

Multimodel ensembles of wheat growth: many models are better than one

PIERRE MARTRE^{1,2}, DANIEL WALLACH³, SENTHOLD ASSENG⁴, FRANK EWERT⁵, JAMES W. JONES⁴, REIMUND P. RÖTTER⁶, KENNETH J. BOOTE⁴, ALEX C. RUANE⁷, PETER J. THORBURN⁸, DAVIDE CAMMARANO⁴, JERRY L. HATFIELD⁹, CYNTHIA ROSENZWEIG⁷, PRAMOD K. AGGARWAL¹⁰, CARLOS ANGULO⁵, BRUNO BASSO¹¹, PATRICK BERTUZZI¹², CHRISTIAN BIERNATH¹³, NADINE BRISSON^{14,15†}, ANDREW J. CHALLINOR^{16,17}, JORDI DOLTRA¹⁸, SEBASTIAN GAYLER¹⁹, RICHIE GOLDBERG⁷, ROBERT F. GRANT²⁰, LEE HENG²¹, JOSH HOOKER²², LESLIE A. HUNT²³, JOACHIM INGWERSEN²⁴, ROBERTO C. IZAURRALDE²⁵, KURT CHRISTIAN KERSEBAUM²⁶, CHRISTOPH MÜLLER²⁷, SOORA NARESH KUMAR²⁸, CLAAS NENDEL²⁶, GARRY O'LEARY²⁹, JØRGEN E. OLESEN³⁰, TOM M. OSBORNE³¹, TARU PALOSUO⁶, ECKART PRIESACK¹³, DOMINIQUE RIPOCHE¹², MIKHAIL A. SEMENOV³², IURII SHCHERBAK¹¹, PASQUALE STEDUTO³³, CLAUDIO O. STÖCKLE³⁴, PIERRE STRATONOVITCH³², THILO STRECK²⁴, IWAN SUPIT³⁵, FULU TAO³⁶, MARIA TRAVASSO³⁷, KATHARINA WAHA²⁷, JEFFREY W. WHITE³⁸ and JOOST WOLF³⁹

¹INRA, UMR1095 Genetics, Diversity and Ecophysiology of Cereals (GDEC), 5 chemin de Beaulieu, F-63 100 Clermont-Ferrand, France, ²Blaise Pascal University, UMR1095 GDEC, F-63 170 Aubière, France, ³INRA, UMR1248 Agrosystèmes et Développement Territorial, F-31 326 Castanet-Tolosan, France, ⁴Agricultural & Biological Engineering Department, University of Florida, Gainesville, FL 32611, USA, ⁵Institute of Crop Science and Resource Conservation, Universität Bonn, D-53 115 Bonn, Germany, ⁶Plant Production Research, MTT Agrifood Research Finland, FI-50 100 Mikkeli, Finland, ⁷National Aeronautics and Space Administration, Goddard Institute for Space Studies, New York, NY 10025, USA, ⁸Commonwealth Scientific and Industrial Research Organization, Ecosystem Sciences, Dutton Park, QLD 4102, Australia, ⁹National Laboratory for Agriculture and Environment, Ames, IA 50011, USA, ¹⁰Consultative Group on International Agricultural Research, Research Program on Climate Change, Agriculture and Food Security, International Water Management Institute, New Delhi 110012, India, ¹¹Department of Geological Sciences and Kellogg Biological Station, Michigan State University, East Lansing, MI 48823, USA, ¹²INRA, US1116 AgroClim, F-84 914 Avignon, France, ¹³Institute of Soil Ecology, Helmholtz Zentrum München, German Research Center for Environmental Health, Neuherberg D-85 764, Germany, ¹⁴INRA, UMR0211 Agronomie, F-78 750 Thiverval-Grignon, France, ¹⁵AgroParisTech, UMR0211 Agronomie, F-78 750 Thiverval-Grignon, France, ¹⁶Institute for Climate and Atmospheric Science, School of Earth and Environment, University of Leeds, Leeds LS29JT, UK, ¹⁷CGIAR-ESSP Program on Climate Change, Agriculture and Food Security, International Centre for Tropical Agriculture, A.A. 6713 Cali, Colombia, ¹⁸Cantabrian Agricultural Research and Training Centre, 39600 Muriedas, Spain, ¹⁹Water & Earth System Science Competence Cluster, c/o University of Tübingen, D-72 074 Tübingen, Germany, ²⁰Department of Renewable Resources, University of Alberta, Edmonton, AB T6G 2E3, Canada, ²¹International Atomic Energy Agency, 1400 Vienna, Austria, ²²School of Agriculture, Policy and Development, University of Reading, RG6 6AR Reading, UK, ²³Department of Plant Agriculture, University of Guelph, Guelph, ON N1G 2W1, Canada, ²⁴Institute of Soil Science and Land Evaluation, Universität Hohenheim, D-70 599 Stuttgart, Germany, ²⁵Department of Geographical Sciences, University of Maryland, College Park, MD 20782, USA, ²⁶Institute of Landscape Systems Analysis, Leibniz Centre for Agricultural Landscape Research, D-15 374 Müncheberg, Germany, ²⁷Potsdam Institute for Climate Impact Research, D-14 473 Potsdam, Germany, ²⁸Centre for Environment Science and Climate Resilient Agriculture, Indian Agricultural Research Institute, New Delhi 110 012, India, ²⁹Department of Primary Industries, Landscape & Water Sciences, Horsham, Vic., 3400, Australia, ³⁰Department of Agroecology, Aarhus University, 8830 Tjele, Denmark, ³¹National Centre for Atmospheric Science, Department of Meteorology, University of Reading, RG6 6BB Reading, UK, ³²Computational and Systems Biology Department, Rothamsted Research, Harpenden, Herts AL5 2JQ, UK, ³³Food and Agriculture Organization of the United Nations, Rome 00153, Italy, ³⁴Biological Systems Engineering, Washington State University, Pullman, WA 99164-6120, USA, ³⁵Earth System Science-Climate Change, Wageningen University, 6700AA Wageningen, The Netherlands, ³⁶Institute of Geographical Sciences and Natural Resources Research, Chinese Academy of Science, Beijing 100101, China, ³⁷Institute for Climate and Water, INTA-CIRN, 1712 Castelar, Argentina, ³⁸Arid-Land Agricultural Research Center, USDA, Maricopa, AZ 85138, USA, ³⁹Plant Production Systems, Wageningen University, 6700AA Wageningen, The Netherlands

Correspondence: Pierre Martre, tel. +33 473 624 351,
fax +33 473 624 457, e-mail: pierre.martre@clermont.inra.fr

†Dr Nadine Brisson passed away in 2011 while this work was being carried out.

Abstract

Crop models of crop growth are increasingly used to quantify the impact of global changes due to climate or crop management. Therefore, accuracy of simulation results is a major concern. Studies with ensembles of crop models can give valuable information about model accuracy and uncertainty, but such studies are difficult to organize and have only recently begun. We report on the largest ensemble study to date, of 27 wheat models tested in four contrasting locations for their accuracy in simulating multiple crop growth and yield variables. The relative error averaged over models was 24–38% for the different end-of-season variables including grain yield (GY) and grain protein concentration (GPC). There was little relation between error of a model for GY or GPC and error for in-season variables. Thus, most models did not arrive at accurate simulations of GY and GPC by accurately simulating preceding growth dynamics. Ensemble simulations, taking either the mean (e-mean) or median (e-median) of simulated values, gave better estimates than any individual model when all variables were considered. Compared to individual models, e-median ranked first in simulating measured GY and third in GPC. The error of e-mean and e-median declined with an increasing number of ensemble members, with little decrease beyond 10 models. We conclude that multimodel ensembles can be used to create new estimators with improved accuracy and consistency in simulating growth dynamics. We argue that these results are applicable to other crop species, and hypothesize that they apply more generally to ecological system models.

Keywords: ecophysiological model, ensemble modeling, model intercomparison, process-based model, uncertainty, wheat (*Triticum aestivum* L.)

Introduction

Global change with increased climatic variability are projected to strongly impact crop and food production, but the magnitude and trajectory of these impacts remain uncertain (Tubiello *et al.*, 2007). This uncertainty, together with the increasing demand for food of a growing world population (Bloom, 2011), has raised concerns about food security and the need to develop more sustainable agricultural practices (Godfray *et al.*, 2010). More confident understanding of global change impacts is needed to develop effective adaptation and mitigation strategies (Easterling *et al.*, 2007). Methodologies to quantify global change impacts on crop production include statistical models (Lobell *et al.*, 2011) and process-based crop simulation models (Porter & Semenov, 2005), which are increasingly used in basic and applied research and to support decision making at different scales (Challinor *et al.*, 2009; Ko *et al.*, 2010; Angulo *et al.*, 2013; Rosenzweig *et al.*, 2013).

Different crop growth and development processes are affected by climatic variability via linear or non-linear relationships resulting in complex and unexpected responses (Trewavas, 2006). It has been argued that such responses can best be captured by process-based crop simulation models that quantitatively represent the interaction and feedback responses of crops to their environments (Porter & Semenov, 2005; Bertin *et al.*, 2010). Wheat is the most important staple crop in the world providing over

20% of the calories and proteins in human diet (FAOSTAT, 2014). It has therefore received much attention from the crop modeling community and over 40 wheat crop models are in use (White *et al.*, 2011). These differ in the processes included in the models and the mechanistic detail used to model individual processes like evapotranspiration or photosynthesis. Therefore, a thorough comparative evaluation of models is essential to understand the reliability of model simulations and to quantify and reduce the uncertainty of such simulations (Rötter *et al.*, 2011).

The Wheat Pilot study (Asseng *et al.*, 2013) of the Agricultural Model Intercomparison and Improvement Project (AgMIP; Rosenzweig *et al.*, 2013) compared 27 wheat models, the largest ensemble of crop models created to date. The models vary greatly in their complexity and in the modeling approaches and equations used to represent the major physiological processes that determine crop growth and development and their responses to environmental factors, see Table S3 in Asseng *et al.* (2013).

An initial study (Asseng *et al.*, 2013) analyzed the variability between crop models in simulating grain yield (GY) under climate change situations without specifically investigating multimodel ensemble estimators considering other end-of-season and in-season variables to better justify their possible application. The present analysis uses the resulting dataset to study how the multimodel ensemble average or median can reproduce in-season and end-of-season observations. In its

simplest and most common form, a multimodel ensemble simulation is produced by averaging the simulations of member models weighted equally (Knutti, 2010). This method has been practiced in climate forecasting (Räisänen & Palmer, 2001; Hagedorn *et al.*, 2005) and in ecological modeling of species distribution (Grenouillet *et al.*, 2011), and it has been shown that multimodel ensembles can give better estimates than any individual model. Such improvement in skill of a multimodel ensemble may be also applicable to crop models. Preliminary evidence suggests that the average of ensembles of simulations is a good estimator of GY for several crops (Palosuo *et al.*, 2011; Rötter *et al.*, 2012; Bassu *et al.*, 2014) and possibly even better than the best individual model across different seasons and sites (Rötter *et al.*, 2012). However, a detailed quantitative analysis of the quality of simulators based on crop model ensembles, compared to individual models is lacking. By looking at outputs of multiple growth variables (both in-season and end-of-season), we would get a broader picture of how ensemble estimators perform and a better understanding of why they perform well compared to individual models. It is important therefore to consider not only GY but also other growth variables. If multimodel ensembles are truly more skillful than the best model in the ensemble, or even simply better than the average of the models, then using ensemble medians or means may be a powerful estimator to evaluating crop response to crop management and environmental factors.

Model evaluations can give quite different results depending on the use of the model that is studied. Here, we investigate the situation where models are applied in environments for which they have not been specifically calibrated, which is typically the situation in global impact studies (Rosenzweig *et al.*, 2014). The model results were compared to measured data from four contrasting growing environments. The modeling groups were provided with weather data, soil characteristics, soil initial conditions, management, and flowering and harvest dates for each site. Although only four locations were tested in the AgMIP Wheat Pilot study, this limitation is partially compensated for by the diversity of the sites ranging from high to low yielding, from short to long season, and irrigated and not irrigated situations.

Two main approaches to evaluate the accuracy and uncertainty of the AgMIP wheat model ensemble were followed. First, we evaluated the range of errors and the average error of the models for multiple growth variables, including both in-season and end-of-season variables. Secondly, we evaluated two ensemble-based models, the mean (e-mean) and the median (e-median) of the simulated values of the

ensemble members. Finally, we studied how the error of e-mean and e-median changed with the size of the ensemble.

Materials and methods

Experimental data

Quality-assessed experimental data from single crops at four contrasting locations representing diverse agro-ecological conditions were used. The locations were Wageningen, The Netherlands (NL; Groot & Verberne, 1991), Balcarce, Argentina (AR; Travasso *et al.*, 2005), New Delhi, India (IN; Naveen, 1986), and Wongan Hills, Australia (AU; Asseng *et al.*, 1998). Typical regional crop management was used at each site. In all experiments, the plots were kept weed-free, and plant protection methods were used as necessary to minimize damage from pests and diseases. Crop management and soil and cultivar information, as given to each individual modeling group, are given in Table 1.

Daily values of solar radiation, maximum and minimum temperature, and precipitation were recorded at weather stations at or near the experimental plots, except for IN solar radiation which was obtained from the NASA POWER dataset of modeled data (Stackhouse, 2006) that extends back to 1983. Daily values of 2-m wind speed (m s^{-1}), dew point temperature ($^{\circ}\text{C}$), vapor pressure (hPa), and relative humidity (%) were estimated for each location from the NASA Modern Era Retrospective-Analysis for Research and Applications (Bosilovich *et al.*, 2011), except for NL wind speed and vapor pressure that were measured on site. Air CO_2 concentration was taken to be 360 ppm at all sites. A weather summary for each site is shown in Table 1 and Fig. 1.

For all sites, end-of-season (i.e. ripeness-maturity) values for GY (t DM ha^{-1}), total aboveground biomass (AGBM_{m} , t DM ha^{-1}), total aboveground nitrogen (AGN_{m} , kg N ha^{-1}), and grain N (GN_{m} , kg N ha^{-1}) were available. From these values, biomass harvest index ($\text{HI} = 100 \times \text{GY}/\text{AGBM}_{\text{m}}$, %), N harvest index ($\text{NHI} = 100 \times \text{GN}_{\text{m}}/\text{AGN}_{\text{m}}$, %), and grain protein concentration ($\text{GPC} = 0.57 \times \text{GN}_{\text{m}}/\text{GY}$, % of grain dry mass) were calculated. In-season measurements included leaf area index (LAI , $\text{m}^2 \text{m}^{-2}$; 15 measurements in total), total aboveground biomass (AGBM , t DM ha^{-1} ; 28 measurements), total aboveground N (AGN , kg N ha^{-1} ; 27 measurements), and soil water content to maximum rooting depth (mm, 28 measurements). Plant-available soil water to maximum rooting depth (PASW , mm) was calculated from the measured soil water content by layer ($\Theta_{\text{v},i}$, vol%), the estimated lower limit of water extraction (LL , vol%), and the thickness of the soil layers (d, m):

$$\text{PASW} = \sum_{i=1}^k d_i \times (\Theta_{\text{v},i} - \text{LL}_i) \quad (1)$$

where k is the number of sampled soil layers.

Based on the critical N dilution curve of wheat (Justes *et al.*, 1994), a N nutrition index (NNI, dimensionless) was calculated to quantify crop N status. Although this curve is empirical, it is based on solid theoretical grounds (Lemaire & Gastal,

Table 1 Details of the experimental sites and experiments provided to the modelers

	Site			
	NL	AR	IN	AU
Site description				
Environment	High-yielding long-season	High/medium-yielding medium-season	Irrigated short-season	Low-yielding rain-fed short-season
Regional representation	Western and northern Europe	Argentina, northern China, western USA	India, Pakistan, southern China	Australia, southern Europe, northern Africa, South Africa, Middle East
Location name	Wageningen ('The Bouwing') The Netherlands	Balcarce Argentina	New Delhi India	Wongan Hills Australia
Coordinates	51°58' N, 05° 37' E	37° 45' S, 58° 18' W	28° 22' N, 77° 7' E	30° 53' S, 116° 43' E
Soil characteristics				
Soil type ^a	Silty clay loam	Clay loam	Sandy loam	Loamy sand
Rooting depth (cm)	200	130	160	210
Apparent bulk density (m ³ m ⁻³)	1.35	1.1	1.55	1.41
Top soil organic matter (%)	2.52	2.55	0.37	0.51
pH	6.0	6.3	8.3	5.7
Maximum plant available soil water (mm to maximum rooting depth)	354	222	109	125
Crop management				
Sowing density (seed m ⁻²)	228	239	250	157
Cultivar				
Name	Arminda	Oassis	HD2009	Gamenya
Vernalization requirement	High	Little	None	Little
Daylength response	High	Moderate	None	Moderate
Ploughed crop residue	Potato (4 t ha ⁻¹)	Maize (7 t ha ⁻¹)	Maize (1.5 t ha ⁻¹)	Wheat/weeds (1.5 t ha ⁻¹)
Irrigation (mm)	0	0	383	0
N application (kg N ha ⁻¹)	120 (ZC30 ^b)/40 (ZC65)	120 (ZC00)	60 (ZC00)/60 (ZC25)	50 (ZC10)
Initial top soil mineral N (kg N ha ⁻¹)	80	13	25	5
Sowing date	21 Oct. 1982	10 Aug. 1992	23 Nov. 1984	12 Jun. 1984
Anthesis date	20 Jun. 1983	23 Nov. 1992	18 Feb. 1985	1 Oct. 1984
Physiological maturity date	1 Aug. 1983	28 Dec. 1992	3 Apr. 1985	16 Nov. 1984
Growing season weather summary				
Cumulative rainfall (mm)	595	336	0	164
Cumulative global radiation (MJ m ⁻²)	2456	2314	2158	2632
Average daily mean temperature (°C)	8.8	13.8	17.5	14.1

^aSaturated soil water content, drainage upper limit and lower limit to water extraction were provided for 10 to 30-cm thick soil layers down to the maximum rooting depth.

^bZC, Zadoks stage (Zadoks *et al.*, 1974) at application is indicated in parenthesis (ZC00, sowing; ZC10, first leaf through coleoptile; ZC25, main shoot and five tillers; ZC30, pseudo stem erection; ZC65, anthesis half-way).

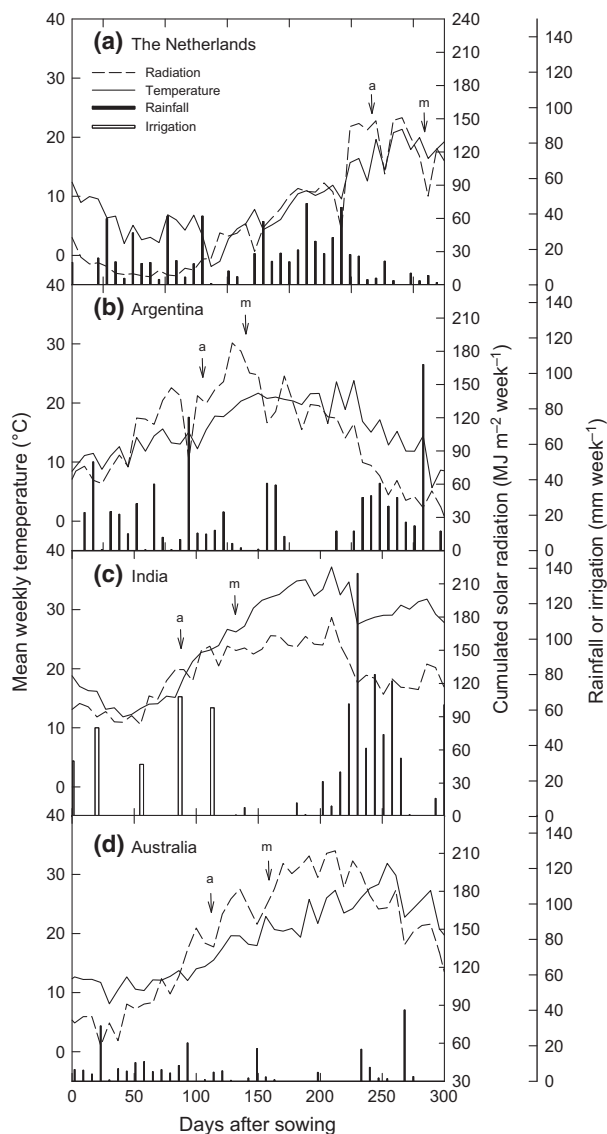


Fig. 1 Weather data at the four studied sites. Mean weekly temperature (solid lines), cumulative weekly solar radiation (dashed lines), cumulative weekly rainfall (vertical solid bars) and irrigation (vertical open bars) in (a) Wageningen, The Netherlands, (b) Balcarce, Argentina, (c) New Delhi, India, and (d) Wongan Hills, Australia. Vertical arrows indicate (a) anthesis and (m) physiological maturity dates.

1997). Climatic conditions can affect growth and N uptake differently, but the NNI reflects these effects in terms of crop N needs (Lemaire *et al.*, 2008; Gonzalez-Dugo *et al.*, 2010). For a given AGBM, NNI was calculated as the ratio between the actual and critical (N_C ; $\text{g N g}^{-1} \text{ DM}$) AGN concentrations defined by the critical N dilution curve (Justes *et al.*, 1994):

$$N_C = 5.35 \times \text{AGBM}^{-0.442} \quad (2)$$

If the NNI value is close to 1 it indicates an optimal crop N status, a value lower than 1 indicates N deficiency and a value higher than 1 indicates N excess.

Models and setup of model intercomparison

The models considered here were the 27 wheat crop models (Table S1) used in the AgMIP Wheat Pilot study (Asseng *et al.*, 2013). All of these models have been described in publications and are currently in use. Not all models simulated all measured variables, either because the models did not simulate them or because they were not in the standard outputs. Of the 27 models, 23 models simulated PASW values, and 20 simulated AGN and GN, and therefore NNI and GPC could be calculated for these 20 models. NHI could be calculated for 19 models.

All modeling groups were provided with daily weather data (i.e. precipitation, minimum and maximum air temperature, mean relative air humidity, dew point temperature, mean air vapor pressure, global radiation, and mean wind speed), basic physical characteristics of soil, initial soil water and N content by layer and crop management information (Table S1). No indication of how to interpret or convert this information into parameter values was given to the modelers. Modelers were provided with observed anthesis and maturity dates for the cultivars grown at each site. Qualitative information on vernalization requirements and daylength responses were also provided. All models were calibrated for phenology to avoid any confounding effects.

In the simulations, phenology parameters were adjusted to reproduce the observed anthesis and maturity dates, but otherwise models were not specifically adjusted to the growth data, which were only revealed to the modelers at the end of the simulation phase of the project. The information provided correspond to the partial model calibration in Asseng *et al.* (2013). Modelers were instructed to keep all parameters except for genotypic coefficients, constant across all four sites. The soil characteristics and initial conditions and crop management were specific to each site but were the same across all models.

The experimental data used in this study have not been used to develop or calibrate any of the 27 models. Experiments at AU and NL were used by one and two models as part of large datasets for model testing in earlier studies, respectively; but no calibration of the models was done. Except for the four Expert-N models which were run by the same group, all models were run by different groups without communication between the groups regarding the parameterization of the initial conditions or cultivar specific parameters. In most cases, the model developers ran their own model.

Model evaluation

Many different measures of the discrepancies between simulations and measurements have been proposed (Bellocchi *et al.*, 2010; Wallach *et al.*, 2013), and each captures somewhat different aspects of model behavior. We concentrated on the root mean squared error (RMSE) and the root mean squared relative error (RMSRE), where each error is expressed as a percentage of the observed value. The RMSE has the advantage of expressing error in the same units as the variable. For

comparing very different environments likely to give a broad range of crop responses, the relative error may be more meaningful than the absolute error as it gives more equal weight to each measurement. However, RMSRE needs to be interpreted with care because it is very sensitive to errors when measured values are small, as occurred for several early-season growth measurements.

RMSE was calculated as the square root of the mean squared error (MSE). MSE for model m and for a particular variable (MSE _{m}) was calculated as:

$$\text{MSE}_m = \frac{1}{N} \sum_{i=1}^N (y_i - \hat{y}_{m,i})^2 \quad (3)$$

where y_i is the value of the i th measurement of this variable, $\hat{y}_{m,i}$ is the corresponding value simulated by model m , and N is the total number of measurements of this variable (i.e. the sum over sites and over sampling dates per site for in-season variables).

RMSRE was calculated as:

$$\text{RMSRE}_m = 100 \times \sqrt{\frac{1}{N} \sum_{i=1}^N \left(\frac{y_i - \hat{y}_{m,i}}{y_i} \right)^2} \quad (4)$$

To assess whether a model that simulates well for one variable also performs well for other variables, Pearson's product-moment correlation between the RMSE or RMSRE value of each model was calculated across the variables. The adjusted two-sided P -values (q -values) resulting from the correction for multiple tests were calculated and reported here.

Multimodel ensemble estimators

We considered two estimators that are based on the ensemble of model simulations. The first ensemble estimator, e-mean, is the mean of the model simulations. The second ensemble estimator, e-median, is the median of the individual model simulations. For each of these ensemble models, e-mean and e-median, we calculated the same criteria as for the individual models, namely MSE, RMSE, and RMSRE.

To explore how e-mean MSE and e-median MSE varied with the number of models in the ensemble, we performed a bootstrap calculation for each value of M' (number of models in the ensemble) from 1 to 27. For each ensemble size M' we drew $B = 25 \times 2^n$ bootstrap samples of M' models with replacement, so the same model might be represented more than once in the sample. n was varied from 1 to 10 and the results were essentially unchanged beyond 3200 (i.e. for $n \geq 7$) bootstrap samples. The results reported here use $n = 9$. The final estimate of MSE for e-mean was then:

$$\text{MSE}_{\text{e-mean}} = \frac{1}{B} \frac{1}{N} \sum_{b=1}^B \sum_{i=1}^N (y_i - \hat{y}_{\text{e-mean},i}^b)^2 \quad (5)$$

where $\hat{y}_{\text{e-mean},i}^b$ is the e-mean estimate in bootstrap sample b of the i th measurements of this variable, given by:

$$\hat{y}_{\text{e-mean},i}^b = \frac{1}{M'} \sum_{m=1}^{M'} \hat{y}_{m,i}^b \quad (6)$$

For e-median the estimate of MSE was calculated as:

$$\text{MSE}_{\text{e-median}} = \frac{1}{B} \frac{1}{N} \sum_{b=1}^B \sum_{i=1}^N (y_i - \hat{y}_{\text{e-median},i}^b)^2 \quad (7)$$

In the case of e-mean, we can calculate the theoretical expectation of MSE analytically as a function of M' . Consider a variable at a particular site. Let μ_i^* represent the true expectation of model simulations for that site (the mean over all possible models), and let $\hat{\mu}_{i,M'}$ represent an e-mean simulation which is based on a sample of models of size M' . The expectation of MSE (expectation over possible samples of M' models) for e-mean is then:

$$\begin{aligned} E(\text{MSE}_{M'}) &= E \left[\frac{1}{N} \sum_{i=1}^N (y_i - \hat{\mu}_{i,M'})^2 \right] \\ &= \frac{1}{N} \sum_{i=1}^N E \left[(y_i - \mu_i^* + \mu_i^* - \hat{\mu}_{i,M'})^2 \right] \\ &= \frac{1}{N} \sum_{i=1}^N \left[(y_i - \mu_i^*)^2 + \frac{\text{var}(\hat{y}_i)}{M'} \right] \end{aligned} \quad (8)$$

where $\text{var}(\hat{y}_i)$ is the variance of the simulated values for the different models. The first term in the sum in (Eqn 8) is the squared bias of e-mean, when e-mean is based on a very large number of models. The second term is the variance of the model simulations divided by M . μ_i^* can be estimated as the average of the simulations over all the models in our study, and $\text{var}(\hat{y}_i)$ can be estimated as the variance of those model simulations.

All calculations and graphs were made using the R statistical software R 3.0.1 (R Core Team, 2013). Pearson's product-moment correlation P -values were adjusted for false discovery rate using the 'LBE' package (Dalmasso *et al.*, 2005), and bootstrap sampling used the R function `sample`.

Results

Evaluation of a population of wheat crop models

In most cases, measured in-season LAI, PASW, AGBM, AGN, and NNI, and end-of-season GY and GPC values were within the range of model simulations (Fig. 2, 3). The main disagreement between measured and simulated values was for LAI at IN, where the median of simulated in-season PASW (Fig. 2g) and AGBM (Fig. 2k) were close to the measured values but most models underestimated LAI (Fig. 2c) and overestimated AGN (Fig. 2o) around anthesis.

Even though measured GY ranged from 2.50 to 7.45 t DM ha⁻¹ across the four sites, the ranges of simulated GY values were similar at the four sites with an average range between minimum and maximum simulations of 1.64 t DM ha⁻¹ (Fig. 3a). The range between minimum and maximum simulations for GPC was also comparable at the four sites, averaging 7.1 percentage points (Fig. 3b). Model errors for GPC were in most cases due to poor simulation of AGN remobilization to

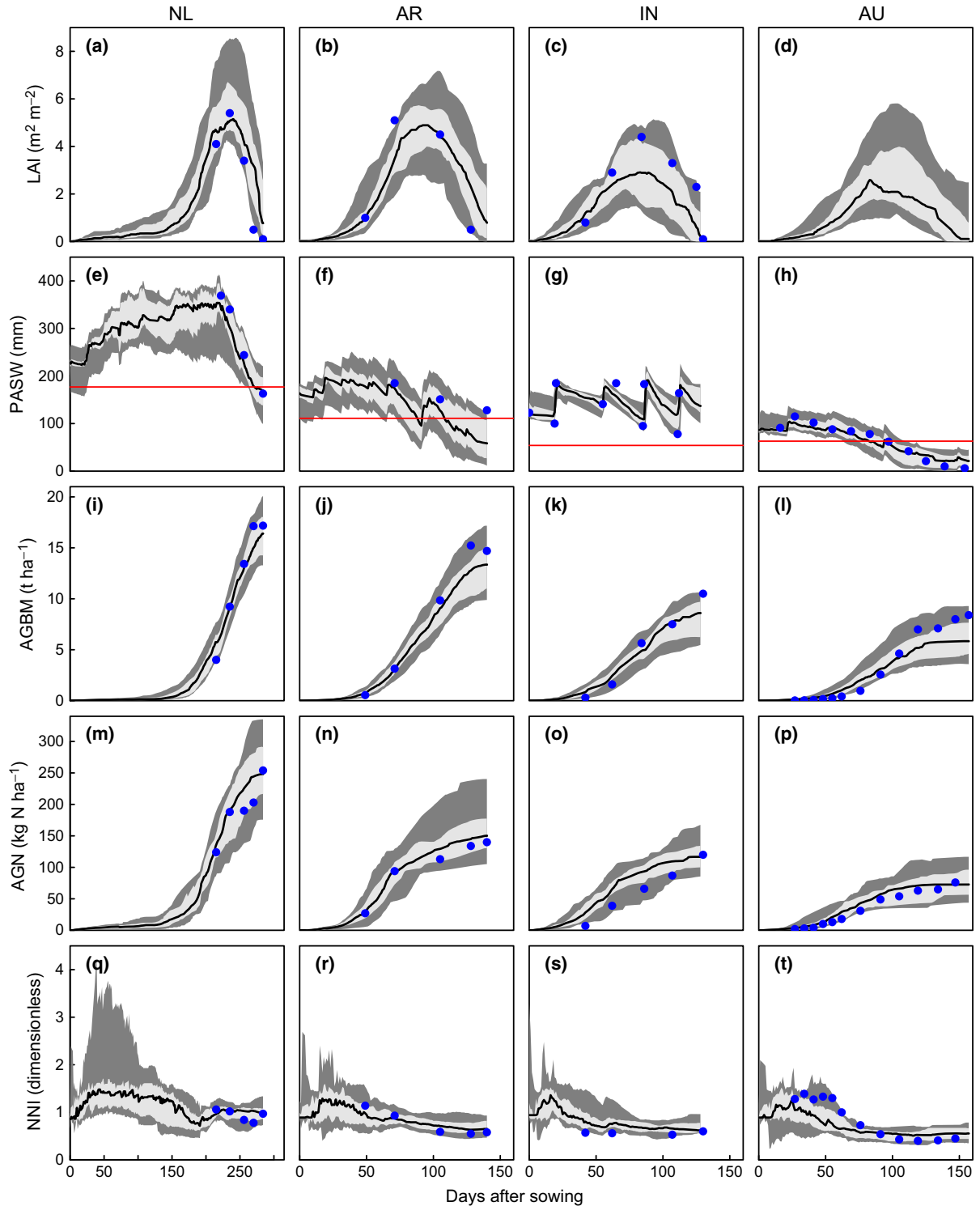


Fig. 2 Measured and simulated values of five in-season wheat crop variables for four sites. (a-d) Leaf area index (LAI), (e-h) plant-available soil water (PASW), (i-l) total aboveground biomass (AGBM), (m-p) total aboveground nitrogen (AGN), and (q-t) nitrogen nutrition index (NNI) vs. days after sowing in The Netherlands (NL), Argentina (AR), India (IN), and Australia (AU). Symbols are single measurements and solid lines are medians of the simulations (i.e. e-median). Dark gray areas indicate the 10th to 90th percentile range and light gray areas the 25th to 75th percentile range of the values generated by different wheat crop models. Twenty-seven models were used to simulate LAI and AGBM, 24 to simulate PASW, 20 to simulate AGN and NNI. In e-h the horizontal red lines indicate 50% soil water deficit.

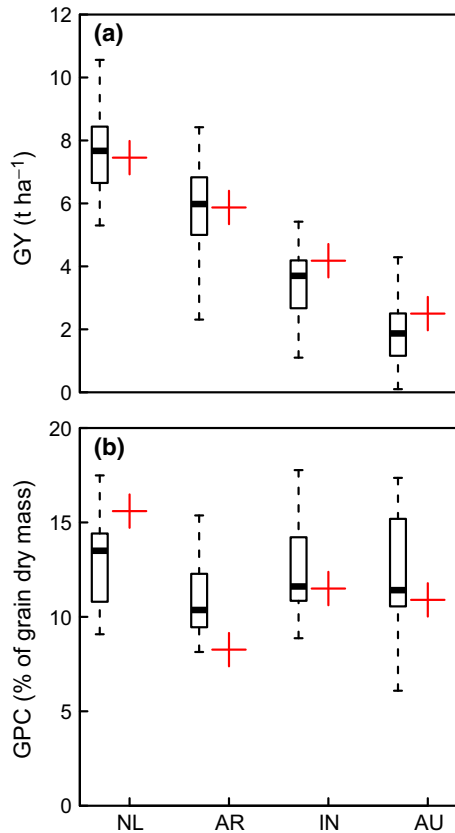


Fig. 3 Measured and simulated values of two major end-of-season wheat crop variables for four sites. Measured (red crosses) and simulated (box plots) values for end-of-season (a) grain yield (GY) and (b) grain protein concentration (GPC) are shown for The Netherlands (NL), Argentina (AR), India (IN), and Australia (AU). Simulations are from 27 different wheat crop models for GY and 20 for GPC. Boxes show the 25th to 75th percentile range, horizontal lines in boxes show medians, and error bars outside boxes show the 10th to 90th percentile range.

grains. Most models overestimated GPC at AR because they overestimated N remobilization to grains, while at NL most models underestimated GPC because they underestimated N remobilization.

The RMSRE averaged over all models was 29% (Fig. 4a and Table S2), and the RMSE average over all models was 1.25 t DM ha⁻¹ for GY (Fig. 4b and Table 2 and Table S3). The uncertainty in simulated GY was large, with RMSRE ranging from 8% to 73% among the 27 models, but 80% of the models had an RMSRE for GY comprised between 14% and 47% (Fig. 4a). For the other end-of-season variables RMSRE ranged from 7% to 60% for HI (averaging 24%), 22% to 61% for GN (averaging 38%), 15% to 52% for NHI (averaging 26%), and 8% to 122% for GPC (averaging 34%; Fig. 4a). For

the in-season variables with multiple measurements per site, the RMSRE ranged from 48% to 1496% for LAI, 37% to 355% for PASW, 41% to 542% for AGBM, 49% to 472% for AGN, and 16% to 104% for NNI (Fig. 4a). The large variability between models occurs because the models have different equations for many functions (as shown in Asseng *et al.* (2013) Table S2 in Supplemental) and different parameter values (Challinor *et al.*, 2014).

Of the three models with the smallest RMSE for GY, only the second-ranked model had RMSE values below the average of all models for all variables considered (Table 2). The other two models had an RMSE substantially higher than the average for at least one variable. The first- and second-ranked models simulated GY closely because of compensating errors. They underestimated LAI around anthesis and final AGBM which was compensated for by overestimating HI. For instance, the first-ranked model simulated that the canopy intercepted 83%, 74% and 51% of the incident radiation around anthesis in AR, IN and NL, respectively, while according to measured LAI values the percentage of radiation interception was close to 93% at the three sites (assuming an extinction coefficient of 0.55, an average value reported for wheat canopies (Sylvester-Bradley *et al.*, 2012)). This model compensated by having unrealistically high HI values that were 19% to 93% higher than measured HI. Theoretical maximum HI has been estimated at 62–64% for wheat (Foulkes *et al.*, 2011), while this model had simulated values up to 69% (in NL). The third-ranked model showed no significant compensation of errors. This model overestimated LAI around anthesis by 16% in AR and NL, but this translated into only a small effect on intercepted radiation, since the canopy intercepted more than 90% of incident radiation based on observed LAI.

Relation between the error for grain yield and that for underlying variables

There was little relation between the errors for different variables (Fig. 4a, b). There were some exceptions, however. Notably, RMSE for AGBM was highly correlated with that for GY, and that for AGN was correlated with GN (Fig. 5). Similarly, RMSE for AGN was highly correlated with that for LAI, PASW, and NNI. Finally, RMSE for NNI was correlated with that for PASW, HI, and GN and to a lesser extent with that for NNI. RMSE for GPC was not significantly correlated with any other variable. Overall, the correlations between RMSRE for different variables were similar to that between RMSE for different variables (Fig. S1).

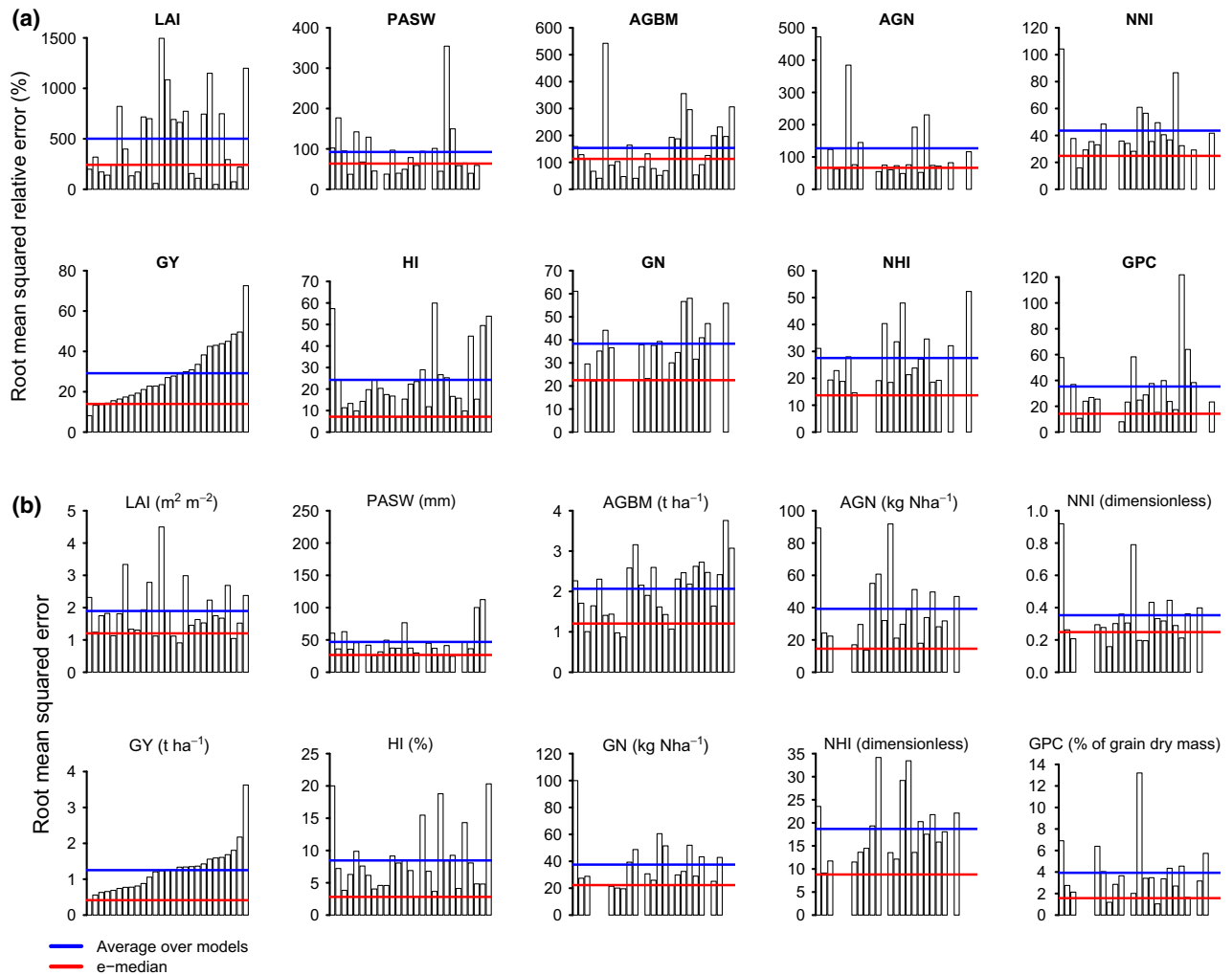


Fig. 4 Wheat crop model errors for in-season and end-of-season variables. (a) Root mean squared relative error (RMSRE) and (b) root mean squared error (RMSE) for in-season leaf area index (LAI), plant-available soil water (PASW), total aboveground biomass (AGBM), total aboveground nitrogen (AGN), nitrogen nutrition index (NNI), and for end-of-season grain yield (GY), biomass harvest index (HI), grain nitrogen yield (GN), nitrogen harvest index (NHI), and grain protein concentration (GPC). Twenty-seven models were used to simulate LAI, AGBM, GY, and HI, 20 to simulate AGN, GN, GPC, and NNI, 24 to simulate PASW, and 19 to simulate NHI. In (a) for GY the models are sorted from left to right in the order of increasing RMSE and this order of models was used to plot all other variables. The horizontal solid blue line shows RMSE or RMRSE averaged over all models and the horizontal red line shows RMSE or RMRSE for the median simulation of all models (e-median).

Multimodel ensemble estimators

Two multimodel ensemble estimators were tested. The first, the e-mean, uses the mean of the simulations of the ensemble members, a common practice in climate ensemble modeling (Knutti, 2010). The second, the e-median, uses the median of the simulations of the ensemble members. The e-median is expected to be less sensitive to outlier simulations than e-mean and therefore provide more robust estimates.

The e-median and e-mean values gave good agreement with measured values in almost all cases, despite

the fact that the simulations of the individual models varied considerably (Fig. 2, 3). For all responses, the RMSRE and RMSE of e-median and e-mean estimators were much lower than the RMSRE and RMSE averaged over all models (Fig. 4). For most variables, e-mean and e-median had similar RMSE and RMSRE values, and their ranking among all models was close (Table 2 and Table S2, S3). The largest difference in ranks was for RMSE for GPC, where e-median was ranked 3 and e-mean was ranked 7.

For most variables, e-mean and e-median were comparable to the best single model for that variable

Table 2 RMSE for in-season and end-of-season variables. Ensemble averages and e-mean and e-median values are based on 27 different models for LAI, AGBM, GY, and HI; 24 for PASW, 20 for AGN, GN, GPC, and NNI; and 19 for NHI. Values for the three best models for GY (based on RMSE) simulation are also given. Data for each individual model are given in Table S4. The numbers in parenthesis indicate the rank of the models (including e-mean and e-median) where 1 indicates the model with the lowest RMSE (i.e. best rank) for that variable. For each variable the model with the lowest RMSE is in bold type.

Estimator	RMSE for in-season variables					RMSE for end-of-season variables				
	LAI (m ⁻² m ⁻²)	PASW (mm)	AGBM (t DM ha ⁻¹)	AGN (kg N ha ⁻¹)	NNI ([◦])	GY (t DM ha ⁻¹)	HI (%)	GN (kg N ha ⁻¹)	NHI (%)	GPC (% of grain DM)
Average over all models	1.90	47	2.07	39	0.35	1.25	8.5	38	18.7	3.93
Model ranked 1 for GY	2.31 (23)	60 (21)	2.26 (17)	89 (21)	0.92 (22)	0.42 (2)	20.0 (28)	100 (22)	23.6 (18)	6.91 (21)
Model ranked 2 for GY	1.24 (7)	36 (9)	1.71 (13)	24 (8)	0.26 (8)	0.56 (4)	7.2 (16)	27 (9)	9.1 (2)	2.75 (9)
Model ranked 3 for GY	1.75 (16)	63 (22)	1.01 (3)	22 (7)	0.21 (4)	0.63 (5)	3.8 (5)	29 (10)	11.7 (5)	2.13 (6)
e-median	1.20 (6)	27 (3)	1.20 (6)	15 (3)	0.25 (7)	0.41 (1)	2.8 (2)	22 (5)	8.8 (1)	1.57 (3)
e-mean	1.29 (8)	27 (5)	1.19 (5)	13 (1)	0.24 (6)	0.49 (3)	2.2 (1)	23 (6)	9.8 (3)	2.32 (7)

(Fig. 4a, b). When e-median was ranked with the other models based on RMSRE, it ranked fourth for GY and third for GPC (Table S2); and first for GY and third for GPC when ranked based on RMSE (Table S3). One way to quantify the overall skill of e-mean and e-median is to consider the sum of ranks over all the variables. The sum of ranks based on RMSE for the 10 variables analyzed in this study was 37 for e-median and 45 for e-mean, while the lowest sum of ranks for an individual model (among the 17 models that simulated all variables) was 53 (Table S2). If we only considered the four variables simulated by all 27 models (i.e. LAI, AGBM, GY, and HI), the sum of ranks for e-median and e-mean was 15 and 17, respectively, while the best sum of ranks for an individual model with these four variables was 28.

To analyze the relationship between the number of models in an ensemble and the RMSE of both e-mean and e-median, we used a bootstrap approach to create a large number of ensembles for different multi-model ensemble sizes M' . For each M' , the RMSE of both e-mean and e-median in each bootstrap ensemble was calculated and averaged over bootstrap samples (Fig. 6). The standard deviation of RMSE for each M' shows how RMSE varies depending on the models that are included in the sample. The bootstrap average for e-mean followed very closely the theoretical expectation of RMSE (Fig. 6). The average RMSE of e-median also decreased with the number of models, in a manner similar to, but not identical to, the average e-mean RMSE. The differences were most pronounced for GPC (Fig. 6j).

Discussion

Working with multimodel ensembles is well-established in climate modeling, but only recently has the necessary international coordination been developed to make this also possible for crop models (Rosenzweig *et al.*, 2013). Here, we examined the performance of an ensemble of 27 wheat models, created in the context of the AgMIP Wheat Pilot study (Asseng *et al.*, 2013). Multiple crop responses, including both end-of-season and in-season growth variables were considered. Among these, GY and GPC are the main determinants of wheat productivity and end-use value. The other variables helped indicate whether models are realistic and consistent in their description of the processes leading to GY and GPC. This provides more comprehensive information on crop system properties beyond GY and is essential for the analysis of adaptation and mitigation strategies to global changes (Challinor *et al.*, 2014).

In only a few cases there were significant correlations between a model's error for one variable and

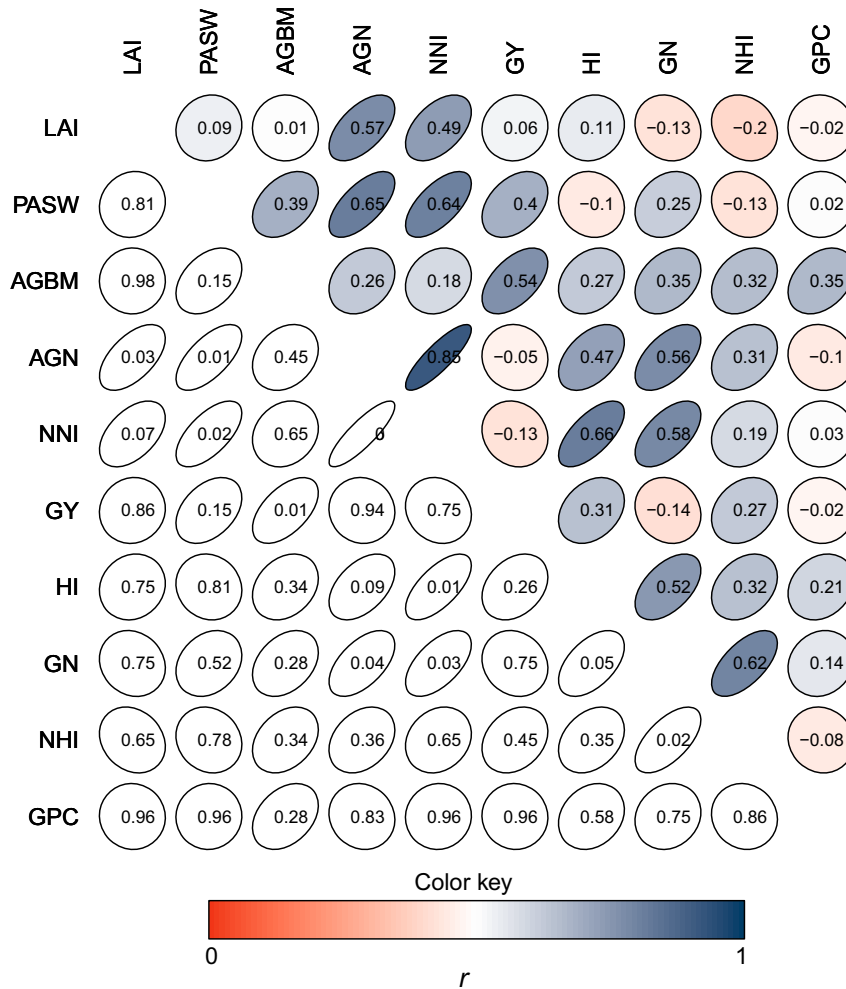


Fig. 5 Correlation matrix for Pearson's product-moment correlation (r) between the root mean squared error of simulated variables. In-season variables: leaf area index (LAI), plant-available soil water (PASW), total aboveground biomass (AGBM), total aboveground nitrogen (AGN), nitrogen nutrition index (NNI). End-of-season variables: grain yield (GY), biomass harvest index (HI), grain nitrogen yield (GN), nitrogen harvest index (NHI), and grain protein concentration (GPC). Twenty-seven models were used to simulate LAI, AGBM, GY, and HI, 20 to simulate AGN, GN, GPC, and NNI, 24 to simulate PASW, and 19 to simulate NHI. The numbers above the diagonal gap are r values and the numbers below are one-sided q -values (adjusted P -values for false discovery rate). The color (for r values only) and the shape of the ellipses indicate the strength of the correlation (the narrower the ellipse the higher the r value) and the direction of each ellipse indicates the sign of the correlation (a right-leaning ellipse indicates a positive correlation and a left-leaning ellipse indicates a negative correlation).

its error for other variables. Several individual models had relatively small errors for GY or GPC and large errors for in-season variables, including two of the three models with the lowest RMSE for GY. These models arrived at accurate simulations of GY or GPC without simulating crop growth accurately and thus got the right answer for, at least in part, the wrong reasons. That is, models can compensate for structural inconsistency. It has been argued that interactions among system components are largely empirical in most crop models (Ahuja & Ma, 2011) and that model error is minimized with different

parameter values for different variables (Wallach, 2011), which would explain why a model might simulate one variable well and not others. However, it remains unclear whether such compensation will be effective in a wide range of environments. The lack of correlation between model errors for different variables shows that one cannot simply evaluate models based on a single variable (response), since evaluation results can be quite different for other variables. It is important then to do crop model ensemble assessment for multiple variables (Challinor *et al.*, 2014), as done in this study.

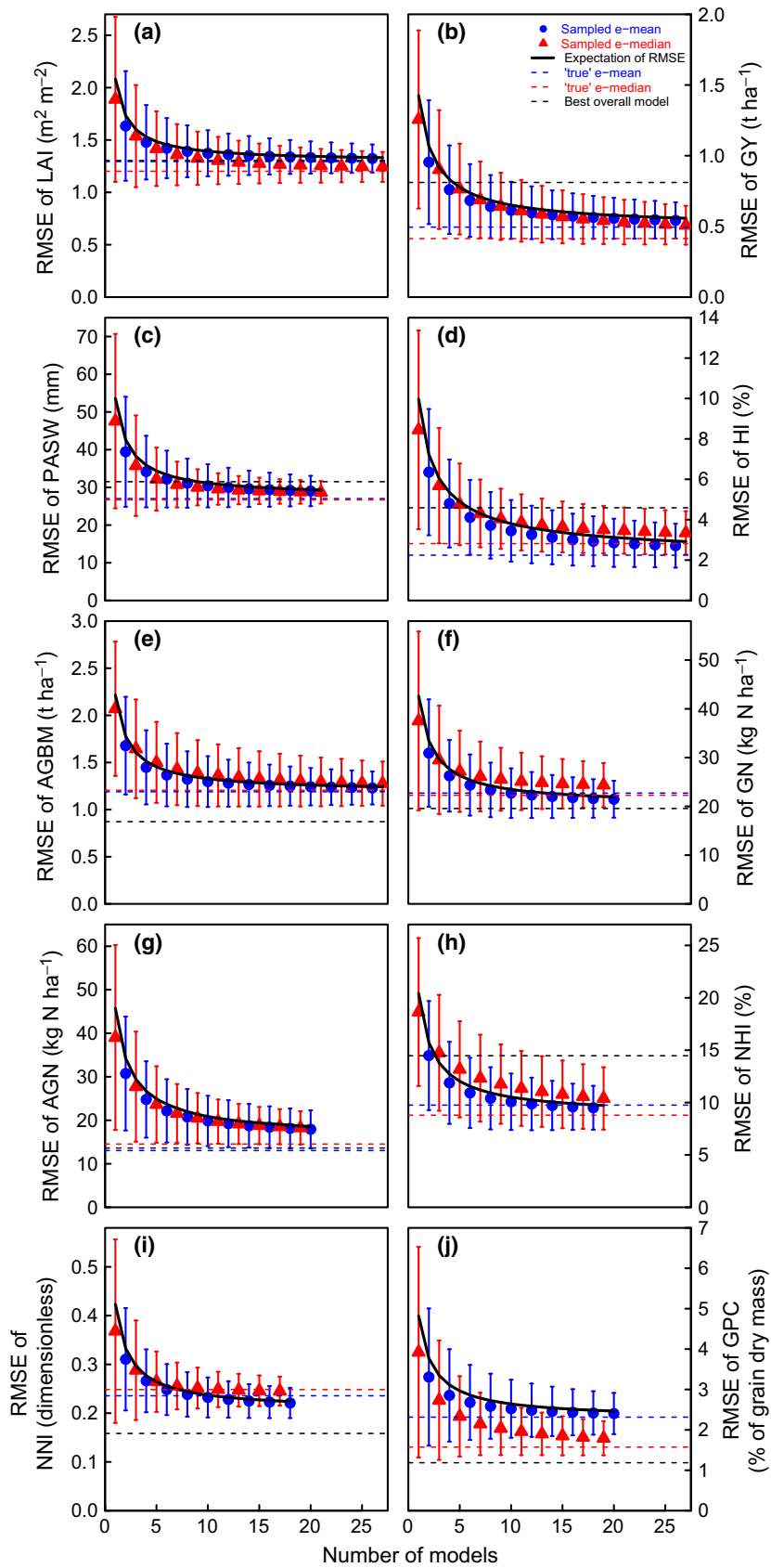


Fig. 6 How the number of models in an ensemble affects error estimates. Average root mean squared error (RMSE) (± 1 s.d.) of e-mean and e-median for in-season (a) leaf area index (LAI), (c) plant-available soil water (PASW), (e) total aboveground biomass (AGBM), (g) total aboveground nitrogen (AGN), and (i) nitrogen nutrition index (NNI) and for end-of-season (b) grain yield (GY), (d) biomass harvest index (HI), (f) grain nitrogen yield (GN), (h) nitrogen harvest index (NHI), and (j) grain protein concentration (GPC) vs. number of models in the ensemble. Values are calculated based on 12,800 bootstrap samples. The solid line is the analytical result for RMSE as a function of sample size (equation (8)). The blue dashed line shows the RMSE for e-mean and the red dashed line the RMSE for e-median of the multimodel ensemble. The black dashed line is the RMSE for the individual model with lowest sum of ranks for RMSE. For visual clarity the RMSE for e-mean is plotted for even numbers of models, and the RMSE for e-median for odd numbers of models.

Compensation of errors may be related to the way models are calibrated. If they are calibrated using only observed variable, e.g. GY, this may give parameter values that lead to unrealistic values of intermediate variables. The calibration insures that any errors in the intermediate variables compensate however, so that GY values are reasonably well-simulated. If final results are not used in calibration, for example if GPC is not used for calibration, then there may be compensation or compounding of the errors in the intermediate variables that lead to GPC.

There does not seem to be any simple relationship between model structure or the approach used to simulate individual processes and model error. Asseng *et al.* (2013) analyzed the response of the 27 crop models used in this study to a short heat shock around anthesis (seven consecutive days with a maximum daily temperature of 35 °C) and found that accounting for heat stress impact does not necessarily result in correctly simulating that effect. Similarly, we found that even closely related models did not necessarily cluster together and no single process could account for model error (data not shown). Therefore, it seems that model performances are not simply related to how a single process is modeled, but rather to the overall structure/parameterization of the model.

The behavior of the median and mean of the ensemble simulations was similar. Both estimators had much smaller errors and better skills than that averaged over models, for all variables. In comparing the sum of ranks of error for all variables, which provides an aggregated performance measure, the e-median was better than e-mean, but most importantly both were superior to even the best performing model in the ensemble. Different measures of performance might give slightly different results, but would not change the fact that e-median and e-mean compare well with even the best models.

E-mean and e-median had small errors in simulating not only end-of-season variables but also in-season variables. This suggests that multimodel ensembles could be useful not only for simulating GY and GPC, but also for relating those results to in-season growth processes. This is important if crop model ensembles are to be useful in exploring the consequences of global

change and the benefits of adaptation or mitigation strategies.

A fundamental question is the origin of the advantage of ensemble predictors over individual models. Two possible explanations relate to compensation among errors in processes descriptions and to more coverage of the possible crop and soil phase spaces. The first possible explanation is that certain models had large errors with compensations to achieve a reasonable yield simulation. In those cases, e-median can supply a better estimate when multiple responses are considered, since it gives reasonable results for all variables. In other cases, it is simply the fact that the errors in the different models tend to compensate each other well, that makes e-median the best estimator over multiple responses. The compensation of errors among models comes, at least in part, from the fact that models do not produce random outputs but are driven by environmental and management inputs and bio-physical processes and therefore they tend to converge to the measured crop response. It is an open question, however, as to whether the superiority of crop model ensemble estimators compared to individual models extends to conditions not tested in this study. Will this still be the case if the models are used to predict the impact of climate change? Or, will multimodel ensembles also be better capable than individual models to simulate the impact of interannual variability in weather at one site?

The second possible explanation relates to phase-space coverage. For climate models, the main reason for the superiority of multimodel ensemble estimators is that better coverage of the whole possible climate phase space leads to greater consistency (Hagedorn *et al.*, 2005). An analogous advantage holds as well for crop model ensembles, they have more associated knowledge and represent more processes than any individual model. Each of the individual models has been developed and calibrated based on a limited dataset. The ensemble simulators are in a sense averaging over these datasets, which gives them the advantage of a much broader database than any individual model and thus reduces the need for site- and varietal-specific model calibration.

The use of ensemble estimators to answer new questions in the future poses specific questions regarding the best procedure for creating an ensemble. Several of these questions have been debated in the climate science community (Knutti, 2010), but not always in a way that is directly applicable to crop models. One question is how performance varies with the number of models in the ensemble. Here we found that the change in ensemble error ($MSE_{M'}$) with the number of model in an ensemble (M') follows the expectation of MSE. Thus when planning ensemble studies, one can estimate the potential reduction in $MSE_{M'}$ and therefore, do a costs vs. benefits analysis for increasing M' . In the ensemble studied here, for all the variables, MSE for an ensemble of 10 models was close to the asymptotic limit for very large M' .

Other questions include how to choose the models in the ensemble, and whether one should weight the models in the ensemble differently, based on past performance and convergence for new situations (Tebaldi & Knutti, 2007). In this respect, the crop modeling community might employ some of the ensemble weighting methods developed by the climate modeling community (Christensen *et al.*, 2010). There are also questions about the possible multiple uses of models. Would it be advantageous to have multiple simulations, based on a diversity of initial conditions (including 'spin-up' periods for models that depend on simulation of changes in soil organic matter) or multiple parameter sets from each model? In any case, the first step is to document the accuracy of multimodel ensemble estimators in specific situations, as done here.

In summary, by reducing simulation error and improving the consistency of simulation results for multiple variables, crop model ensembles could substantially increase the range of questions that could be addressed. A lack of correlation between end-of-season and in-season errors in the individual models indicates that further work is needed to improve the representation of the dynamics of growth and development processes leading to GY in crop models. This is crucial for their application under changed climatic or management conditions.

Most of the physical and physiological processes that are simulated in wheat models are the same as for other crops. In fact, several of the models in this study have a generic structure so that they can be applied to various crops, and for some of them the differences between crops are simply in the parameter values. It is thus reasonable to expect that the results obtained here for wheat are broadly applicable to other crop species. It would be worthwhile to study whether these results also apply more generally to biological and ecological system models.

Acknowledgements

P.M. is grateful to the INRA metaprogram 'Adaptation of Agriculture and Forests to Climate Change' and Environment and Agronomy Division for supporting several stays at the University of Florida during this work.

References

- Ahuja LR, Ma L (2011) A synthesis of current parameterization approaches and needs for further improvements. In: *Methods of Introducing System Models into Agricultural Research* (eds Ahuja LR, Ma L), pp. 427–440. American Society of Agronomy, Crop Science Society of America, Soil Science Society of America, Madison, WI.
- Angulo C, Rötter R, Lock R, Enders A, Fronzek S, Ewert F (2013) Implication of crop model calibration strategies for assessing regional impacts of climate change in Europe. *Agricultural and Forest Meteorology*, **170**, 32–46.
- Asseng S, Keating BA, Fillery IRP *et al.* (1998) Performance of the APSIM-wheat model in Western Australia. *Field Crops Research*, **57**, 163–179.
- Asseng S, Ewert F, Rosenzweig C *et al.* (2013) Uncertainty in simulating wheat yields under climate change. *Nature Climate Change*, **3**, 827–832.
- Bassu S, Brisson N, Durand J-L *et al.* (2014) How do various maize crop models vary in their responses to climate change factors? *Global Change Biology*, **20**, 2301–2320.
- Bellocchi G, Rivington M, Donatelli M, Matthews K (2010) Validation of biophysical models: issues and methodologies. A review. *Agronomy for Sustainable Development*, **30**, 109–130.
- Bertin N, Martre P, Genard M, Quilot B, Salon C (2010) Under what circumstances can process-based simulation models link genotype to phenotype for complex traits? Case-study of fruit and grain quality traits. *Journal of Experimental Botany*, **61**, 955–967.
- Bloom DE (2011) 7 Billion and Counting. *Science*, **333**, 562–569.
- Bosilovich MG, Robertson FR, Chen JY (2011) Global energy and water budgets in MERRA. *Journal of Climate*, **24**, 5721–5739.
- Challinor AJ, Wheeler T, Hemming D, Upadhyaya HD (2009) Ensemble yield simulations: crop and climate uncertainties, sensitivity to temperature and genotypic adaptation to climate change. *Climate Research*, **38**, 117–127.
- Challinor A, Martre P, Asseng S, Thornton P, Ewert F (2014) Making the most of climate impacts ensembles. *Nature Climate Change*, **4**, 77–80.
- Christensen JH, Kjellström E, Giorgi F, Lenderink G, Rummukainen M (2010) Weight assignment in regional climate models. *Climate Research*, **44**, 179–194.
- Dalmasso C, Broët P, Moreau T (2005) A simple procedure for estimating the false discovery rate. *Bioinformatics*, **21**, 660–668.
- Easterling WE, Aggarwal PK, Batima P *et al.* (2007) Food, fibre and forest products. In: *Climate Change 2007: Impacts, Adaptation and Vulnerability. Contribution of Working Group II to the Fourth Assessment Report of the intergovernmental Panel on Climate Change* (eds Parry ML, Canziani OF, Palutikof JP, Van De Linden P, Hanson CE), pp. 273–313. Cambridge University Press, Cambridge, UK.
- FAOSTAT (2014) Food and Agricultural organization of the United Nations (FAO). FAO Statistical Databases. Available at: faostat3.fao.org (accessed 28 October 2014).
- Foulkes MJ, Slafer GA, Davies WJ *et al.* (2011) Raising yield potential of wheat. III. Optimizing partitioning to grain while maintaining lodging resistance. *Journal of Experimental Botany*, **62**, 469–486.
- Godfray HCJ, Beddington JR, Crute IR *et al.* (2010) Food security: the challenge of feeding 9 billion people. *Science*, **327**, 812–818.
- Gonzalez-Dugo V, Durand J-L, Gastal F (2010) Water deficit and nitrogen nutrition of crops. A review. *Agronomy for Sustainable Development*, **30**, 529–544.
- Grenouillet G, Buisson L, Casajus N, Lek S (2011) Ensemble modelling of species distribution: the effects of geographical and environmental ranges. *Ecography*, **34**, 9–17.
- Groot JJR, Verberne ELJ (1991) Response of wheat to nitrogen fertilization, a data set to validate simulation models for nitrogen dynamics in crop and soil. In: *Nitrogen Turnover in the Soil-Crop System. Modelling of Biological Transformations, Transport of Nitrogen and Nitrogen Use Efficiency. Proceedings of a Workshop* (eds Groot JJR, De Willigen P, Verberne ELJ), pp. 349–383. Institute for Soil Fertility Research, Haren, The Netherlands.
- Hagedorn T, Doblas-Reyes FJ, Palmer TN (2005) The rationale behind the success of multi-model ensembles in seasonal forecasting – I. Basic concept. *Tellus*, **57A**, 219–233.
- Justes E, Mary B, Meynard JM, Machet JM, Thelier-Huche L (1994) Determination of a critical nitrogen dilution curve for winter wheat crops. *Annals of Botany*, **74**, 397–407.

- Knutti R (2010) The end of model democracy? An editorial comment. *Climatic Change*, **102**, 395–404.
- Ko J, Ahuja L, Kimball B *et al.* (2010) Simulation of free air CO₂ enriched wheat growth and interactions with water, nitrogen, and temperature. *Agricultural and Forest Meteorology*, **150**, 1331–1346.
- Lemaire G, Gastal F (1997) N uptake and distribution in plant canopies. In: *Diagnosis of the Nitrogen Status in Crops* (ed. Lemaire G), pp. 3–43. Springer Verlag, Berlin, Germany.
- Lemaire G, Jeuffroy M-H, Gastal F (2008) Diagnosis tool for plant and crop N status in vegetative stage: theory and practices for crop N management. *European Journal of Agronomy*, **28**, 614–624.
- Lobell DB, Schlenker W, Costa-Roberts J (2011) Climate trends and global crop production since 1980. *Science*, **333**, 616–620.
- Naveen N (1986) Evaluation of soil water status, plant growth and canopy environment in relation to variable water supply to wheat. Unpublished PhD, IARI, New Delhi.
- Palosuo T, Kersebaum KC, Angulo C *et al.* (2011) Simulation of winter wheat yield and its variability in different climates of Europe: a comparison of eight crop growth models. *European Journal of Agronomy*, **35**, 103–114.
- Porter JR, Semenov MA (2005) Crop responses to climatic variation. *Philosophical Transactions of the Royal Society of London B Biological Sciences*, **360**, 2021–2035.
- R Core Team (2013) *R: A Language and Environment for Statistical Computing*. R Foundation for Statistical Computing, Vienna, Austria.
- Räisänen J, Palmer TN (2001) A probability and decision-model analysis of a multimodel ensemble of climate change simulations. *Journal of Climate*, **14**, 3212–3226.
- Rosenzweig C, Elliott J, Deryng D *et al.* (2014) Assessing agricultural risks of climate change in the 21st century in a global gridded crop model intercomparison. *Proceedings of the National Academy of Sciences*, **111**, 3268–3273.
- Rosenzweig C, Jones JW, Hatfield JL *et al.* (2013) The Agricultural Model Intercomparison and Improvement Project (AgMIP): protocols and pilot studies. *Agricultural and Forest Meteorology*, **170**, 166–182.
- Rötter RP, Carter TR, Olesen JE, Porter JR (2011) Crop-climate models need an overhaul. *Nature Climate Change*, **1**, 175–177.
- Rötter RP, Palosuo T, Kersebaum KC *et al.* (2012) Simulation of spring barley yield in different climatic zones of Northern and Central Europe: a comparison of nine crop models. *Field Crops Research*, **133**, 23–36.
- Stackhouse P (2006) Prediction of worldwide energy resources. Available at: <http://power.larc.nasa.gov> (accessed 28 October 2014).
- Sylvester-Bradley R, Riffkin P, O'leary G (2012) Designing resource-efficient ideotypes for new cropping conditions: wheat (*Triticum aestivum* L.) in the High Rainfall Zone of southern Australia. *Field Crops Research*, **125**, 69–82.
- Tebaldi C, Knutti R (2007) The use of the multi-model ensemble in probabilistic climate projections. *Philosophical Transactions of the Royal Society A: Mathematical, Physical and Engineering Sciences*, **365**, 2053–2075.
- Travasso MI, Magrin GO, Rodríguez R, Grondona MO (2005) Comparing CERES-wheat and SUCROS2 in the Argentinean Cereal Region. In: *MODSIM 2005 International Congress on Modelling and Simulation* (eds Zerger A, Argent RM), pp. 366–369. Modelling and Simulation Society of Australia and New Zealand. Available at: <http://www.mssanz.org.au/MODSIM95/Vol%201/Travasso.pdf>. (accessed 28 October 2014)
- Trewavas A (2006) A brief history of systems biology: “every object that biology studies is a system of systems”. Francois Jacob (1974). *Plant Cell*, **18**, 2420–2430.
- Tubiello FN, Soussana J-F, Howden SM (2007) Crop and pasture response to climate change. *Proceedings of the National Academy of Sciences*, **104**, 19686–19690.
- Wallach D (2011) Crop Model Calibration: a Statistical Perspective. *Agronomy Journal*, **103**, 1144–1151.
- Wallach D, Makowski D, Jones JW, Brun F (2013) *Working with Dynamic Crop Models. Methods Tools and Examples for Agriculture and Environment*. Academic Press, London.
- White JW, Hoogenboom G, Kimball BA, Wall GW (2011) Methodologies for simulating impacts of climate change on crop production. *Field Crops Research*, **124**, 357–368.
- Zadoks JC, Chang TT, Konzak CF (1974) A decimal code for the growth stages of cereals. *Weed Research*, **14**, 415–421.

Supporting Information

Table S1. Name, reference and source of the 27 wheat crop models used in this study.

Table S2. Root mean square relative error (RMSRE) for in-season and end-of-season variables.

Table S3. Root mean square error (RMSE) for in-season and end-of-season variables.

Figure S1. Correlation matrix for Pearson's product-moment correlation (r) between the root mean squared relative error of simulated variables.