# Supplementary Information to "Global gridded crop models underestimate yield losses from climatic extremes

Cornelia Auer[1,20], Kobe De Maeyer[2,20], Christoph Müller[1], Michael Hoehle[3], Jacob Schewe[1], Jonas Jaegermeyr[4,5,1], Juraj Balkovic[6], Thiago Berton Ferreira[7], Babacar Faye[8], Christian Folberth[6], Jose R. Guarin[9], Stefanie Heinicke[1], Gerrit Hoogenboom[7], Toshichika Iizumi[10], Atul K. Jain[11], Tzu-Shun Lin[12], Wenfeng Liu[13,14], Oleksandr Mialyk[15], Masashi Okada[16], Sam S. Rabin[12], Chenzhi Wang[17], Heidi Webber[17,18], Florian Zabel[19]

[1] *Potsdam Institute for Climate Impact Research, Potsdam, Germany*
[2] *Utrecht University, Netherlands*
[3] *University of Greifswald, Germany*
[4] *Columbia University, Climate School, New York, NY 10025, USA*
[5] *NASA Goddard Institute for Space Studies, New York, NY 10025, USA*
[6] *Biodiversity and Natural Resources Program, International Institute for Applied Systems Analysis (IIASA), Laxenburg, Austria*
[7] *University of Florida, USA*
[8] *Universite du Sine Saloum El Hadj Ibrahima NIASS (USSEIN), Kaolack, Senegal*
[9] *NASA Goddard Institute for Space Studies, New York, NY 10025, USA*
[10] *National Agriculture and Food Research Organization, Tsukuba, Japan*
[11] *Department of Climate, Meteorology and Atmospheric Sciences (CLiMAS), University of Illinois, Urbana-Champaign*
[12] *NSF National Center for Atmospheric Research, Boulder, CO, USA*
[13] *State Key Laboratory of Efficient Utilization of Agricultural Water Resources, China Agricultural University, Beijing 100083, China*
[14] *Center for Agricultural Water Research in China, College of Water Resources and Civil Engineering, China Agricultural University, Beijing 100083, China*
[15] *University of Twente, the Netherlands*
[16] *National Institute for Environmental Studies, Tsukuba, Japan*
[17] *Leibniz Centre for Agricultural Landscape Research (ZALF), Germany*
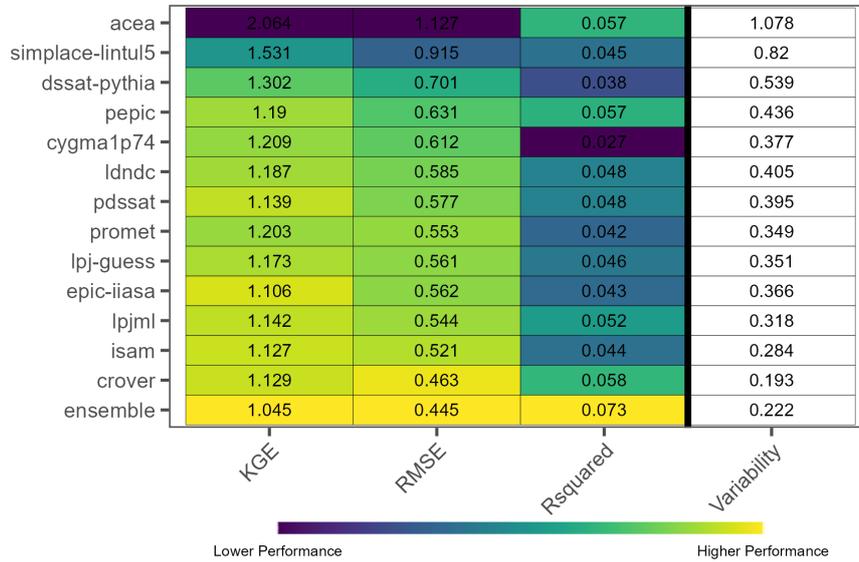[18] *BTU Cottbus, Cottbus, Germany*
[19] *Department of Environmental Sciences, University of Basel, Switzerland*
[20] *These authors contributed equally: Cornelia Auer, Kobe De Maeyer*
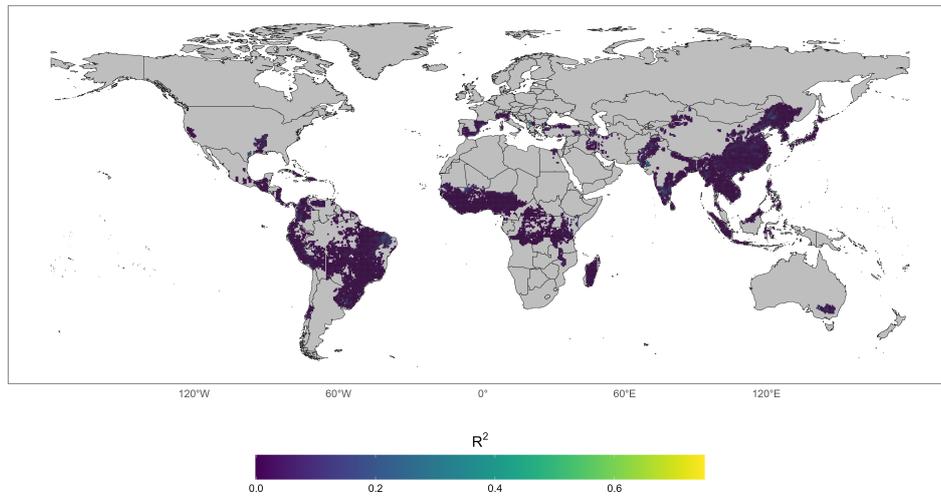
July 21, 2025

# 1 General model performance

As an alternative to the general model performance table shown in Figure 1 in the main text, where the metric values were averaged using a harvest area weighted mean, we provide here the same table using an unweighted mean (see Supplementary Figure 1). This approach does not account for the differences in the harvest area between grid cells and results in significantly lower $R^2$ values. This indicates that especially low producing areas (e.g. on the African continent but not only) are poorly represented and lead to lower overall performance.

| | KGE | RMSE | Rsquared | Variability |
|---|---|---|---|---|
| acea | 2.064 | 1.127 | 0.057 | 1.078 |
| simplace-lintul5 | 1.531 | 0.915 | 0.045 | 0.82 |
| dssat-pythia | 1.302 | 0.701 | 0.038 | 0.539 |
| pepic | 1.19 | 0.631 | 0.057 | 0.436 |
| cygma1p74 | 1.209 | 0.612 | 0.027 | 0.377 |
| ldndc | 1.187 | 0.585 | 0.048 | 0.405 |
| pdssat | 1.139 | 0.577 | 0.048 | 0.395 |
| promet | 1.203 | 0.553 | 0.042 | 0.349 |
| lpj-guess | 1.173 | 0.561 | 0.046 | 0.351 |
| epic-iiasa | 1.106 | 0.562 | 0.043 | 0.366 |
| lpjml | 1.142 | 0.544 | 0.052 | 0.318 |
| isam | 1.127 | 0.521 | 0.044 | 0.284 |
| crover | 1.129 | 0.463 | 0.058 | 0.193 |
| ensemble | 1.045 | 0.445 | 0.073 | 0.222 |

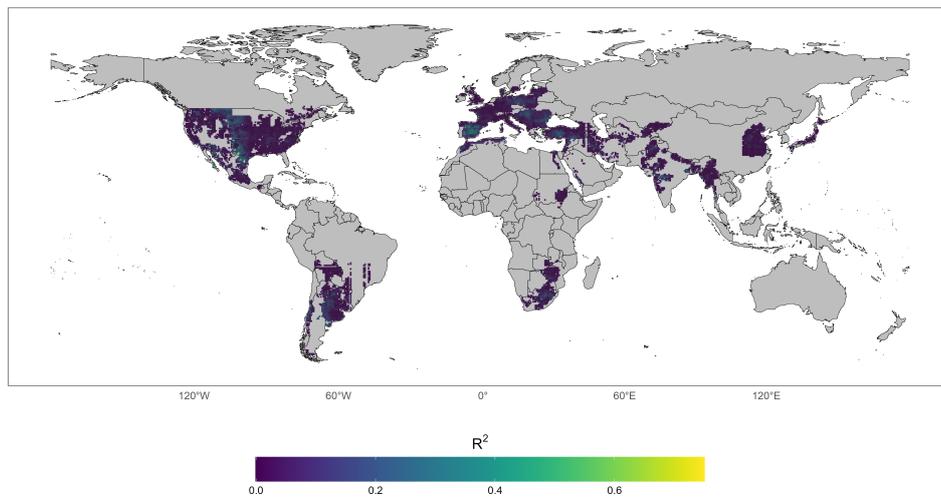Lower Performance        Higher Performance

Supplementary Figure 1: Global assessment of overall crop model performance by comparing the Iizumi et al. benchmark data GDHY (Iizumi and Sakai [2020]) against model results across the time period 1981-2016. Results are aggregated over all four crops (soy, maize, winter wheat, and first growing season rice) and based on the absolute detrending. Table comparing the individual models and model ensemble median for multiple metrics: $R^2$, Root Mean Square Error (RMSE), Kling-Gupta Efficiency (KGE), and variability (measured by standard deviation). The metrics are computed per grid cell and then averaged into a single value by taking an unweighted mean. Models with low variability tend to exhibit favorable KGE and RMSE values alongside moderate $R^2$ scores (relative to models with higher variability). Please note, the same table is provided in the main text, without being weighted by harvest area. This suggests that low-yielding regions, often coinciding with regions in income countries, e.g. on the African continent, are inadequately represented in the models, which contributes to lower overall performance metrics.
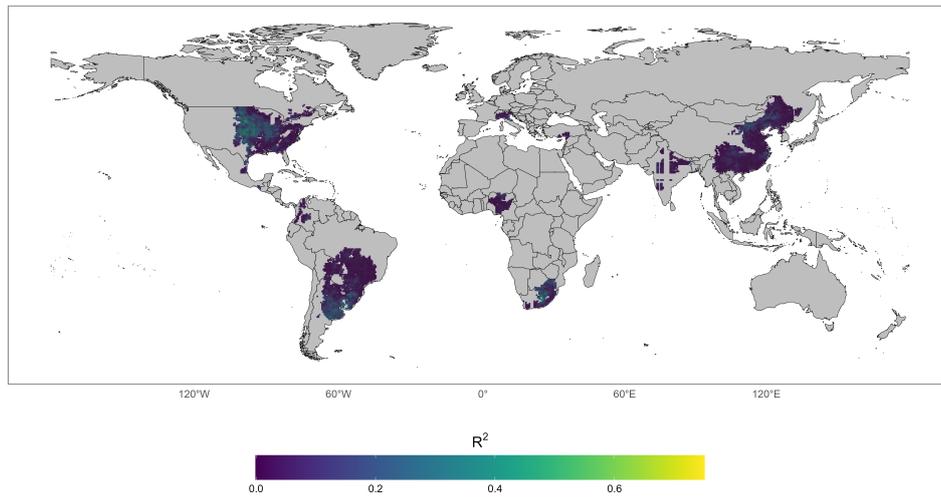
Additionally we provide here the global performance maps for each main crop type- soy, maize, winter wheat, and first growing season rice- separately instead of aggregated over crop types as presented in Figure 1 from the main analysis.
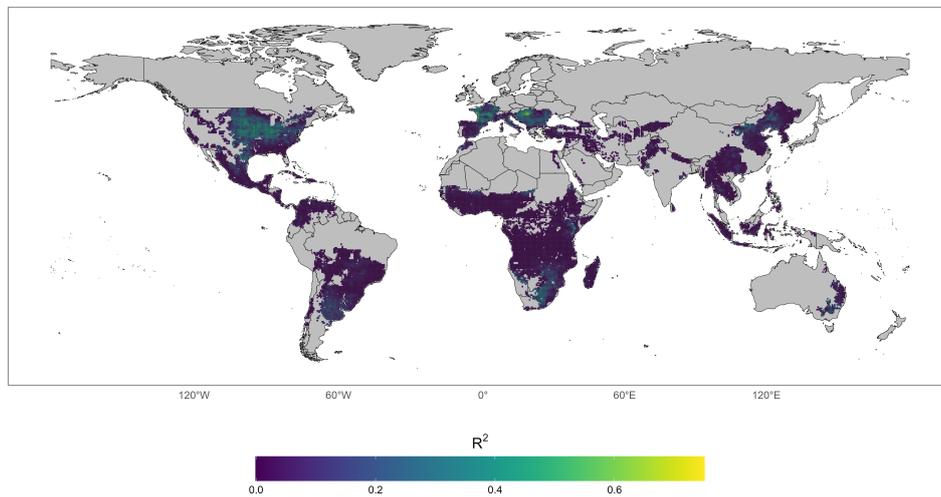


Supplementary Figure 2: RICE (first growing season) – Global assessment of general crop model performance (1981–2016), comparing Iizumi et al. benchmark data GDHY (Iizumi and Sakai [2020]), using absolute detrended yield values (see Eq. 6). Global $R^2$ map for the model ensemble median and benchmark data.



Supplementary Figure 3: WINTER WHEAT– Global assessment of geneal crop model performance (1981–2016), comparing Iizumi et al. benchmark data GDHY (Iizumi and Sakai [2020]), using absolute detrended yield values (see Eq. 6). Global $R^2$ map for the model ensemble median and benchmark data.

Supplementary Figure 4: SOY – Global assessment of general crop model performance (1981–2016), comparing Iizumi et al. benchmark data GDHY (Iizumi and Sakai [2020]), using absolute detrended yield values (see Eq. 6). Global $R^2$ map for the model ensemble median and benchmark data.



Supplementary Figure 5: MAIZE – Global assessment of general crop model performance (1981–2016), comparing Iizumi et al. benchmark data GDHY (Iizumi and Sakai [2020]), using absolute detrended yield values (see Eq. 6). Global $R^2$ map for the model ensemble median and benchmark data.
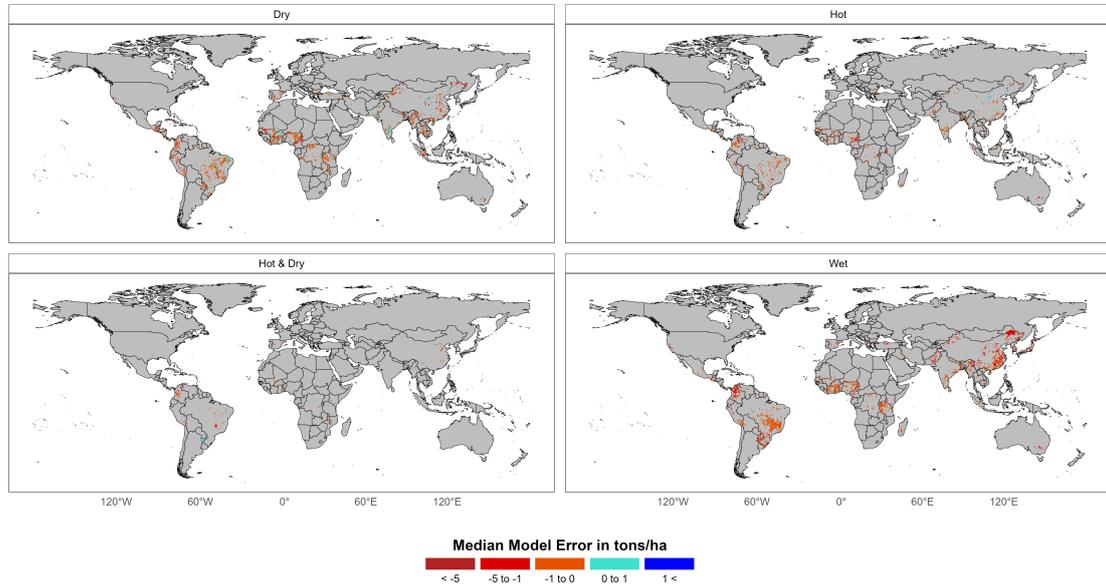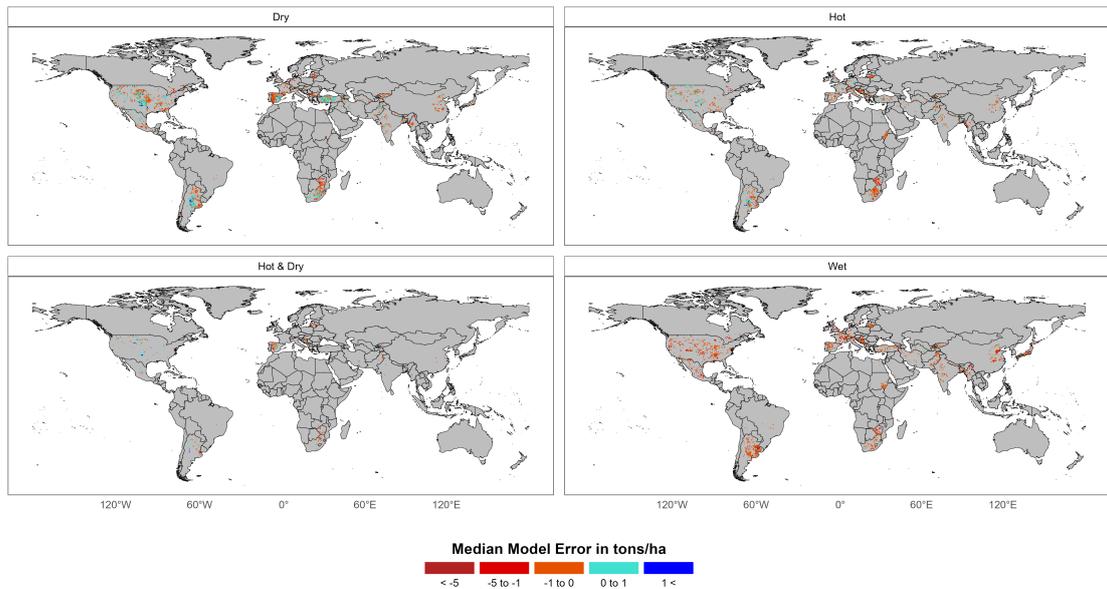
# 2 Climatic extremes spatial performance

Here we assess the spatial distribution of median errors of the model ensemble for different climatic extremes per crop instead of aggregated over crops as presented in Figure 2 from the main analysis.



Supplementary Figure 6: RICE (first growing season) – Spatial distribution of model ensemble median errors under different extreme weather conditions: dry, hot, compound hot & dry, and wet. Results are based on detrended yield anomalies (Eq. 6). Colors indicate the median absolute difference between the model ensemble median and benchmark data from GDHY (Iizumi and Sakai [2020]), with orange to red denoting underestimation and turquoise to blue indicating overestimation.

Supplementary Figure 7: WINTER WHEAT – Spatial distribution of model ensemble median errors under different extreme weather conditions: dry, hot, compound hot & dry, and wet. Results are based on detrended yield anomalies (Eq. 6). Colors indicate the median absolute difference between the model ensemble median and benchmark data from GDHY (Iizumi and Sakai [2020]), with orange to red denoting underestimation and turquoise to blue indicating overestimation.



Supplementary Figure 8: SOY – Spatial distribution of model ensemble median errors under different extreme weather conditions: dry, hot, compound hot & dry, and wet. Results are based on detrended yield anomalies (Eq. 6). Colors indicate the median absolute difference between the model ensemble median and benchmark data from GDHY (Iizumi and Sakai [2020]), with orange to red denoting underestimation and turquoise to blue indicating overestimation.
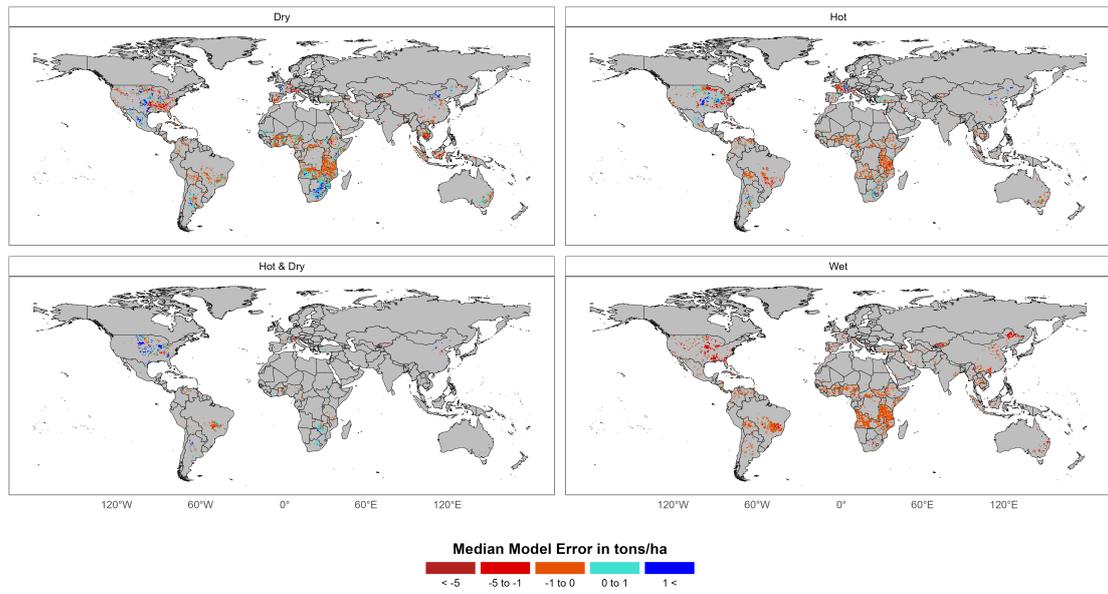
Supplementary Figure 9: MAIZE– Spatial distribution of model ensemble median errors under different extreme weather conditions: dry, hot, compound hot & dry, and wet. Results are based on detrended yield anomalies (Eq. 6). Colors indicate the median absolute difference between the model ensemble median and benchmark data from GDHY (Iizumi and Sakai [2020]), with orange to red denoting underestimation and turquoise to blue indicating overestimation.
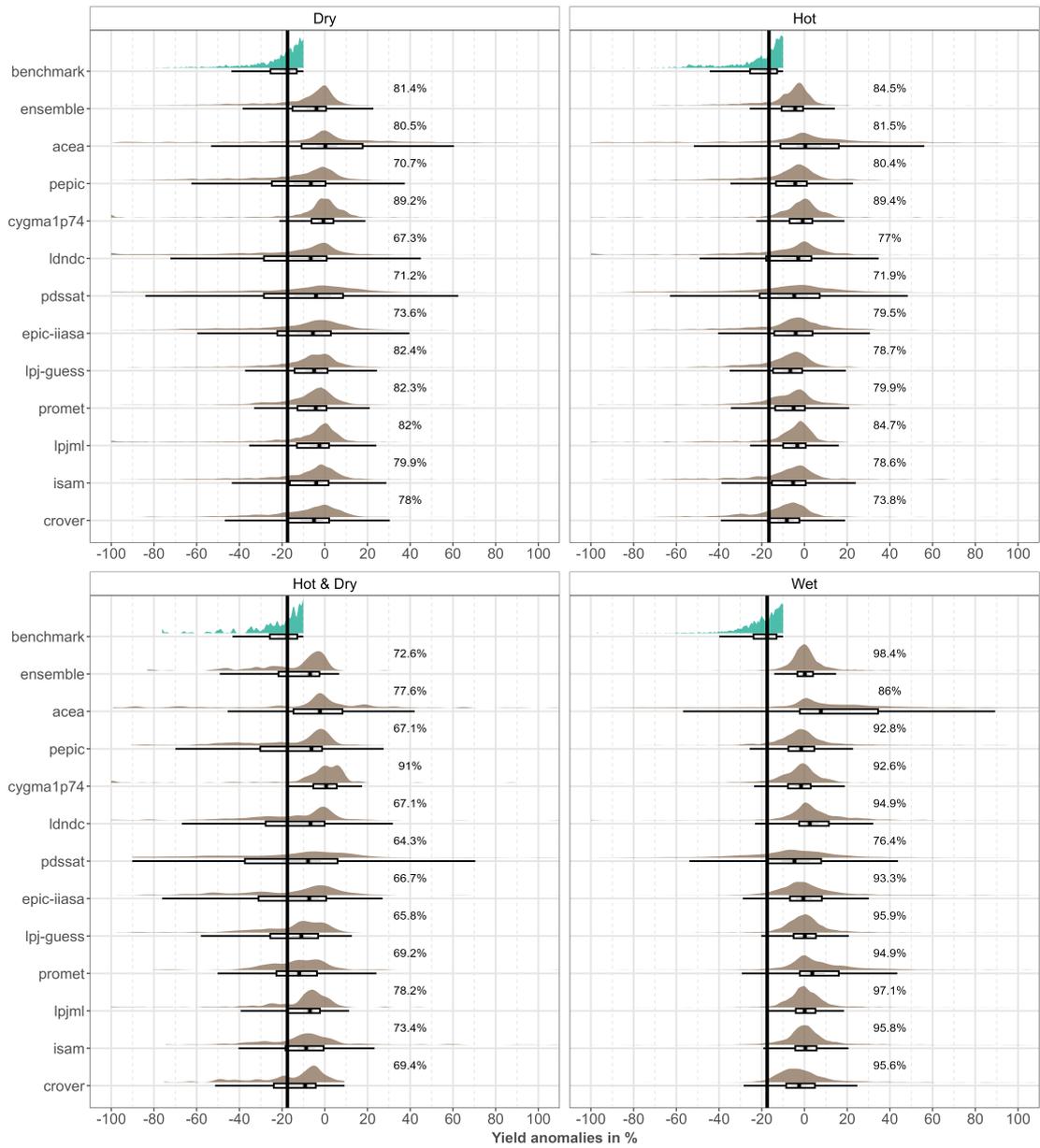
# 3 Raincloud plots: relative detrended yield anomalies

Four main crop types - soy, maize, winter wheat, and first growing season rice- are assessed individually below. The Supplementary Figures 10, 11, 12, 13, depict raincloud plots of the distribution of relative detrended yield anomalies (in %) of the Iizumi et al. benchmark data GDHY (Iizumi and Sakai [2020]) (red) and of individual and model ensemble (blue) across the extremes hot, dry, compound hot & dry, and wet. Positive yield anomalies show a model misinterpretation, where extremes interpreted to having positive impacts on crop growth.) Boxplots below the distribution plots depict the median and interquartile range (IQR). These plots are the pendant to the Figure 3 in the main text, there aggregated over all crop types.

Supplementary Figure 10: RICE (first growing season) – Raincloud plots for relative detrended yield anomalies (red=benchmark, gray=models) across extremes. Boxplots show median/IQR. Extreme conditions for rice seems to be a specific challenge. Anomaly distributions show greater leftward shifts for hot than dry extremes, hot & dry extremes are better represented. Wet event distributions typically centre around zero, except crover and pdssat models better capture wet impacts.

Note: medians may be skewed due to considerable number of outliers.

Supplementary Figure 11: WINTER WHEAT – Raincloud plots for relative detrended yield anomalies (red=benchmark, blue=models) across extremes. Boxplots show median/IQR. Models much better catch the signal of the extremes, except for wet extremes. For dry, hot, dry & hot mostly all models have their interquartile range in the negative area. The models show high variability. For hot & dry two models even overestimate the impact (simplace-lintul5 and epic-iiasa). Wet event distributions typically centre around zero, except pdssat model better captures the wet impact. lpj-guess and simplace-lintul5 misinterpret the signal and even show positive impact (median) by the extreme. Note: medians may be skewed due to considerable number of outliers.

Supplementary Figure 12: SOY – Raincloud plots for relative detrended yield anomalies (red=benchmark, blue=models) across extremes. Boxplots show median/IQR. Models show exhibit exceptional behaviour among each other for soy. This is especially expressed in the hot & dry case, where median values are broadly distributed over the negative range. While for dry and hot extremes underestimation dominates, in the case of compound hot & dry extremes most models overestimate the impact (when looking at the median). Wet event distributions typically centre around zero, except pdssat model better captures the wet impact.
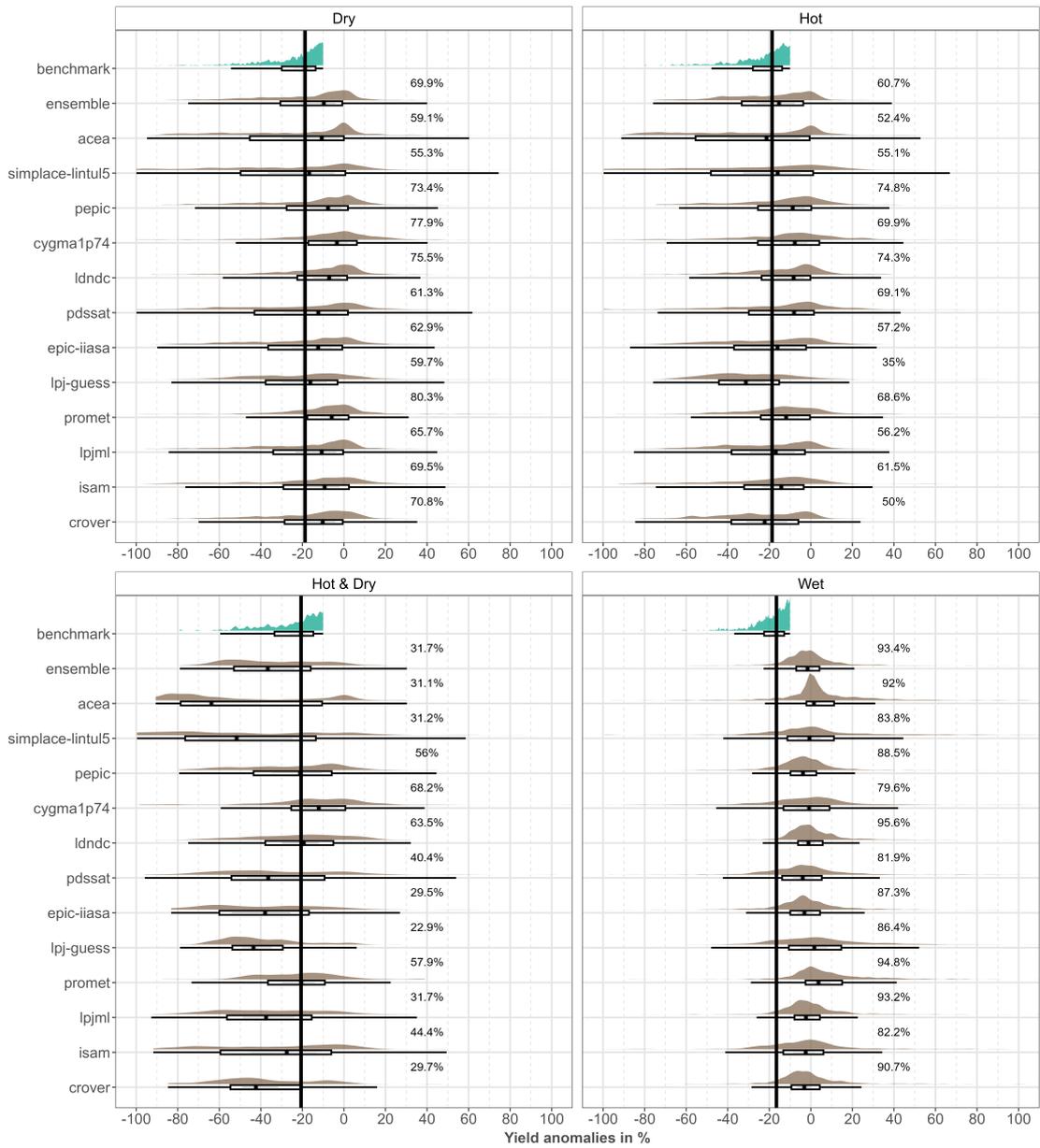
Supplementary Figure 13: MAIZE – Raincloud plots for relative detrended yield anomalies (red=benchmark, blue=models) across extremes. Boxplots show median/IQR. Maize has very few data points, still some trends can be detected. Dry and hot extremes are underestimated, but the negative signal is often detected. For compound hot & dry extremes more than half of the models have a median showing higher impacts than the benchmark data. Wet event distributions typically centre around zero, and models hardly catch the signal.

# 4 Raincloud plots: relative detrended yield anomalies for top five producers

Four main crop types - soy, maize, winter wheat, and first growing season rice- are assessed individually below. Instead of using all data as done earlier, we only look at extremes occurring in the top five producers for each crop respectively here. The Supplementary Figures 14, 15, 16, 17, 18, depict raincloud plots of the distribution of relative detrended yield anomalies (in %) of the Iizumi et al. benchmark data GDHY (Iizumi and Sakai [2020]) (red) and of individual and model ensemble (blue) across the extremes hot, dry, compound hot & dry, and wet. Positive yield anomalies show a model misinterpretation, where extremes interpreted to having positive impacts on crop growth.) Boxplots below the distribution plots depict the median and interquartile range (IQR). These plots are the pendant to the Figure 3 in the main text, there aggregated over all crop types and for all data points.

Supplementary Figure 14: Raincloud plots for relative detrended yield anomalies (red=benchmark, blue=models) of aggregated crop types of Iizumi et al. benchmark data GDHY (Iizumi and Sakai [2020]) and individual models or model ensemble median for climate extremes. Boxplots show median/IQR. The percentages shown alongside the distributions indicate the proportion of underestimation cases for each model and the ensemble. Only extreme events occurring in the top five producing countries—China (23.5% of global production), USA (23.1%), Brazil (7.14%), Argentina (4.39%), and India (2.99%)—are included in this figure. Consistent with the main findings, models tend to underestimate the effects of dry, hot, and wet extremes even when restricted to the top-producing countries.

Supplementary Figure 15: RICE (first growing season) – Raincloud plots for relative detrended yield anomalies (red=benchmark, blue=models) of Iizumi et al. benchmark data GDHY (Iizumi and Sakai [2020]) and individual models or model ensemble median for climate extremes. Boxplots show median/IQR. The percentages shown alongside the distributions indicate the proportion of underestimation cases for each model and the ensemble. Only extreme events occurring in the top five producing countries—China (30.7% of global production), India (14.2%), Indonesia (6.98%), Thailand (6.71%), and Brazil (6.35%)—are included in this figure. Consistent with the main findings, models tend to underestimate the effects of climate extremes even when restricted to the top-producing countries.

15

Top 5 producers: USA (14.1%), FRA (12.7%), CHN (11.4%), DEU(8.53%), TUR(5.97%)

Supplementary Figure 16: WINTER WHEAT – Raincloud plots for relative detrended yield anomalies (red=benchmark, blue=models) of Iizumi et al. benchmark data GDHY (Iizumi and Sakai [2020]) and individual models or model ensemble median for climate extremes. Boxplots show median/IQR. The percentages shown alongside the distributions indicate the proportion of underestimation cases for each model and the ensemble. Only extreme events occurring in the top five producing countries—USA (14.1% of global production), France (12.7%), China (11.4%), Germany (8.53%), Türkiye (5.97%)—are included in this figure. Consistent with the main findings, models tend to underestimate the effects of dry, hot, and wet extremes even when restricted to the top-producing countries.

Supplementary Figure 17: SOY – Raincloud plots for relative detrended yield anomalies (red=benchmark, blue=models) of Iizumi et al. benchmark data GDHY (Iizumi and Sakai [2020]) and individual models or model ensemble median for climate extremes. Boxplots show median/IQR. The percentages shown alongside the distributions indicate the proportion of underestimation cases for each model and the ensemble. Only extreme events occurring in the top five producing countries—USA (41.4% of global production), Brazil (28.2%), Argentina (16.8%), China (7.8%), and India (1.37%)—are included in this figure. Consistent with the main findings, models tend to underestimate the effects of dry and wet extremes even when restricted to the top-producing countries. For hot and compound hot and dry extremes, some models generally overestimate the impacts instead.
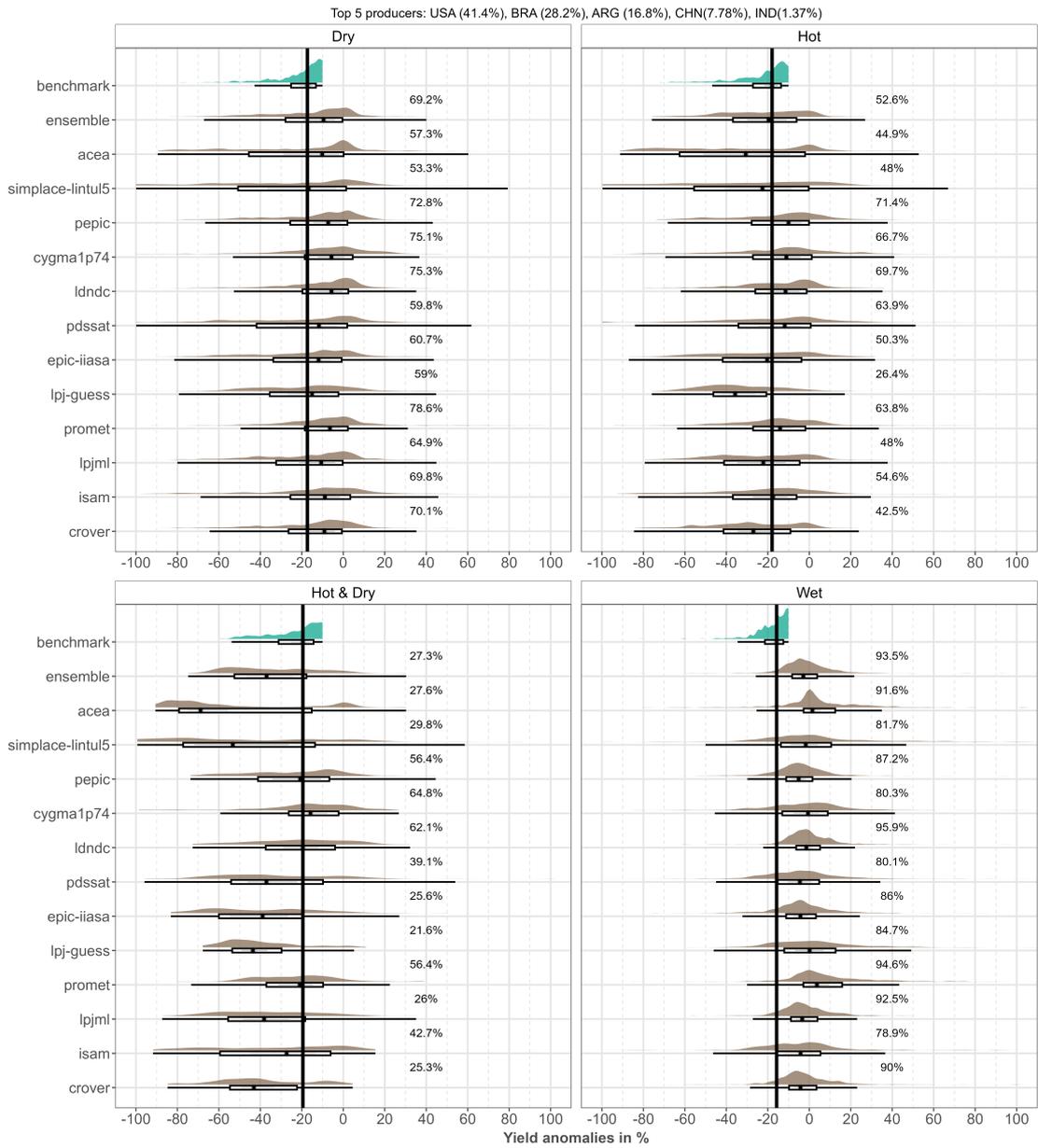
Supplementary Figure 18: MAIZE – Raincloud plots for relative detrended yield anomalies (red=benchmark, blue=models) of Iizumi et al. benchmark data GDHY (Iizumi and Sakai [2020]) and individual models or model ensemble median for climate extremes. Boxplots show median/IQR. The percentages shown alongside the distributions indicate the proportion of underestimation cases for each model and the ensemble. Only extreme events occurring in the top five producing countries—USA (40.6% of global production), China (16.4%), Brazil (7.91%), Argentina (4.35%), and France (4.24%)—are included in this figure. The models show a mixed response: while some generally underestimate the effects of climatic extremes, others tend to overestimate them—except in the case of wet extremes, where all models systematically underestimate the impacts.

# 5 Raincloud plots: absolute detrended yield anomalies

Four main crop types - first growing season rice, winter wheat, soy, and maize - are assessed together and individually below. The Supplementary Figures 19, 20, 21, 22,23 depict raincloud plots of the absolute detrended distribution of yield anomalies (in tons/ha) of the Iizumi et al. benchmark data GDHY (Iizumi and Sakai [2020]) (red) and of individual and model ensemble (blue) across the extremes hot, dry, compound hot & dry, and wet. Positive yield anomalies show a model misinterpretation, where extremes are assigned to having positive impacts on crop growth.) Boxplots below the distribution plots depict the median and interquartile range (IQR).



Supplementary Figure 19: Raincloud plots for absolute detrended yield anomalies (red=benchmark, blue=models) of aggregated crop types of Iizumi et al. benchmark data GDHY (Iizumi and Sakai [2020]) and individual models or model ensemble median (tons/ha) for climate extremes. Boxplots show median/IQR. Although the median shows that in most cases the models are able to capture a negative signal by the respective extreme they are seldomly able to simulate the magnitude.
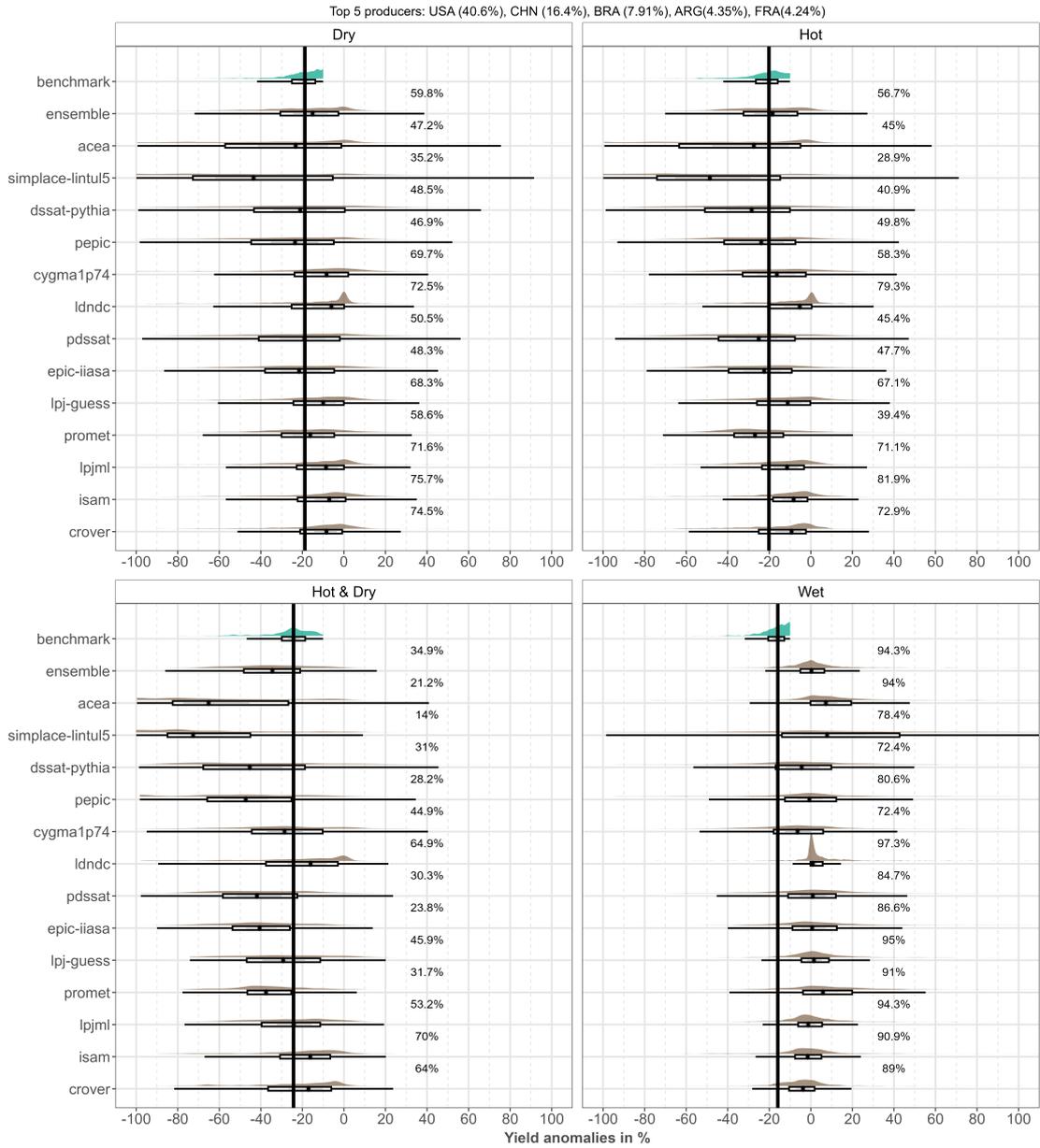
Supplementary Figure 20: RICE (first growing season) – Raincloud plots for absolute detrended yield anomalies (red=benchmark, blue=models) of Iizumi et al. benchmark data GDHY (Iizumi and Sakai [2020]) and individual models or model ensemble median (tons/ha) for climate extremes. Boxplots show median/IQR. For dry and hot extremes the anomaly distributions except for pepic mostly centre around zero. For wet as well, with the exception of acea which misinterprets the signal as something positive. For hot & dry combined the distribution of nearly half of the models shifts to the left and thus, better represents the negative impact of the signal.

Supplementary Figure 21: WINTER WHEAT – Raincloud plots for absolute detrended yield anomalies (red=benchmark, blue=models) of Iizumi et al. benchmark data GDHY (Iizumi and Sakai [2020]) and individual models or model ensemble median (tons/ha) for climate extremes. Boxplots show median/IQR. Some models exhibit strong outliers for wheat under all extreme types. Median values range from over to underestimation, with lpjml performing having close agreement of its median with the median of the benchmark for dry, hot, hot & dry. Wet extremes are consequently underestimated and anomalies centre again for most models around zero.
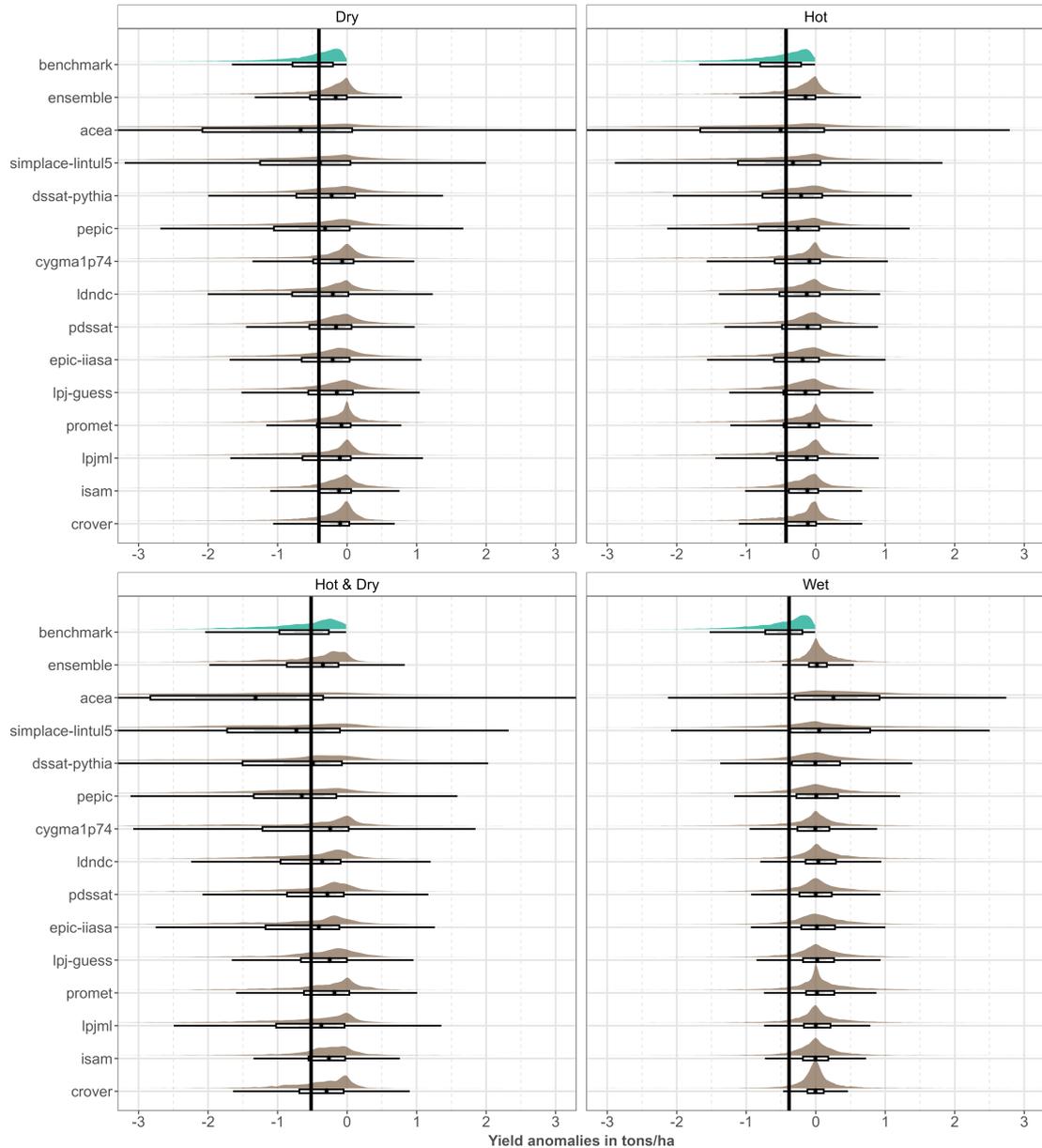
Supplementary Figure 22: SOY – Raincloud plots for absolute detrended yield anomalies (red=benchmark, blue=models) of Iizumi et al. benchmark data GDHY (Iizumi and Sakai [2020]) and individual models or model ensemble median (tons/ha) for climate extremes. Boxplots show median/IQR. The heterogeneity we saw for the relative anomalies above in Figure 12 is not as expressed for dry and hot extremes here with an absolute error. Many models have a median that is very close to the median of the benchmark data. lpjml consequently overestimates the impact (except for wet) and has an unusual high variance. Hot & dry has a many models where the median even exceeds the benchmark median, and again for wet extremes the distribution of the models mostly centres around zero, except this time for lpjml that has tendency to better detect the negative signal.
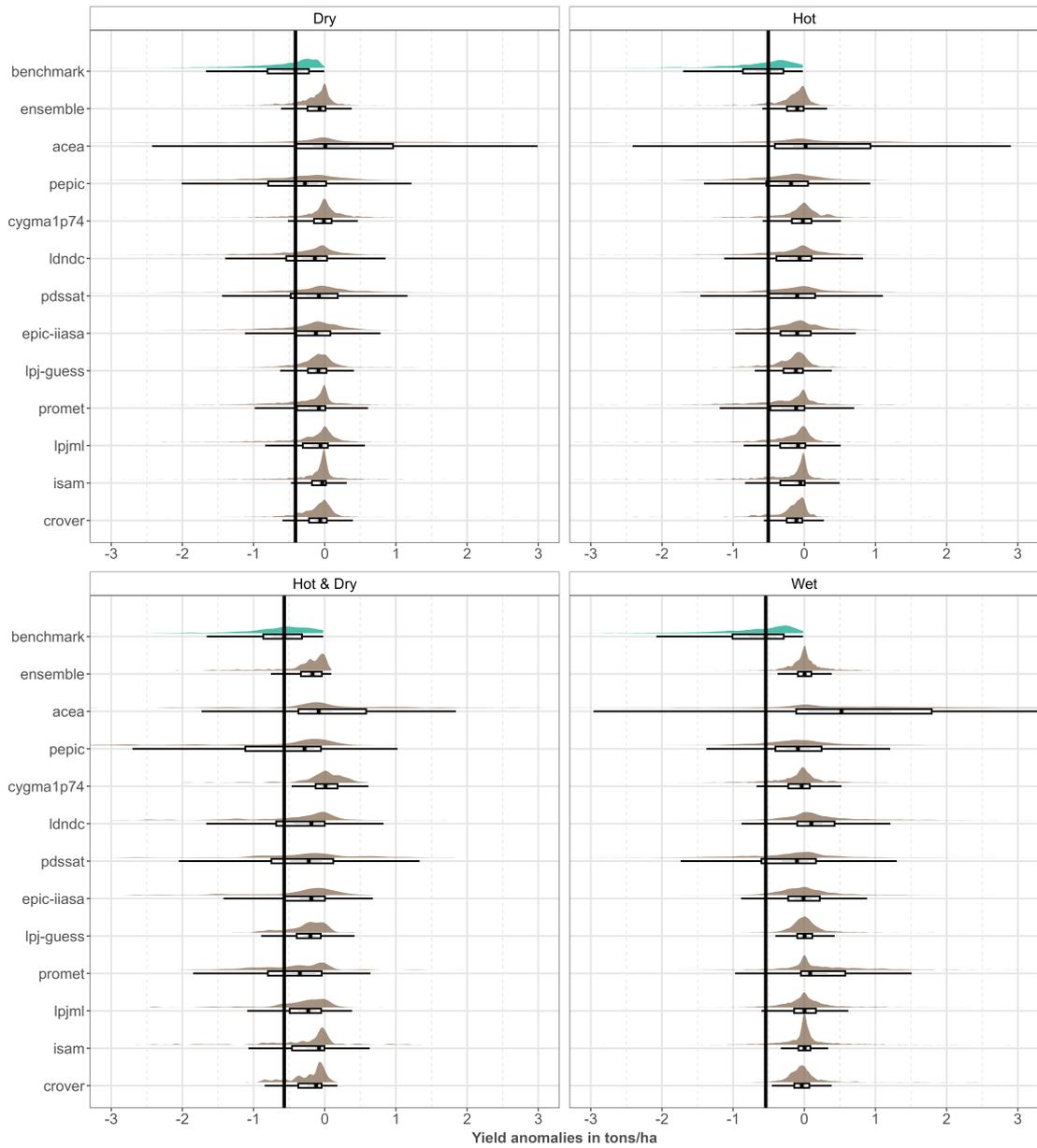
Supplementary Figure 23: MAIZE – Raincloud plots for absolute detrended yield anomalies (red=benchmark, blue=models) of Iizumi et al. benchmark data GDHY (Iizumi and Sakai [2020]) and individual models or model ensemble median (tons/ha) for climate extremes. Boxplots show median/IQR. Compared to the other crops we see that already the benchmark data for maize has a much broader range in terms of absolute anomalies. We also see that the models exhibit a broad range of outliers across dry, hot, compound hot & dry extremes. For many models the median in the case of dry extremes is close to that one of the benchmark data. Less so for hot extremes. For compound hot & dry extremes the models exhibit high variance and strong outliers in the positive and negative direction. For wet extremes the distribution of the models mostly centres around zero, except for acea and simplace-lintul5 misinterpreting the signal in the positive direction.

# 6 Raincloud plots: relative model error

To also depict the extent of the error, we display for four main crop types - first growing season rice, winter wheat, soy, and maize - the distribution of the individual and ensemble relative model error in the Supplementary Figures 24, 25, 26, 27,28. The relative model error is quantified as the difference between observed (Iizumi et al. benchmark data GDHY (Iizumi and Sakai [2020])) and simulated yield anomalies in percent (i.e. based on relative detrending) across the extremes hot, dry, compound hot & dry, and wet. Positive relative model errors represent an overestimation of the negative impact. Boxplots below the distribution plots depict the median and interquartile range (IQR).

Supplementary Figure 24: Raincloud plot illustrating the distribution of individual and ensemble relative model error values across different extreme events, aggregated over four main crop types: first growing season rice, winter wheat, soy, and maize. The relative model error is quantified as the difference between observed and simulated yield anomalies in percent (i.e. based on relative detrending). Red areas indicate underestimation by the model (percentage denoted in red), while grey areas represent overestimation. Red percentage numbers indicate data points with underestimation. Boxplots below the distribution plots depict the median and interquartile range. For dry, hot, and wet extremes, all models underestimate the impact more frequently than they overestimate it. Notably, wet extremes are the least accurately captured by the models. The ensemble consistently ranks among the lower-performing results across all event types. For compound hot & dry extremes, all models demonstrate reduced underestimation. Models with high variance exhibit a median error close to zero, however, still coupled with strong overestimation.

Supplementary Figure 25: RICE (first growing season) – Raincloud plots for model error across extremes (relative detrended anomaly difference: benchmark - model, in %), with red=underestimation, gray=overestimation. Boxplots show median/IQR, red percentage numbers indicate data points with underestimation. All extreme types lead to a rather homogeneous picture of underestimation. .

Supplementary Figure 26: WINTER WHEAT – Raincloud plots for model error across extremes (relative detrended anomaly difference: benchmark - model, in %), with red=underestimation, gray=overestimation. Boxplots show median/IQR, red percentage numbers indicate data points with underestimation. Dry, hot, and compound hot & dry extremes lead to a rather homogeneous picture of underestimation. Only simplace-lintul5 and epic-iiasa overestimate more often for compound hot & dry extremes. For wet extremes we get a rather homogeneous picture of underestimation again.

Supplementary Figure 27: SOY – Raincloud plots for model error across extremes (relative detrended anomaly difference: benchmark - model, in %), with red=underestimation, gray=overestimation. Boxplots show median/IQR, red percentage numbers indicate data points with underestimation. In general, the models show smaller proportions of underestimation. For compound hot & dry extremes, the models more often overestimated than underestimated. For wet extremes we get a rather homogeneous picture of underestimation.

Supplementary Figure 28: MAIZE – Raincloud plots for model error across extremes (relative detrended anomaly difference: benchmark - model, in %), with red=underestimation, gray=overestimation. Boxplots show median/IQR, red percentage numbers indicate data points with underestimation. dry extremes lead to slightly less underestimated data points than hot extremes. Compound hot & dry extremes lead to a very mixed picture - while some models keep the trend to underestimate the impact others even overestimate it (simplace lintul5, acea, pepic, dssat-pythia, epic-iiasa). The ensemble nearly overestimates (52.1%) as often as it underestimates (47.9%). For wet extremes the amount of underestimated data points ranks from 73.7% for cygma1p74 to 96.2% for the ensemble.

# 7 Raincloud plots: absolute model error

To also depict the extent of the error, we display for four main crop types - first growing season rice, winter wheat, soy, and maize - the distribution of the individual and ensemble absolute model error in the Supplementary Figures 29, 30, 31, 32,33. The absolute model error is quantified as the difference between observed (Iizumi et al. benchmark data GDHY (Iizumi and Sakai [2020])) and simulated yield anomalies in tons/ha (i.e. based on absolute detrending) across the extremes hot, dry, compound hot & dry, and wet. Positive absolute model errors represent an overestimation of the negative impact. Boxplots below the distribution plots depict the median and interquartile range (IQR).



Supplementary Figure 29: Raincloud plots for model error across extremes aggregated for four main crop types soy, maize, winter wheat, and first growing season rice (absolute detrended anomaly difference: benchmark - model, in tons/ha), with red=underestimation, gray=overestimation. Boxplots show median/IQR, red percentage numbers indicate data points with underestimation. Interestingly, the median absolute error is quite close in case of dry and hot extremes for most models. Simplace-lintul5 is close to zero and acea overestimates more than it underestimates in the median. For compound hot & dry extremes the results are more heterogeneous: some models shift with the median error close to zero (simplace-lintul5, dssat-pythia, pepic, cygma1p74, lndc, pdssat, epic-iiasa, lpjml).

Supplementary Figure 30: RICE (first growing season) – Raincloud plots for model error across extremes (absolute detrended anomaly difference: benchmark - model, in tons/ha), with red=underestimation, gray=overestimation. Boxplots show median/IQR, red percentage numbers indicate data points with underestimation. For all four extremes the models exhibit a homogeneous picture of underestimation, especially pronounced again for wet extremes. Acea is special out as with a wide range of underestimation errors (IQR) and a higher median error.

Supplementary Figure 31: WINTER WHEAT – Raincloud plots for model error across extremes (absolute detrended anomaly difference: benchmark - model, in tons/ha), with red=underestimation, gray=overestimation. Boxplots show median/IQR, red percentage numbers indicate data points with underestimation. Compared to the other crops, wheat results look more balanced, and many models have reduced data points of underestimation for dry and hot extremes, which flips for compound hot & dry, here acea, simplace-lintul5, epic-iiasa and lpjml even overestimate the loss more often than they underestimate it. Ldndc is perfectly in the middle. Wet extremes are again consequently underestimated.

Supplementary Figure 32: SOY – Raincloud plots for model error across extremes (absolute detrended anomaly difference: benchmark - model, in tons/ha), with red=underestimation, gray=overestimation. Boxplots show median/IQR, red percentage numbers indicate data points with underestimation. For soy under dry or hot extremes models have a much lower tendency to underestimate the impact. In the case of compound hot & dry only ldndc underestimates the impact. Wet extremes are again consequently underestimated.

Supplementary Figure 33: MAIZE – Raincloud plots for model error across extremes (absolute detrended anomaly difference: benchmark - model, in tons/ha), with red=underestimation, gray=overestimation. Boxplots show median/IQR. Although nearly all models underestimate dry or hot extremes, over and underestimation are more balanced. In the case of compound hot & dry extremes more models overestimate the extreme. However, the distributions become here exceptionally wide and uncertainty seems high. Wet extremes are again consequently underestimated.

# 8 Decomposing the Kling-Gupta Efficiency

As the Kling-Gupta efficiency (KGE, Eq. 17) consists of three components: alpha (variability), beta (bias) and r (correlation), we decompose the composite values (see Figure 4 in the Main text) plotting regional heatmaps for alpha (Supplementary Figure 34), beta (Supplementary Figure 35), and r (Supplementary Figure 36) separately. Further, we provide the regional heatmap for the hit rate (Eq. 18) metric (Supplementary Figure 37). Note that all metrics here are based on absolute detrending (Eq. 6). We also create a global map showing to which region each grid cell is allocated (Figure 38).

Supplementary Figure 34: Alpha of Kling Gupta Efficiency indicating model variability. Values $\geq 1$ indicate a higher variability in the respective model simulation data than the variability in the benchmark data. The closer to 1, the better the model performance as the variabilities are more equal. We note that the models generally show higher variability than the benchmark data across all regions and extreme types. For Middle Africa, Central America, Western Europe, and Southern Africa we see that many models have a variability that is about 3 times or sometimes even up to 10 times higher than the benchmark data.

**Dry**

| | crover | isam | lpjml | promet | lpj-guess | epic-iiasa | pdssat | ldndc | cygma1p74 | pepic | dssat-pythia | simplace-lintul5 | acea | ensemble |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Western Europe | 1.86 | 3.08 | 1.4 | 1.98 | 2 | 1.53 | 1.32 | 2.23 | 0.8 | 1.67 | -0.14 | 0.14 | 0.58 | 1.67 |
| Eastern Europe | 0.09 | 0.89 | 0.72 | 0.9 | 0.66 | 0.36 | 0.41 | 0.5 | 0.75 | -0.07 | -0.14 | -0.42 | -0.77 | 0.47 |
| Southern Europe | 1.03 | 1.58 | 0.9 | 1.44 | 1.22 | 1.07 | 1.32 | 1.21 | 1.32 | 1.18 | 1.17 | 0.67 | 0.64 | 1.2 |
| Northern Europe | 3.25 | 1.45 | 2.26 | 1.35 | 3.38 | 2.3 | 2.09 | 2.53 | | 2.41 | | 2.67 | 2.23 | 2.21 |
| Australia and New Zealand | 0.5 | 0.6 | -0.81 | 0.57 | -0.05 | 0.11 | 0.61 | -0.22 | 0.09 | -0.87 | -0.18 | -0.94 | -3.23 | 0.01 |
| Northern America | 0.36 | 0.78 | 0.71 | 0.78 | 0.49 | 0.46 | 0.44 | 0.48 | 0.23 | 0.21 | -0.19 | 0.28 | -0.8 | 0.41 |
| Central America | 0.32 | -0.42 | -0.47 | 0.6 | -0.9 | -0.24 | 0.31 | -0.35 | -0.75 | -0.46 | -0.65 | -0.74 | -2.3 | -0.24 |
| South America | 0.91 | 0.8 | 0.52 | 0.82 | 0.83 | 0.67 | 0.69 | 0.89 | 1.22 | 0.27 | 1.2 | 0.02 | -0.3 | 0.71 |
| Western Asia | 0.44 | 0.66 | -0.35 | 0.32 | 0.15 | 0.23 | 0.18 | 0.03 | 0.99 | 0 | -0.95 | -0.48 | -1 | 0.16 |
| Central Asia | 0.03 | 0.59 | 0.85 | 0.96 | 1.08 | 0.54 | 0.59 | 0.7 | 0.6 | -0.41 | 0.2 | 1.1 | -0.77 | 0.63 |
| Southern Asia | 0.76 | 0.66 | 0.3 | 0.76 | 0.57 | 0.43 | 0.47 | 0.68 | 0.63 | 0.25 | 0.59 | 0.6 | 0.34 | 0.62 |
| Eastern Asia | 0.79 | 0.2 | 0.44 | 0.82 | 0.76 | 0.44 | 0.82 | 0.51 | 0.36 | 0.64 | 0.35 | 0.6 | 0.44 | 0.59 |
| South-Eastern Asia | 1.04 | 1.15 | 0.98 | 0.77 | 0.94 | 0.92 | 0.95 | 0.88 | 1.24 | 0.81 | 0.87 | 0.43 | 0.51 | 0.97 |
| Northern Africa | 0.9 | 0.68 | 1.1 | 0.91 | 0.71 | 0.89 | 1.44 | 0.75 | 1.13 | -0.52 | 0.57 | -2.66 | -1.44 | 0.86 |
| Western Africa | 1.24 | 0.88 | 0.89 | 1.14 | 0.76 | 0.93 | 0.98 | -0.17 | 1.38 | -0.34 | 0.83 | -1.83 | -2.5 | 0.91 |
| Middle Africa | 0.78 | -0.23 | 0.48 | 0.45 | 0.11 | 0.46 | 1.09 | -2.28 | 0.96 | -0.81 | 0.71 | -4.67 | -3.55 | 0.3 |
| Eastern Africa | 0.8 | 0.3 | 0.63 | 0.56 | 0.3 | 0.36 | 0.59 | -0.39 | 0.8 | -0.68 | 0.63 | -1.42 | -3.2 | 0.41 |
| Southern Africa | 0.89 | -1.19 | -0.7 | 0 | -0.34 | -0.59 | -0.08 | -0.64 | 1.07 | -1.35 | -0.33 | -0.62 | -2.88 | -0.39 |

**Hot**

| | crover | isam | lpjml | promet | lpj-guess | epic-iiasa | pdssat | ldndc | cygma1p74 | pepic | dssat-pythia | simplace-lintul5 | acea | ensemble |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Western Europe | 1.61 | 2.1 | 1.45 | 0.46 | 2.01 | 1 | 0.42 | 1.58 | 1.63 | 1.91 | 0.1 | -0.22 | 1 | 1.25 |
| Eastern Europe | -0.13 | 0.55 | 0.18 | 0.5 | 0.02 | -0.1 | -0.19 | 0.37 | 0.58 | -0.69 | -0.63 | -0.47 | -1.42 | 0.12 |
| Southern Europe | 1.37 | 1.46 | 1.4 | 1.34 | 1.19 | 1 | 1.3 | 1.39 | 1.22 | 1.13 | 1.18 | 0.96 | 1.05 | 1.28 |
| Northern Europe | 1.54 | 1.62 | 1.59 | 0.75 | 4.95 | 1.55 | 1.43 | 0.65 | | 1.82 | | -0.19 | 1.76 | 1.55 |
| Australia and New Zealand | 0.47 | 0.79 | 0.4 | 1.01 | 0.54 | 0.51 | 0.8 | 0.36 | 0.79 | -0.35 | 0.54 | -0.22 | -1.15 | 0.49 |
| Northern America | 0.65 | 1.06 | 0.12 | 0.67 | 0.63 | -0.13 | 0.35 | 0.87 | -0.22 | 0.1 | -0.71 | -0.26 | -1.57 | 0.3 |
| Central America | -0.13 | -0.81 | -0.91 | 0.17 | -0.15 | -0.79 | -0.23 | -1.39 | -0.79 | -1.03 | -1.47 | -2.4 | -4.62 | -0.68 |
| South America | 0.89 | 1.05 | 0.67 | 0.85 | 0.72 | 0.88 | 0.94 | 0.94 | 0.93 | 0.52 | 1 | 0.02 | -0.18 | 0.82 |
| Western Asia | 0.26 | 0.57 | 0.31 | 0.68 | 0.26 | 0.47 | 0.57 | 0.4 | 0.92 | 0.47 | 0.22 | -0.54 | -0.18 | 0.46 |
| Central Asia | -0.35 | 0.51 | 0.87 | 1.32 | 1.09 | 0.71 | 0.83 | 0.85 | 0.72 | -0.7 | 0.44 | 1.43 | -0.75 | 0.74 |
| Southern Asia | 0.86 | 0.99 | 0.87 | 0.83 | 0.77 | 0.67 | 0.81 | 0.96 | 0.97 | 0.63 | 0.57 | 0.75 | 0.67 | 0.81 |
| Eastern Asia | 0.67 | -0.12 | 0.03 | 0.32 | 0.54 | 0.18 | 0.82 | 0.33 | 0.39 | 0.62 | 0.13 | 0.27 | 0.32 | 0.42 |
| South-Eastern Asia | 0.95 | 0.99 | 0.81 | 0.94 | 1.18 | 0.83 | 0.88 | 0.81 | 0.79 | 0.35 | 0.81 | 1.09 | -0.2 | 0.9 |
| Northern Africa | 0.71 | 0.54 | 0.73 | 1 | 0.63 | 0.81 | 1.18 | 0.69 | 0.79 | 0.88 | 0.09 | 0.08 | 0.43 | 0.81 |
| Western Africa | 1.19 | 0.82 | 1.06 | 1.25 | 0.74 | 1.24 | 0.67 | 0.32 | 1.2 | 0.61 | 1.12 | -0.08 | -1.56 | 1 |
| Middle Africa | 0.9 | 0.83 | 1 | 0.91 | 1.84 | 1.2 | 1.24 | 0.43 | 0.89 | 1.78 | 1.14 | 0.69 | 2.67 | 1.03 |
| Eastern Africa | 0.81 | 0.64 | 0.82 | 0.59 | 0.61 | 0.85 | 0.81 | 0.5 | 0.81 | 0.5 | 0.78 | 0.25 | -1.44 | 0.7 |
| Southern Africa | 0.77 | -0.32 | 0.1 | 0.27 | 0.4 | 0.19 | 0.53 | 0.35 | 0.92 | -0.79 | 0.33 | -0.11 | -0.88 | 0.22 |

**Hot & Dry**

| | crover | isam | lpjml | promet | lpj-guess | epic-iiasa | pdssat | ldndc | cygma1p74 | pepic | dssat-pythia | simplace-lintul5 | acea | ensemble |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Western Europe | 1.36 | 1.61 | 1.1 | 0.4 | 0.85 | 0.5 | 0.32 | 1.24 | -0.03 | 0.86 | -0.76 | -0.3 | -0.24 | 0.67 |
| Eastern Europe | -0.17 | 0.15 | -0.14 | 1.49 | -0.05 | 0.21 | 0.55 | 0.36 | 0.31 | -1.54 | -0.57 | -0.8 | -1.89 | 0.08 |
| Southern Europe | 1.42 | 1.43 | 1.05 | 1.29 | 1.07 | 0.88 | 1.08 | 1.32 | 1.05 | 1.27 | 0.38 | 0.28 | 0.39 | 1.13 |
| Northern Europe | 1.52 | 3.29 | 2.99 | 0.67 | 4.08 | 1.82 | 2.98 | 2.18 | | 6.25 | | -0.79 | 3.12 | 2.73 |
| Australia and New Zealand | 0.67 | 1.15 | 0.26 | 1.18 | 0.41 | 0.43 | 1.15 | 0.37 | 0.23 | -0.6 | 0.85 | -0.6 | -2.3 | 0.51 |
| Northern America | 0.5 | 0.88 | -0.42 | 0.62 | 0.53 | -0.81 | -0.34 | 0.32 | -0.85 | -0.43 | -2.2 | -0.61 | -3.01 | -0.22 |
| Central America | -0.37 | -1.44 | -2.67 | -0.81 | -2.03 | -1.73 | -0.63 | -2.9 | -0.96 | -1.06 | -3.63 | 8.02 | -7.44 | -1.78 |
| South America | 0.91 | 0.91 | 0.86 | 0.94 | 0.49 | 0.48 | 0.68 | 0.57 | 1.02 | -0.09 | 0.63 | -0.89 | -1.69 | 0.65 |
| Western Asia | -0.08 | 0.39 | -0.51 | 0.65 | -0.44 | 0.14 | 0.47 | -0.55 | 1.42 | 0.19 | -1.92 | -1.07 | -1.46 | 0.02 |
| Central Asia | -0.76 | -0.02 | 0.05 | 1.68 | 1.06 | 0.28 | 0.41 | 0.43 | 0.32 | -1.4 | -0.06 | 1.14 | -1.17 | 0.33 |
| Southern Asia | 0.71 | 0.75 | 0.49 | 0.76 | 0.58 | 0.54 | 0.3 | 0.82 | 0.98 | 0.24 | 0.71 | 0.63 | -0.01 | 0.63 |
| Eastern Asia | 0.51 | -0.78 | -0.79 | -0.07 | 0.43 | -0.7 | 0.94 | -0.63 | 0.47 | 0.12 | -0.43 | -0.62 | -1.3 | -0.17 |
| South-Eastern Asia | 0.93 | 0.92 | 0.66 | 0.99 | 1.18 | 0.91 | 0.85 | 0.33 | 0.9 | 0.22 | 0.8 | 0.93 | -2.25 | 0.8 |
| Northern Africa | 0.8 | 0.29 | 0.92 | 0.93 | 1.2 | 1.04 | 0.97 | 1.23 | 0.42 | 0.95 | 0.51 | 1.22 | 0.2 | 0.91 |
| Western Africa | 1.07 | 0.32 | 1.08 | 1.04 | 0.56 | 1.01 | 0.75 | -1.09 | 1.27 | -1.31 | 0.79 | -2.38 | -5.x | 0.63 |
| Middle Africa | 0.95 | 0.43 | 0.85 | 0.62 | 0.84 | 0.34 | 0.48 | -1.84 | 0.88 | -0.37 | 0.03 | -3.18 | -3.79 | 0.31 |
| Eastern Africa | 1.01 | 0.45 | 0.85 | 0.39 | 0.52 | 0.72 | 0.7 | -0.2 | 0.9 | -0.83 | 0.76 | -1.79 | -4 | 0.57 |
| Southern Africa | 0.59 | -0.58 | -1.04 | -0.06 | -0.37 | -0.51 | -0.12 | -0.96 | 1.9 | -0.41 | -0.29 | -0.06 | -1.92 | -0.31 |

**Wet**

| | crover | isam | lpjml | promet | lpj-guess | epic-iiasa | pdssat | ldndc | cygma1p74 | pepic | dssat-pythia | simplace-lintul5 | acea | ensemble |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Western Europe | 3.33 | 3.14 | 2.9 | 4.05 | 3.52 | 2.52 | 1.97 | 3.38 | 2.46 | 2.62 | 3.52 | 3.91 | 3.54 | 3.03 |
| Eastern Europe | 0.84 | 0.87 | 0.7 | 1.41 | 0.73 | 0.85 | 0.75 | 0.88 | 0.52 | 0.61 | 1.49 | 1.37 | 0.9 | 0.78 |
| Southern Europe | 1.88 | 1.77 | 1.77 | 1.68 | 1.83 | 2.03 | 1.74 | 1.7 | 2.23 | 1.66 | 2.74 | 1.99 | 2 | 1.81 |
| Northern Europe | 3.25 | 2.95 | 2.21 | 2.42 | 3.81 | 2.35 | 2.27 | 2.42 | | 2.39 | | 5.01 | 2.62 | 2.56 |
| Australia and New Zealand | 1.34 | 1.55 | 1.95 | 2.05 | 2.17 | 1.88 | 1.81 | 1.83 | 1.14 | 2.19 | 1.66 | 2.04 | 3.01 | 1.82 |
| Northern America | 1.65 | 1.62 | 1.75 | 1.71 | 1.48 | 1.56 | 1.51 | 1.55 | 1.4 | 1.39 | 2.15 | 2.52 | 2.39 | 1.63 |
| Central America | 0.81 | 0.86 | 0.81 | 1.26 | 0.82 | 0.78 | 1.41 | 1.74 | 0.12 | 0.88 | 0.68 | 1.08 | 0.72 | 0.84 |
| South America | 1.18 | 1.13 | 1.24 | 1.27 | 1.27 | 1.27 | 0.99 | 1.25 | 1.3 | 1.18 | 1.27 | 1.69 | 1.46 | 1.18 |
| Western Asia | 1.09 | 1.21 | 1.53 | 1.13 | 1.36 | 1.38 | 1.38 | 1.12 | 1.55 | 1.1 | 1.35 | 0.91 | 1.45 | 1.24 |
| Central Asia | 1.04 | 0.99 | 1.15 | 0.8 | 0.85 | 1.16 | 0.94 | 0.94 | 0.83 | 1.17 | 0.73 | 0.94 | 1.09 | 0.92 |
| Southern Asia | 1.11 | 1.11 | 1.35 | 1.1 | 1.15 | 1.12 | 1.13 | 1.32 | 1.03 | 1.12 | 1.15 | 0.99 | 1.28 | 1.11 |
| Eastern Asia | 1.37 | 1.45 | 1.35 | 2.03 | 1.47 | 1.4 | 1.31 | 1.74 | 0.95 | 1.35 | 1.81 | 1.62 | 1.97 | 1.42 |
| South-Eastern Asia | 1.41 | 1.52 | 1.53 | 1.56 | 1.49 | 1.47 | 1.47 | 1.65 | 1.31 | 1.5 | 1.3 | 1.55 | 1.75 | 1.51 |
| Northern Africa | 1.12 | 1.46 | 0.73 | 1.19 | 0.74 | 0.9 | 0.9 | 1.81 | 0.84 | 0.86 | 1.5 | 2.13 | 2.19 | 0.74 |
| Western Africa | 1.15 | 1.08 | 1.17 | 1.31 | 1.17 | 0.91 | 0.88 | 1.3 | 0.96 | 0.76 | 0.85 | 0.34 | 1.81 | 1.1 |
| Middle Africa | 0.98 | 1.34 | 1.14 | 1.63 | 2.8 | 1.98 | 1.84 | 3.43 | 1.11 | 2.73 | 1.09 | 9.09 | 12.88 | 1.62 |
| Eastern Africa | 0.69 | 0.67 | 0.63 | 0.78 | 0.63 | 0.72 | 0.87 | 0.92 | 0.71 | 0.93 | 0.59 | 1.98 | 1.27 | 0.71 |
| Southern Africa | 1.27 | -0.65 | 0.64 | 0.29 | 0.91 | 0.27 | 0.45 | 0.33 | 1.83 | -0.94 | -0.07 | -0.11 | -2.56 | 0.31 |

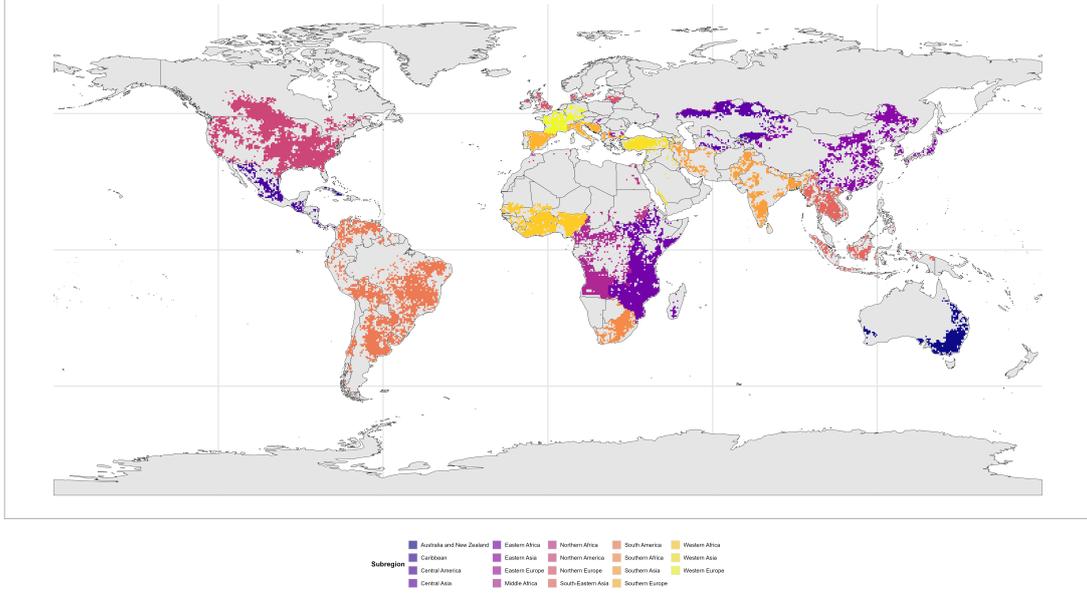Overestimation bias — No bias — Underestimation bias

Supplementary Figure 35: Beta of Kling-Gupta Efficiency indicating model bias in t/ha. Values close to 0 indicate lower model bias, positive values indicate an underestimation of climate extreme impacts on crops, negative values indicate overestimation. Generally the models tend to underestimate, especially for extreme wetness conditions and except for hot & dry, where we see a more balanced picture of over-and underestimation. We note that across all extremes, the models systematically underestimate climate extremes in (Western) Europe with especially large biases under extreme dry and extreme wet conditions. In contrast with the other models, Acea shows more frequent overestimation than underestimation. Almost all models overestimate in Central America except under wet extremes. In Southern Africa, biases are generally lower, especially under wet extremes the contrast with other regions is remarkable, with some models even overestimating the impacts.

Supplementary Figure 36: r of Kling-Gupta Efficiency indicating correlation between benchmark and simulated data. In Northern America, (Eastern) Europe, and Southern Africa we observe the most positive correlation values. The simulation data are negatively correlated with the benchmark observations in Northern Africa, especially under compound hot & dry extremes.

**Dry**

**Hot**

**Hot & Dry**

**Wet**

None captured | All captured

Supplementary Figure 37: Hit rate (dimensionless fraction) indicating the frequency of climate extremes captured by the respective model i.e. events where the model shows at least a negative detrended yields. The models typically capture about 70% of the dry or hot events, while less than 50% of the wet events. In general, the models seem to capture hot events slightly better than dry events. For compound hot & dry events we notice even higher hit rates close to 100%, except for regions close to the equator such as South America, Middle Africa, and Southern Asia. In general, the models show large regional differences with some models (e.g. ldndc under dry events) having the highest hit rates for a certain region (Western and Eastern Africa) but the lowest for another region (Western Europe). Regions where all models show consistently high hit rates are Australia & New Zealand for dry events and Northern America for compound hot & dry events.

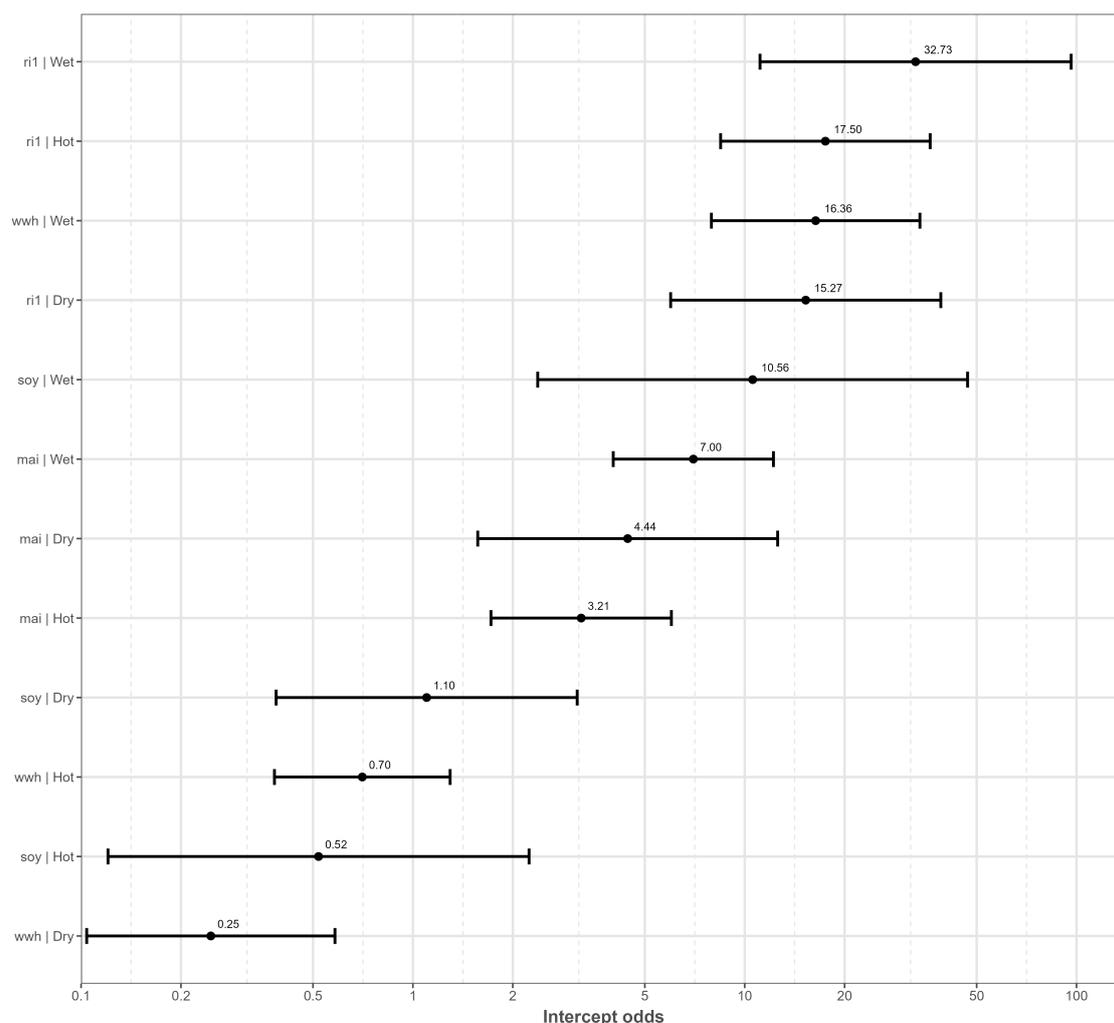Supplementary Figure 38: Regions used for the regional heatmaps

# 9 Model characteristics analysis: baseline (intercept) odds from GLMMs

In Supplementary Figure 39 we visualize the baseline (intercept) odds of underestimating yield losses under climatic extremes estimated from the Generalized Linear Mixed Models (GLMMs) we fitted separately for each crop (maize, rice, wheat, and soy) and different climatic extreme types (dry, wet, hot). See Eq. 22 for more details on how we defined the GLMM. The visualized intercept odds should be used along with the odds ratios shown in Figure 5 in the Main text to infer the absolute effects of the model characteristics on the odds of underestimation. From the intercept odds, one can calculate the probability of underestimation per crop and extreme type under the hypothetical absence of any model characteristic. For instance, for soybean under extreme wet conditions, a crop model without any of the analyzed characteristics is estimated to have an odds of underestimation of 10.56. This corresponds to a probability of underestimation of:

$$p_{\text{baseline}} = \frac{10.56}{10.56 + 1} = 91.3\%$$

Now suppose the presence of a model characteristic, for example, the inclusion of soil carbon and nitrogen cycling, which has an estimated odds ratio of 0.39. This means the odds decrease by a factor of 0.39. The new probability of underestimation is then:

$$p_{\text{soilCN}} = \frac{10.56 \times 0.39}{1 + (10.56 \times 0.39)} = 80.5\%$$

Supplementary Figure 39: Forest plot of the intercept (baseline) odds of underestimating the yield losses under specific climatic extremes inferred from the Generalized Linear Mixed Model (GLMM, Eq. 22). The errorbars represent the 90% confidence intervals and the numbers above the dot the mean point estimate. The odds values indicate the probability of underestimation relative to the probability of not underestimating the impacts of climatic extremes in the hypothetical absence of any of the crop model characteristics analyzed in this study. The odds ratios presented in Figure 5 from the Main analysis, indicate the estimated factors by which these baseline odds would change when including the corresponding model characteristic. The odds of underestimation is highest for rice under extreme wet conditions, while wheat under extreme dry conditions shows the lowest odds and hence probability of underestimation.

# 10 Methods

# 11 Model characteristics

Table 1 summarizes model characteristics relevant for simulating crop development and yield under climate change following e.g. tables in Müller et al. [2017] S1-S4. In this paper, the information on the characteristics is coded as follows: '1' indicates the characteristic is implemented, '0' means it is not included, and '-1' denotes the characteristic is not relevant for the model.

**Soil carbon and nitrogen cycling**
Simulating soil carbon and nitrogen (CN) cycling is crucial for representing nutrient availability and long-term soil fertility, which directly influence crop growth and yield.

- 1: Soil carbon and nitrogen cycling is explicitly simulated, allowing for dynamic nutrient feedbacks on crop growth.

- 0: Soil carbon and nitrogen cycling is not simulated; nutrient availability is statically or externally prescribed.

**Waterlogging stress**
Waterlogging occurs when excess soil water limits oxygen availability to roots, leading to reduced crop growth and potential yield loss.

- 1: Waterlogging stress is simulated, reducing crop growth under excessive soil moisture conditions.

- 0: Waterlogging is not simulated; crops do not experience stress from excess water.

**Explicit photosynthesis**
We distinguish two approaches of modeling light utilization: Photosynthesis-based (PS) models simulate physiological processes directly, while radiation use efficiency (RUE) models use empirical relationships.

- 1: Light utilization is modeled through explicit photosynthesis processes, capturing detailed physiological responses.

- 0: Light utilization is modeled using empirical radiation use efficiency relationships.

- -1: Not applicable; the model does not simulate this process.

**Complex temperature stress**
Temperature stress can impact crops in multiple ways, from basic reductions in photosynthesis and biomass to more complex effects on reproduction, phenology, and yield components.

- 1: Complex temperature stress is simulated, affecting multiple processes such as reproduction, phenology, and harvest index.

- 0: Only simple temperature stress is simulated, impacting photosynthesis and biomass accumulation.

**Variable roots**
Root development can be sensitive to temperature, influencing the plant's ability to access water and nutrients. Further, water availability (under drought or waterlogging conditions) can limit root expansion and function. Both, temperature and water stress are key drivers of root growth.

- 1: Root development responds dynamically to temperature or water stress. These signals affect root growth rates, distribution and function.

- 0: Root development is influenced by other stresses, but not by temperature or water stress.

- -1: Root development is not explicitly modeled in relation to stress.

**Dynamic harvest modeling**
The harvest index (HI) determines the proportion of biomass allocated to yield, and can be fixed or dynamically adjusted in response to environmental and management factors.

- 1: Harvest index and biomass allocation are dynamically simulated, responding to environmental and management changes.

- 0: Harvest index is fixed.

- -1: Not explicitly modeled.

| Model | Soil carbon & nitrogen cycling | Waterlogging stress | Explicit photosynthesis | Complex temperature stress | Variable roots | Dynamic harvest modeling |
|---|---|---|---|---|---|---|
| acea | 0 | 1 | -1 | 1 | 1 | 1 |
| crover | 0 | 0 | 1 | 0 | -1 | 0 |
| cygma1p74 | 0 | 1 | 0 | 0 | 0 | 1 |
| dssat-pythia | 1 | 1 | 0 | 0 | 1 | 1 |
| epic-iiasa | 1 | 1 | 0 | 0 | 1 | 1 |
| isam | 1 | 0 | 1 | 1 | 1 | 1 |
| ldndc | 1 | 0 | 1 | 0 | 0 | 1 |
| lpj-guess | 1 | 0 | 1 | 0 | 0 | 1 |
| lpjml | 1 | 0 | 1 | 0 | 0 | 0 |
| pepic | 1 | 1 | 0 | 0 | 1 | 1 |
| promet | 0 | 0 | 1 | 1 | 1 | 1 |
| simplace-lintul5 | 1 | 0 | 0 | 0 | 1 | 1 |
| pdssat | 1 | 0 | 1 | 1 | 1 | 1 |

Supplementary Table 1: Model comparison across key characteristics, as described in Subsection 11, for the applied ISIMIP3a model setup ISIMIP. Values indicate: 1 = characteristic is present; 0 = alternative representation; –1 = characteristic not applicable for the model.

# References

Toshichika Iizumi and Toru Sakai. The global dataset of historical yields for major crops 1981–2016. *Scientific Data*, 7(1), March 2020. ISSN 2052-4463. doi: 10.1038/s41597-020-0433-7. URL http://dx.doi.org/10.1038/s41597-020-0433-7.

Christoph Müller, Joshua Elliott, James Chryssanthacopoulos, Almut Arneth, Juraj Balkovic, Philippe Ciais, Delphine Deryng, Christian Folberth, Michael Glotter, Steven Hoek, Toshichika Iizumi, Roberto C. Izaurralde, Curtis Jones, Nikolay Khabarov, Peter Lawrence, Wenfeng Liu, Stefan Olin, Thomas A. M. Pugh, Deepak K. Ray, Ashwan Reddy, Cynthia Rosenzweig, Alex C. Ruane, Gen Sakurai, Erwin Schmid, Rastislav Skalsky, Carol X. Song, Xuhui Wang, Allard de Wit, and Hong Yang. Global gridded crop model evaluation: benchmarking, skills, deficiencies and implications. *Geoscientific Model Development*, 10(4):1403–1422, April 2017. ISSN 1991-959X. doi: 10.5194/gmd-10-1403-2017. URL https://gmd.copernicus.org/articles/10/1403/2017/. Publisher: Copernicus GmbH.