there was not enough
537 flexibility to eliminate bias for both development stages. Models with large bias for the
538 calibration data tended to have large MSE values for the evaluation data (Supplementary
539 figure S2). This suggests that the parameters to estimate should include one parameter that is
540 nearly additive (i.e. that adds an amount that is nearly the same for all environments) for each
541 observed output, and that is not too limited in the allowed range of values.

542 The calibration choices here suggest other recommendations for calibration. One
543 concerns the specific choice for the objective function. Among the models that defined a
544 likelihood or a sum of squares criterion, all but one assumed that all model errors had equal
545 variance and were independent. One should probably take into account unequal variances and
546 correlation of simulation errors for BBCH30 and BBCH55 in the same field. A second
547 recommendation concerns the software. There does not seem to be any clear connection
548 between the software used for calibration and the predictive quality of the resulting calibrated
549 model. Various different software solutions were used by the best predicting models, but
550 largely the same software solutions were also found among the models with the largest
551 prediction errors. A problem that may arise concerns the test for convergence to the parameter

18

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

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

## Acknowledgements

# References

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

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

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

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

Asseng, Senthold, Martre, P., Maiorano, A., Rötter, R. P., O'Leary, G. J., Fitzgerald, G. J., … Ewert, F. (2019). Climate change impact and adaptation for wheat protein. *Global Change Biology*, *25*(1), 155–173. https://doi.org/10.1111/gcb.14481

Beven, K. (2006). A manifesto for the equifinality thesis. *Journal of Hydrology*, *320*(1–2), 18–36.

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

Bindi, M., Palosuo, T., Trnka, M., & Semenov, M. (2015). Modelling climate change impacts on crop production for food security. *Climate Research*, *65*(September), 3–5. https://doi.org/10.3354/cr01342

Campbell, B. M., Hansen, J., Rioux, J., Stirling, C. M., Twomlow, S., & (Lini) Wollenberg, E. (2018). Urgent action to combat climate change and its impacts (SDG 13): transforming agriculture and food systems. *Current Opinion in Environmental Sustainability*, *34*, 13–20. https://doi.org/10.1016/J.COSUST.2018.06.005

Carberry, P. S., Liang, W., Twomlow, S., Holzworth, D. P., Dimes, J. P., McClelland, T., … Keating, B. A. (2013). Scope for improved eco-efficiency varies among diverse cropping systems. *Proceedings of the National Academy of Sciences of the United States of America*, *110*(21), 8381–8386. https://doi.org/10.1073/pnas.1208050110

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

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

Doherty, J. E., Hunt, R. J., & Tonkin, M. J. (2010). *Approaches to highly parameterized inversion: A guide to using PEST for model-parameter and predictive-uncertainty*

642    *analysis: U.S. Geological Survey Scientific Investigations Report 2010–5211*. Retrieved
643    from http://pubs.usgs.gov/sir/2010/5211

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

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

649  Godfray, H. C. J., Beddington, J. R., Crute, I. R., Haddad, L., Lawrence, D., Muir, J. F., …
650    Toulmin, C. (2010). Food security: the challenge of feeding 9 billion people. *Science*
651    *(New York, N.Y.)*, *327*(5967), 812–818. https://doi.org/10.1126/science.1185383

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

656  He, J., Jones, J. W., Graham, W. D., & Dukes, M. D. (2010). Influence of likelihood function
657    choice for estimating crop model parameters using the generalized likelihood uncertainty
658    estimation    method.    *Agricultural    Systems*,    *103*(5),    256–264.
659    https://doi.org/10.1016/j.agsy.2010.01.006

660  Hochman, Z., Gobbett, D. L., & Horan, H. (2017). Climate trends account for stalled wheat
661    yields    in    Australia    since    1990.    *Global    Change    Biology*,    *23*(5),    2071–2081.
662    https://doi.org/10.1111/gcb.13604

663  Houska, T., Kraft, P., Chamorro-Chavez, A., & Breuer, L. (2015). SPOTting Model
664    Parameters Using a Ready-Made Python Package. *PLOS ONE*, *10*(12), e0145180.
665    https://doi.org/10.1371/journal.pone.0145180

666  Hunt, J. R., Lilley, J. M., Trevaskis, B., Flohr, B. M., Peake, A., Fletcher, A., … Kirkegaard,
667    J. A. (2019). Early sowing systems can boost Australian wheat yields despite recent
668    climate change. *Nature Climate Change*, *9*(3), 244–247. https://doi.org/10.1038/s41558-
669    019-0417-9

670  Hussain, J., Khaliq, T., Ahmad, A., & Akhtar, J. (2018). Performance of four crop model for
671    simulations of wheat phenology, leaf growth, biomass and yield across planting dates.
672    *PLOS ONE*, *13*(6), e0197546. https://doi.org/10.1371/journal.pone.0197546

673  IPCC. (2014). *Climate Change 2014: Synthesis Report. Contribution of Working Groups I, II*
674    *and III to the Fifth Assessment Report of the Intergovernmental Panel on Climate*
675    *Change*. (P. K. Pachauri & A. Meyer, Eds.). Geneva: IPCC.

676  Johnson, S. G. (n.d.). The NLopt nonlinear-optimization package.

677  Kobayashi, K., & Salam, M. U. (2000). Comparing simulated and measured values using
678    mean squared deviation and its components. *Agronomy Journal*, *92*, 345–352.

679  Li, G., Yu, M., Fang, T., Cao, S., Carver, B. F., & Yan, L. (2013). Vernalization requirement
680    duration in winter wheat is controlled by TaVRN-A1 at the protein level. *The Plant*
681    *Journal*:    *For    Cell    and    Molecular    Biology*,    *76*(5),    742–753.
682    https://doi.org/10.1111/tpj.12326

Maiorano, A., Martre, P., Asseng, S., Ewert, F., Müller, C., Rötter, R. P., … Zhu, Y. (2017). Crop model improvement reduces the uncertainty of the response to temperature of multi-model ensembles. *Field Crops Research*, *202*. https://doi.org/10.1016/j.fcr.2016.05.001

Maiorano, Andrea, Martre, P., Asseng, S., Ewert, F., Müller, C., Rötter, R. P., … Zhu, Y. (2016). Crop model improvement reduces the uncertainty of the response to temperature of multi-model ensembles. *Field Crops Research*. https://doi.org/10.1016/j.fcr.2016.05.001

Mullen, K., Ardia, D., Gil, D., Windover, D., & Cline, J. (2011). "DEoptim": An R Package for Global Optimization by Differential Evolution. *Journal of Statistical Software*, *40*, 1–26.

Parry, M. ., Rosenzweig, C., Iglesias, A., Livermore, M., & Fischer, G. (2004). Effects of climate change on global food production under SRES emissions and socio-economic scenarios. *Global Environmental Change*, *14*(1), 53–67. https://doi.org/10.1016/j.gloenvcha.2003.10.008

Piao, S., Liu, Q., Chen, A., Janssens, I. A., Fu, Y., Dai, J., … Zhu, X. (2019). Plant phenology and global climate change: current progresses and challenges. *Global Change Biology*, gcb.14619. https://doi.org/10.1111/gcb.14619

Poeter, E. P., Hill, M. C., Banta, E. R., Mehl, S., & Christensen, S. (2005). *UCODE_2005 and Six Other Computer Codes for Universal Sensitivity Analysis, Calibration, and Uncertainty Evaluation: U.S. Geological Survey Techniques and Methods 6-A11.*

Poeter, Eileen P., & Hill, M. C. (1999). UCODE, a computer code for universal inverse modeling. *Computers & Geosciences*, *25*(4), 457–462. https://doi.org/10.1016/S0098-3004(98)00149-6

Porter, J. R., Howden, M., & Smith, P. (2017). Considering agriculture in IPCC assessments. *Nature Climate Change*, *7*(10), 680–683. https://doi.org/10.1038/nclimate3404

Rezaei, E. E., Siebert, S., & Ewert, F. (2015). Intensity of heat stress in winter wheat—phenology compensates for the adverse effect of global warming. *Environmental Research Letters*, *10*(2), 024012. https://doi.org/10.1088/1748-9326/10/2/024012

Rezaei, E. E., Siebert, S., Hüging, H., & Ewert, F. (2018). Climate change effect on wheat phenology depends on cultivar change. *Scientific Reports*, *8*(1), 4891. https://doi.org/10.1038/s41598-018-23101-2

Seber, G. A. F., & Wild, C. J. (1989). *Nonlinear regression*. New York: Wiley .

Seidel, S. J., Palosuo, T., Thorburn, P., & Wallach, D. (2018). Towards improved calibration of crop models – Where are we now and where should we go? *European Journal of Agronomy*, *94*, 25–35. https://doi.org/10.1016/J.EJA.2018.01.006

Svystun, T., Bhalerao, R. P., & Jönsson, A. M. (2019). Modelling Populus autumn phenology: The importance of temperature and photoperiod. *Agricultural and Forest Meteorology*, *271*, 346–354. https://doi.org/10.1016/J.AGRFORMET.2019.03.003

Vrugt, J. A., ter Braak, C. J. F., Diks, C. G. H., Robinson, B. A., Hyman, J. M., & Higdon, D. (2009). Accelerating Markov Chain Monte Carlo Simulation by Differential Evolution with Self-Adaptive Randomized Subspace Sampling. *International Journal of Nonlinear*

22

725 *Sciences and Numerical Simulation*, *10*(3).
726     https://doi.org/10.1515/IJNSNS.2009.10.3.273

727 Vrugt, Jasper A. (2016). Markov chain Monte Carlo simulation using the DREAM software
728     package: Theory, concepts, and MATLAB implementation. *Environmental Modelling &*
729     *Software*, *75*, 273–316. https://doi.org/10.1016/J.ENVSOFT.2015.08.013

730 Wallach, D., Makowski, D., Jones, J. W., & Brun, F. (2019). *Working with Dynamic Crop*
731     *Models: Methods, Tools and examples for Agriculture and Environment*. London, U.K.:
732     Academic Press.

733 Wang, E., Martre, P., Zhao, Z., Ewert, F., Maiorano, A., Rötter, R. P., … Asseng, S. (2017).
734     The uncertainty of crop yield projections is reduced by improved temperature response
735     functions. *Nature Plants*, *3*. https://doi.org/10.1038/nplants.2017.102

736 Went, F. W. (2003). The Effect of Temperature on Plant Growth. *Annual Review of Plant*
737     *Physiology*, *4*(1), 347–362. https://doi.org/10.1146/annurev.pp.04.060153.002023

738 Wheeler, T., & von Braun, J. (2013). Climate change impacts on global food security. *Science*
739     *(New York, N.Y.)*, *341*(6145), 508–513. https://doi.org/10.1126/science.1239402

740 Yuan, S., Peng, S., & Li, T. (2017). Evaluation and application of the ORYZA rice model
741     under different crop managements with high-yielding rice cultivars in central China.
742     *Field Crops Research*, *212*, 115–125. https://doi.org/10.1016/J.FCR.2017.07.010

743 Zadoks, J. C., Chzang, T. T., & Konzak, C. F. (1974). A decimal code for the growth stages
744     of cereals. *Weed Research*, *14*(6), 415–421. https://doi.org/10.1111/j.1365-
745     3180.1974.tb01084.x

746

747