1 **Supporting Information for**

2 **An analysis of methodological and spatial differences in global cropping**
3 **systems models and maps**

4

5

6 Weston Anderson[1*], Liangzhi You[1,2*], Stanley Wood[3], Ulrike Wood-Sichra[1], Wenbin Wu[2]

7

8 [1] International Food Policy Research Institute
9 2033 K Street, NW
10 Washington, DC 20006, USA

11

12 [2] Key Laboratory of Agri-informatics,
13 Ministry of Agriculture / Institute of Agricultural Resources and Regional Planning
14 Chinese Academy of Agricultural Sciences,
15 Beijing 100081, China

16

17 [3] Bill&Melinda Gates Foundation
18 500 Fifth Avenue North
19 Seattle, WA 98109, USA

20

21 *For Correspondence: Weston.b.Anderson@gmail.com, L.You@cgiar.org, +1-202-862-8168(phone), +1-202-467-
22 4439(Fax)

23

24

25 Introduction

26

27     This supporting information provides details on the mathematical formulation of the downscaling

28 procedures for each model discussed in the main text, details of the Gaussian filter analysis, the

29 intermediate output for each analysis (the pixel-wise comparisons and the original cropping systems

30 maps for each product), and the full analysis for maize and rice. The data is intended to provide readers

31 with the complete details of each model specification and each analysis for those interested in a specific

32 model methodology or output.

33

## Appendix S1: Model Methodology

### S1.1 M3 distribution of crop harvested areas and yields:

In each grid cell that had agricultural inventory data the map of crop area was calculated as follows:

$$fcrop_i = fcropland_i \left( \frac{crop_k}{cropland_k} \right) \quad\quad (S1)$$

Where $fcrop_i$ is the harvested area of a specific crop in pixel $i$, $fcropland_i$ is the fraction of pixel $i$ designated as cropland, $crop_k$ is the harvested area of a specific crop in statistical reporting unit $k$ and $cropland_k$ is the amount of cropland in statistical reporting unit $k$. Yield was distributed uniformly across each grid cell as equivalent to the yield reported in the statistical reporting unit as a whole.

### S1.2 MIRCA distribution of harvested area:

MIRCA primarily reconciles the differences between Siebert et al., 2007 dataset of areas equipped for irrigation (AEI), Cropland extent of Ramankutty et al., 2008, and the harvested area (HA) maps of the M3 dataset (Monfreda et al., 2008) to provide a monthly cropping map for irrigated and rainfed crops. The priorities used to reconcile inconsistencies between the datasets is outlined in Table A1. MIRCA first produces a Condensed Crop Calendar of harvested area for each sub-crop $c$ and SRU $k$ ($HAccc_{c,k}$). A sub-crop is used to represent multi-cropping systems or different sub-groups of a crop that grow at different points in the year. For a complete description of how the Condensed Cropping Calendar is produced, see Portman et al., 2008.

Table S1: Priorities in distributing the Condensed Crop Calendars to monthly growing area grids

| Priority | Dataset | Goal |
|---|---|---|
| 1 | Area equipped for irrigation (Siebert et al., 2007) | In each month and grid cell the sum of irrigated crop-specific areas is lower than or equal to the area equipped for irrigation |
| 2 | Cropland extent (Ramankutty et al., 2008) | In each grid cell and month the sum of crop-specific irrigated and rainfed areas is lower than or equal to the cropland extent |
| 3 | Harvested crop area (Monfreda et al., 2008) | In each grid cell and for each crop class the annual sum of the irrigated and rainfed harvested crop area is equal to the total harvested area of the specific crop |

*Irrigated Crops:*

2

The MIRCA process is a production system specific cell-wise approach to disaggregating harvested area by month. Area equipped for irrigation and cropland extent both include fallow land in their definition, so these classes need not be used completely so long as the annual harvested area is disaggregated and designated as either rainfed or irrigated. Area equipped for irrigation is prioritized over cropland extent and harvested area in the following process:

1. Calculate the irrigated harvested area (IHA) for subcrop $c$ in cell $i$ of month $m$

$$IHA_{c,i,m} = \frac{(HA_{c,i,m} \times fAEI_{i,m})}{\sum subcrop\ types_{i,m}} \tag{S2}$$

Where $fAEI_{i,m}$ is the fraction of pixel $i$ equipped for irrigation in month $m$

2. Assign irrigated harvested areas to the monthly minimum of area equipped for irrigation and harvested area for subcrop $c$ in cell $i$ of month $m$

$if\ AEI > 0\ and\ CroplandExtent > 0\ and\ HA_{c,i,m,} > 0,\ IHA_{c,i,m} = \min(AEI_{c,i,m}, HA_{c,i,m})$ (S3)

3. If there still exists $HAccc_{ck}$ distribute irrigated growing areas to those cells that have cropland extent greater than zero and that are equipped for irrigation even if no $HA_{c,i,m}$ exists

$if\ AEI > 0\ and\ CroplandExtent > 0,\ IHA_{c,i,m} = remaining\ HAccc_{ck}$ (S4)

4. If there still exists $HAccc_{c,k}$, distribute it to areas within cells that are equipped for irrigation, even if the cropland extent (and therefore $HA_{c,i,m,}$) is zero.

$if\ AEI > 0\ and\ CroplandExtent = 0,\ IHA_{c,i,m} = remaining\ HAccc_{ck}$ (S5)

_Rainfed Crops:_
Following the distribution of irrigated crop areas, rainfed crops were distributed. Rainfed annual crops were treated differently than rainfed permanent crops. Annual crops were allowed to grow on areas equipped for irrigation so long as they were available, while permanent crops were not.

5. Calculate the rainfed harvested area for crop $c$ in pixel $i$ of month $m$ ($RHA_{c,i,m}$) by distributing rainfed crops to areas in which available cropland extent exceeds available area equipped for irrigation:

$if\ CroplandExtent > AEI,\ RHA_{c,i,m} = \min(HA_{c,i,m}, CroplandExtent_i)$ (S6)

91    6.  If there still exists HAccc$_{c,k}$, expand suitable areas in cells with more cropland extent than area

92        equipped for irrigation to 95% of the cell, leaving room to account for infrastructure.

93    $if\ CroplandExtent > AEI,\ RHA_{c,i,m,MIRCA} = \min(HA_{c,i,m}, (0.95 * Area_i))$        (S7)

94

95    7.  If there still exists HAccc$_{c,k}$, expand suitable areas in cells with either cropland extent or area

96        equipped for irrigation to 95% of the cell, leaving room to account for infrastructure.

97    $if\ CroplandExtent > 0\ or\ AEI > 0,\ RHA_{c,i,m,MIRCA} = \min(HA_{c,i,m}, (0.95 * Area_i))$    (S8)

98

99    The total harvested area of all rainfed and irrigated crops is therefore the sum of the IHA calculated in

100   steps A2-A5 and the RHA calculated in steps A6-A8.

101

102   **S1.3 SPAM harvested area and yield distribution**

103        The SPAM model distributes available crop statistics using a cross-entropy approach that

104   incorporates ancillary data on crop price, market access, biophysical suitability and expert

105   elicitation. Shannon (1948) first introduced the concept of information entropy to measure the

106   uncertainty of expected information in a system. Jaynes (1957) adopted the concept of

107   information entropy and further proposed the principle of maximum entropy in statistical

108   inference: the least informative probability distribution can be found by maximizing the entropy.

109   In other words, without information to the contrary, all possible states of a system are equally

110   likely. With respect to the concept's application to SPAM, Golan, Judge and Miller (1996)

111   presents the formation of maximum entropy (or minimum cross-entropy) principle for use in

112   parameter estimation problems.

113

114   *Harvested Area calculation*

115        SPAM distributes statistical information from allocation unit (e.g a country or a province)

116   by using the cropping intensity of crop *j* in production system *l* to convert the reported harvested

117   areas ($HA_{j,l}$) to physical areas ($CropArea_{jl}$) as:

118

119   $CropArea_{jl} = \dfrac{HA_{jl}}{CroppingIntensity_{jl}}\quad \forall j, l$                (S9)

120

121 SPAM next defines the area allocated to pixel *i* for crop *j* in production system *l* ($A_{ijl}$) using the

122 share of the total physical area for crop *j* in production system *l (Share$_{jl}$)* and the physical area

123 ($CropArea_{jl}$) as:

124

125 $A_{ijl} = CropArea_{jl} \times Share_{jl} \times s_{ijl}$   $\forall i, j, l$ (S10)

126

127 The minimum cross-entropy approach employed by the SPAM model calculates the area

128 shares for crop *j* of pixel *i* in production system *l* as:

129

130 $\min_{\{s_{ijl}\}} \left[ CE(s_{ijl}, \pi_{ijl}) = \sum_i \sum_j \sum_l s_{ijl} \ln(s_{ijl}) - \sum_i \sum_j \sum_l s_{ijl} \ln(\pi_{ijl}) \right]$ (S11)

131
132 *Subject to the following constraints:*

133 $\sum_{i \in k} s_{ijl} = 1$   $\forall j, i, k$ (S12)

134 $\sum_j \sum_l CropArea_{jl} \times Share_{jl} \times s_{ijl} \leq CroplandExtent_i$   $\forall i$ (S13)

135 $CropArea_{jl} \times Share_{jl} \times s_{ijl} \leq CropSuitableArea_{ijl}$   $\forall i, j, l$ (S14)

136 $\sum_{i \in k} \sum_l CropArea_{jl} \times Share_{jl} \times s_{ijl} = SubCropArea_{jk}$   $\forall k, j \in J$ (S15)

137 $\sum_{l \in L} CropArea_{jl} \times Share_{jl} \times s_{ijl} \leq AEI_i$   $\forall i$ (S16)

138 $1 \geq s_{ijl} \geq 0$   $\forall i, j, l$ (S17)

139

140 Where *l* may be irrigated, rainfed high-input, rainfed subsistence or rainfed low input.

141 *CroplandExtent$_i$* is the total extent of cropland for pixel *i* and $CropSuitableArea_{ijl}$ is the area

142 suitable for crop *j* at input level *l* in pixel *i*. $SubCropArea_{jk}$ is the crop area statistics for crop *j*

143 in subnational SRU *k*. *AEI$_i$* is the area equipped for irrigation in pixel *i*. *J* is a set of commodities

144 for which sub-national production statistics exist and *L* is a set of commodities within pixel *i* that

145 are irrigated.  $\pi_{ijl}$ represents the prior estimate of area shares for crop *j* at input level *l* in pixel *i*.

146 The prior is developed using expert elicitation where available and elsewhere is

147 calculated based on potential unit revenue, $Rev_{ijl}$.

148

149 $Rev_{ijl} = Share_{jl} \times Price_j \times Access_{ij} \times SuitableYield_{ijl}$ (S18)

150

151 Where $Price_j$ is the price of crop $j$ $Access_{ij}$ is a measure of the physical accessibility of the

152 market for crop $j$ from pixel $i$. $SuitableYield_{ijl}$ is the agro-climatically suitable yield for crop $j$

153 at input level $l$ in pixel $i$. Then the prior allocation of crop area is estimated using irrigated area

154 and cropland as follows.

155

156 $PriorArea_{ijl} = AEI_i \times \frac{Rev_{ijl}}{\sum_j Rev_{ijl}}$   $\forall j, i, \ \forall l = irrigated$ (S19)

157 $PriorArea_{ijl} = \left( CroplandExtent_i - AEI_i - PriorArea_{ij,subsistence} \right) \times \frac{Rev_{ijl}}{\sum_l \sum_j Rev_{ijl}}$   $\forall j, i, \ \forall l =$

158 $rainfed$ (S20)

159

160 In the case of subsistence farming, the revenue measure is replaced by a measure of population

161 density. The subsistence part of the sub-national crop area is then pre-allocated using rural

162 population density as a weight.

163

164 $PriorArea_{ij,subsistence} = SubCropArea_{jk} \times Percent_{jl} \times \frac{Pop_i}{\sum_{i \in k} Pop_i}$   $\forall j, i, l=$subsistence (S21)

165

166 After this pre-allocation, the prior is calculated by normalizing the allocated areas over the whole

167 allocation unit:

168

169 $\pi_{ijl} = \frac{PriorArea_{ijl}}{\sum_i PriorArea_{ijl}}$   $\forall i, j, l$ (S22)

170

171

172 <u>Yield Calculation</u>

173     The calculation of yield is based on the statistical yield information for crop $j$ within production

174 system $l$ for each SRU $k$. First, the average potential yield, $\bar{Y}_{jl}$, is calculated as:

175

176 $\bar{Y}_{jlk} = \frac{\sum_i Suitability_{ijl} \times A_{ijl}}{\sum_i A_{ijl}}$   $\forall k$ (S23)

177 $Y_{ijl} = \frac{Suitability_{ijl} \times CropYield_{jlk}}{\bar{Y}_{jlk}}$ (S24)

178

179  Where $CropYield_{jlk}$ is the statistical yield reported for crop $j$ in production system $l$ within SRU

180  $k$. Then the production of crop $j$ in production system $l$, and pixel $i$, $Prod_{ijl}$, could be calculated

181  as the following:

182

183  $Prod_{ijl} = (A_{ijl} \times CropingIntensity_{jl}) \times Y_{ijl}$                                    (S25)

184

185  **S1.4 GAEZ distribution of harvested area and yield**:

186        The GAEZ model distributes available statistical data using an iterative rebalancing

187  procedure that converges to the same answer as the cross entropy approach used by SPAM

188  (Fischer et al., 2006). The GAEZ formulation includes distance to market, population density,

189  ruminant livestock density, farming system zone and producer price by crop as a means of

190  further disaggregating available national or sub-national statistics.

191

192  The iterative rebalancing algorithm used by the model is documented in detail in Fischer et al.

193  (2006), but generally works by using multipliers (separated into rainfed $R$ and irrigated $I$) for

194  area ($\lambda_j^R$ and $\lambda_j^I$) by crop $j$, for cropping intensity ($\rho^R$ and $\rho^I$), and yield ($\mu_j^R$ and $\mu_j^I$) by crop $j$.

195  The algorithm updates the multipliers iteratively such that all constraints are met, and in the

196  process produces grid-cell specific allocations of harvested area and production for rain-fed and

197  irrigated land. $\rho^R$ and $\rho^I$ provide a measure of the discrepancy between the potential for multi-

198  cropping and actual cropping intensity, while $\mu_j^R$ and $\mu_j^I$ represent the gap between actual and

199  potential crop yields.

200

201  The GAEZ model uses a two-step nested process within each iteration of the rebalancing

202  algorithm, by which land is broadly allocated into two sets of crops: Set $I_1$, for which the spatial

203  distribution layer ($\varepsilon_{ij}$) exists for pixels $i$ and crops $j$, and Set $I_2$, for which the spatial distribution

204  layer does not exist. $\varepsilon_{ij}$ is defined as a subset of the M3 dataset for selected crops in countries for

205  which more than 50% of the data was derived from sub-national statistics. Shares of land are

206  distributed to each set of irrigated ($SetShare_1^I$ and $SetShare_2^I$) and rainfed ($SetShare_1^R$ and

207  $SetShare_2^R$) crops as follows:

208

$$SetShare_1^I = \frac{\sum_{j \in I_1^I} CroppingIntensity_i^I \times Yield_{ij}^I \times Price_j \times \lambda_j^I \times RelativeYield_{ij}^I}{\sum_{j \in I_1^I \cup I_2^I} CroppingIntensity_i^I \times Yield_{ij}^I \times Price_j \times \lambda_j^I \times RelativeYield_{ij}^I} \quad (S26)$$

209 $SetShare_2^I = 1 - SetShare_1^I$ (S27)

210

211 $I_1^I = \{(j \in I_1) \wedge (\varepsilon_{ij} > 0) \wedge (RelativeYield_{ij}^I \geq \gamma_j^I)\}$ (S28)

212 $I_2^I = \{(j \in I_1) \wedge (RelativeYield_{ij}^I \geq \gamma_j^I)\}$ (S29)

213

214 Where $RelativeYield_{ij}^I$ is defined in equation (S60), and $\gamma_j^I$ is the crop allocation relative yield

215 threshold for irrigated crop $j$. Similarly for rainfed crops:

216

$$SetShare_1^R = \frac{\sum_{j \in I_1^R} CroppingIntensity_i^R \times Yield_{ij}^R \times Price_j \times \lambda_j^R \times RelativeYield_{ij}^R}{\sum_{j \in I_1^R \cup I_2^R} CroppingIntensity_i^R \times Yield_{ij}^R \times Price_j \times \lambda_j^R \times RelativeYield_{ij}^R} \quad (S30)$$

217 $SetShare_2^R = 1 - SetShare_1^R$ (S31)

218 $I_1^R = \{(j \in I_1) \wedge (\varepsilon_{ij} > 0) \wedge (RelativeYield_{ij}^R \geq \gamma_j^R)\}$ (S32)

219 $I_2^R = \{(j \in I_1) \wedge (RelativeYield_{ij}^R \geq \gamma_j^R)\}$ (S33)

220

221 Where $RelativeYield_{ij}^R$ is defined in equation (A61), $\gamma_j^R$ is the crop allocation relative yield

222 threshold for rainfed crop $j$.

223

224 In the second step of the process, rainfed crop-specific area shares ($Share_{ij}^R$) and irrigated crop-

225 specific physical area shares ($Share_{ij}^I$) are calculated for each grid cell $i$ and crop $j$ in the set of

226 grid cells within a specific SRU $k$ as follows:

227

228 $Share_{ij}^I = SetShare_1^I \times \frac{(\varepsilon_{ij}^I \times \lambda_j^I)}{\sum_{k \in I_1^I}(\varepsilon_{ik}^I \times \lambda_k^I)} , \ j \in I_1^I$ (S34)

229 $Share_{ij}^I = 0 , \ j \in I_1 \ \wedge \ j \notin I_1^I$ (S35)

230 $Share_{ij}^I = SetShare_2^I \times \frac{CroppingIntensity_{ij}^I \times Yield_{ij}^I \times Price_j \times \lambda_j^I \times RelativeYield_{ij}^I}{\sum_{k \in I_2^I} CroppingIntensity_{ik}^I \times Yield_{ik}^I \times Price_k \times \lambda_k^I \times RelativeYield_{ik}^I} , j \in I_2^I$ (S36)

8

231 $Share_{ij}^I = 0$ , $j \in I_2 \wedge j \notin I_2^I$ (S37)

232

233 And similarly for rainfed crop areas:

234 $Share_{ij}^R = SetArea_1^R \times \frac{(\varepsilon_{ij}^R \times \lambda_j^R)}{\sum_{k \in I_1^R}(\varepsilon_{ik}^R \times \lambda_k^R)}$ , $j \in I_1^R$ (S38)

235 $Share_{ij}^R = 0$ , $j \in I_1 \wedge j \notin I_1^R$ (S39)

236 $Share_{ij}^R = SetShare_2^R \times \frac{CroppingIntensity_{ij}^R \times Yield_{ij}^R \times Price_j \times \lambda_j^R \times RelativeYield_{ij}^R}{\sum_{k \in I_2^R} CroppingIntensity_{ik}^R \times Yield_{ik}^R \times Price_k \times \lambda_k^R \times RelativeYield_{ik}^R}$ , $j \in I_2^R$ (S40)

237 $Share_{ij}^R = 0$ , $j \in I_2 \wedge j \notin I_2^R$ (S41)

238

239 The crop-specific area shares are then used to calculate irrigated harvested areas ($HA_{ij}^I$) and

240 rainfed harvested areas ($HA_{ij}^R$) for each pixel $i$ and crop $j$ in the set of grid cells within a specific

241 SRU $k$ as follows:

242

243 $HA_{ij}^I = fCroplandExtent_i^I \times GridCellArea_i \times \left( \frac{\rho^I \sum_{k \in crops} Share_{ij}^I \times CroppingIntensity_{ij}^I}{\sum_{k \in crops} Share_{ij}^I} \right)$ (S42)

244 $HA_{ij}^R = fCroplandExtent_i^R \times GridCellArea_i \times \left( \frac{\rho^R \sum_{k \in crops} Share_{ij}^R \times CroppingIntensity_{ij}^R}{\sum_{k \in crops} Share_{ij}^R} \right)$ (S43)

245

246 The constraints in the model are as follows:

247

248 Grid cell proportion of rainfed and irrigated land:

249 $fCroplandExtent_i^R = fCroplandExtent_i^T - fCroplandExtent_i^I$, $\forall i$ (S44)

250

251 Where $fCroplandExtent_i^T$ is the proportion of total cropland in grid cell $i$ (IIASA dataset,

252 developed as part of GAEZ) and $fCroplandExtent_i^I$ is the proportion of each cropland grid cell

253 that is equipped for irrigation (Seibert et al., 2007) dataset.

254

255 Cropland extent by grid cell:

256 $CroplandExtent_i = fCroplandExtent_i^T \times GridCellArea_i$, $\forall i$ (S45)

257

258      Where $fCroplandExtent_i^T$ is the total share of cropland (irrigated and rainfed) in each pixel

259

260      Grid cell cropping intensity for annual crops:

261      $CroppingIntensity_i^I = \rho^I \times CroppingFactor_i^I$          (S46)

262      $CroppingIntensity_i^R = \rho^R \times CroppingFactor_i^R$          (S47)

263

264      Where $CroppingFactor_i^I$ and $CroppingFactor_i^R$ correspond to the cultivation intensity class

265      factor of irrigated and rainfed annual crops, respectively.

266

267      Total irrigated harvested area by crops:

268      $HA_j^I = \propto_j^I \times HA_j, \quad \forall j$          (S48)

269      $HA_j^R = \left(1 - \propto_j^I\right) \times HA_j, \quad \forall j$          (S49)

270

271      Where $\propto_j^I$ is the proportion of harvested area that is irrigated for crop $j$.

272

273      Harvested area and cropland extent (irrigated and rainfed) by grid cell:

274      $HA_i^I = CroppingIntensity_i^I \times fCroplandExtent_i^I \times GridCellArea_i, \quad \forall i$      (S50)

275      $HA_i^R = CroppingIntensity_i^R \times fCroplandExtent_i^R \times GridCellArea_i, \quad \forall i$      (S51)

276

277      Harvested area by pixel and crop:

278      $HA_{ij}^I = CroppingIntensity_i^I \times Share_{ij}^I \times fCroplandExtent_i^I \times GridCellArea_i, \forall i, j \in annual\ crops$    (S52)

279      $HA_{ij}^I = CroppingIntensity^P \times Share_{ij}^I \times fCroplandExtent_i^I \times GridCellArea_i, \forall i, j \in perennial\ crops$ (S53)

280      $HA_{ij}^R = CroppingIntensity_i^R \times Share_{ij}^R \times fCroplandExtent_i^R \times GridCellArea_i, \quad \forall i, j \in annual\ crops$    (S54)

281      $HA_{ij}^R = CroppingIntensity^P \times Share_{ij}^R \times fCroplandExtent_i^R \times GridCellArea_i, \forall i, j \in perennial\ crops$ (S55)

282

283      Where $CroppingIntensity^P$ is the cropping intensity of perennial crops.

284      $HA_j^I = \sum_i HA_{ij}^I, \quad \forall j$          (S56)

285      $HA_j^R = \sum_i HA_{ij}^R, \quad \forall j$          (S57)

286

287      Grid cell irrigated yield:

288     $Yield_{ij}^{I} = \mu_j^I \times PotentialYield_{ij}^{I,high}, \ \forall i,j$                                                (S58)

289     $Yield_{ij}^{R} = \mu_j^R \times \left( (1 - \psi_{ij}^R) \times PotentialYield_{ij}^{R,low} + \left( \psi_{ij}^R \times PotentialYield_{ij}^{R,high} \right) \right), \ \forall i,j$     (S59)

290

291     Where $PotentialYield_{ij}^{I,high}$, $PotentialYield_{ij}^{R,high}$ and $PotentialYield_{ij}^{R,low}$ are the potential

292     yield of crop j on grid cell i in a high input irrigated system, high input rainfed system and low

293     input rainfed system, respectively. Potential yield information is derived from the GAEZ v3.0

294     database (Fischer et al., 2013). $\psi_{ij}^R$ is a spatial location factor used to reflect differences in

295     management intensity and input use derived from remote sensing, household survey data,

296     information on farm size or market orientation of a household.

297

298     Grid cell relative yield factor:

299     $RelativeYield_{ij}^{I} = \left. PotentialYield_{ij}^{I,high} \middle/ \max_{k \in gridcell} \left( PotentialYield_{kj}^{I,high} \right) \right.$                   (S60)

300     $RelativeYield_{ij}^{R} = \left. PotentialYield_{ij}^{R,high} \middle/ \max_{k \in gridcell} \left( PotentialYield_{kj}^{R,high} \right) \right.$                   (S61)

301

302     Further information on the formulation used in the model may be found in the GAEZ v.3.0

303     documentation main text with details in Appendix A8 (Fischer et al., 2013). Details on the

304     iterative rebalancing algorithm may be found in Fischer et al. (2006).

**Appendix S2: Comparison of the Downscaling Methodologies**

*Cropland extent delineation*

305

306

307      As a first step towards delineating crop specific harvested area and yield, each cropping

308      system model defined a spatially explicit layer of cropland extent, representing the proportion of

309      cropland in each 5-minute pixel globally. Because each subsequent step in the modeling process

310      relies on the definition of cropland extent, the degree to which each pair of cropland extent

311      products agree represents an upper bound of inter-model agreement on the spatial distribution of

312      crop physical areas.

313      M3, MIRCA and SPAM all rely on the same base dataset for cropland extent:

314      Ramankutty et al., (2008), which is an extension of Leff et al., (2004). Leff et al., (2004)

315      synthesize satellite-derived land cover data and agricultural census data worldwide to assess the

316      distribution of major crops across a global 5 arc minute grid in terms of the proportion of the

317      total harvested area of each of the crops in each administrative unit. Following and improving on

318      this work, Ramankutty et al. (2008) developed a new global land cover data set for croplands and

319      pasture circa 2000 (at the same 5 arc minute resolution of the original dataset) by combining

320      Boston University's MODIS-derived land cover data (Friedl et al., 2002) and SPOT

321      VEGETATION based GLC2000 (Bartholome and Belward, 2005). Ramankutty et al. (2008)

322      apply a multiple linear regression model to relate the combined satellite derived datasets to the

323      agricultural statistics using a least squares error framework. The optimization is applied

324      separately to six different regions of the world.

325      The cropland extent developed by Ramankutty et al. (2008) is used directly by M3 and

326      with modifications by MIRCA and SPAM. By combining Ramankutty et al. (2008) and the

327      global map of irrigation areas (GMIA), MIRCA produced a global dataset of monthly growing

328    areas of 26 irrigated crops on the same 5 arc minutes grid. SPAM similarly reconciles the GMIA

329    map of irrigated areas and Ramunkutty cropland extent by setting the cropland extent to be at

330    least equal to the irrigated area in a preprocessing step. For more information on each of these

331    methodologies, see Appendix S1.

332        GAEZ uses GLC2000 data and GMIA, but also considers a global land cover

333    categorization (IFPRI, 2002), which is based on a reinterpretation of the Global Land Cover

334    Characteristics Database v.2.0 (EROS Data Center, 2000), a layer of forest land from the Forest

335    Resources Assessment of FAO (FAO, 2001), the IUCN-WCMC protected areas inventory

336    (WPDA, 2009) and an estimate of land required for housing and infrastructure for the year 2000

337    derived from FAO-SDRN, based on LANDSCAN 2003 that were calibrated to UN 2000

338    population figures (Fischer et al., 2008; Bhaduri et al., 2002; Dobson et al., 2000). GAEZ runs a

339    cross-sectional regression on the land cover distributions to derive weights, which are then

340    applied in an iterative adjustment procedure to match estimated reference values such that the

341    geographic and statistical data are consistent.

342

343    *Suitability Constraints*

344        GAEZ and SPAM further constrain potential crop distribution using biophysical and

345    socioeconomic suitability prior to allocating the harvested area of each crop. M3 and MIRCA do

346    not consider suitability criteria. SPAM directly uses the suitable area product from GAEZ,

347    meaning that despite using different cropland extent products to constrain the distribution of

348    crops, the two models use identical constraints on biophysically suitable land. The GAEZ

349    suitability product integrates an extensive set of edaphic and climatic factors into its biophysical

350    suitability analysis to produce a potential yield estimate and a suitability index by production

351  system and crop. Further information on the suitability index analysis developed as part of the

352  GAEZ model may be found in Fischer et al. (2013).

353  In addition to biophysical suitability criteria, both SPAM and GAEZ consider the

354  socioeconomic factors that often constrain or encourage crop production. As a means of

355  differentiating between low, medium and high input or management conditions, GAEZ divides

356  the land into land use types. Land use types are derived using information on road infrastructure,

357  livestock density, population density and distance to market. Low input, for example, relies on

358  available human/livestock labor while high input is market oriented, using improved varieties,

359  fertilizer, pesticides and machinery. Similar to GAEZ, SPAM explicitly models different

360  production systems, which include high-input irrigated, high-input rainfed, low –input rainfed

361  and subsistence (always low-input rainfed). SPAM includes data on crop prices and market

362  access to construct a realistic market scenario in which there are not only biophysical barriers to

363  producing crops, but social economic forces as well (see Appendix S1 for an explicit

364  mathematical formulation).

365

366  *Distribution of harvested area and yield*

367  Perhaps the largest methodological differences between M3, MIRCA, SPAM and GAEZ

368  are in the approaches used to downscale statistical data reported at administrative unit level into

369  gridcell specific values. M3 uses the most straightforward method, allocating each crop evenly

370  across potential cropland in each statistical reporting unit as the proportion of harvested area

371  occupied by the crop to total harvested area in that reporting unit. Crop yield in each grid cell is

372  assigned as being the same as the yield reported for the statistical unit as a whole. This approach

373  implicitly assumes both environmental conditions and management/production systems are

374 uniform across the cropland extents of each statistical reporting unit, or that there is insufficient

375 information to characterize the spatial variations of crop production within a statistical unit. As a

376 result, the distinct tolerances of individual crops to those spatial patterns are not incorporated in

377 the downscaling procedure. This approach does not, furthermore, acknowledge the very

378 significant differences between the yield levels of irrigated and rainfed production systems, nor

379 of commercial and smallholder producers within these sometimes large and highly diverse

380 statistical reporting units.

381       MIRCA primarily focuses on reconciling the differences between information derived

382 from sub-national crop production statistics, M3 crop distributions and the Siebert et al. (2005)

383 irrigated areas database. MIRCA deals only with harvested area and essentially uses the relative

384 share of rainfed and irrigated cropland within each grid cell to break out M3 total crop areas into

385 gridcell-specific rainfed and irrigated areas. The MIRCA model derives cropping intensity,

386 which is used to convert harvested to physical area. MIRCA also includes use of numerous

387 checks and adjustments to reconcile differences between the cropland area of Ramankutty et al.

388 (2008) and the irrigated area estimates of Siebert et al. (2005) within each gridcell, given that

389 total cropland area should at all times be greater than or equal to the irrigated cropland area.

390 Similar to M3, MIRCA does not consider any form of suitability in its downscaling procedure.

391       The downscaling approaches of GAEZ and SPAM are predicated on the importance of

392 attempting to take explicit account of available evidence of the spatial variation of production

393 conditions within the cropland extent and of the significantly different yield resulting from each

394 of those systems. Both GAEZ and SPAM use an approach that produces a result mathematically

395 equivalent to that of a cross-entropy formulation, but GAEZ uses an iterative rebalancing

396 procedure to adjust weighting factors until all constraints in the model are met, while SPAM uses

397    a cross-entropy formulation. SPAM uses explicit cropping intensity from statistics and expert

398    opinion to calculate cropping intensity while GAEZ derives a cropping intensity factor through

399    the rebalancing procedure (see Appendix S1). Although the two models incorporate similar

400    information (see Table 2, main text), the manner in which the information is used to constrain the

401    model differs (see Appendix S1 for details on the mathematical formulation of each model).

402    Additionally, GAEZ differs from SPAM in that it uses a "location factor" to incorporate spatially

403    explicit information including geo-referenced household survey data. The prior in SPAM is used

404    to capture spatially explicit information as well but the model does not include household survey

405    data, instead leveraging the field presence of the CGIAR network to incorporate an extensive

406    dataset of expert elicitations.

407

408    **Input Data and Model Interdependencies**

409    The major determinants of the potential reliability of downscaling efforts are (a) the quality

410    of the cropland extent dataset indicating the physical extent and area intensity of cropland (e.g.,

411    share of cropland area in each 5 arc minute grid cell), and (b) the resolution and reliability of the

412    sub-national crop statistics. Each model builds on a common set of available data as well as

413    previous work in cropping systems modeling. Table 2 (main text) illustrates both the broad

414    linkages and increasing sets of input data and assumptions that each of the M3, MIRCA, GAEZ

415    and SPAM datasets relies upon.

416

417    *National and Sub-National Statistics*

418    All four datasets draw on FAOSTAT national data to provide control totals for cropland area,

419    the harvested area, and yields of specific crops, while also spending considerable efforts to

420   collect sub-national crop statistics to allow as detailed as possible disaggregation of national

421   totals within sub-national administrative boundaries. Since MIRCA relies on M3 to provide its

422   input data on the spatial allocation of the total area and average yield, it relies initially on the

423   same sources of subnational crop statistics. The GAEZ model uses data from FAOSTAT as a

424   constraint at the national level and – similar to MIRCA - uses the M3 sub-national statistics for

425   select crops in countries that have sub-national statistics covering more than 50% of the country.

426   SPAM relies on a separate collection of sub-national statistical data sources, focusing on

427   increased coverage in developing countries.

428       M3 reports a total of 22,106 statistical reporting units globally, of which 56 were national,

429   2,299 were first level sub-national disaggregation (e.g., US state level), and 19,751 were second

430   level (e.g. US county level) reporting units. SPAM reports 24,507 statistical units of which 251

431   were national, 2,758 were first level, and 21,498 were second level. SPAM focused its data

432   collection efforts particularly in developing countries. For example in Africa the M3 and SPAM

433   data sets were developed using around 300 and 4,150 second-level statistical reporting units

434   respectively.

435

436   *Extent of Irrigation*

437       Those models that distinguish between rainfed and irrigated cultivations (MIRCA, SPAM

438   and GAEZ) all use the GMIA v4.0, released in 2007 (Siebert et al., 2005), to identify the location

439   and area intensity of irrigated production. However, the MIRCA and SPAM teams compiled

440   information in the national and sub-national shares of different production systems and cropping

441   intensities independently. MIRCA and SPAM both draw on FAO's AQUASTAT and national

442   databases for gaining greater insights into national and crop-specific irrigation extents and

443    practices, but MIRCA relies on a richer collection of national data, including a more complete

444    collection of national/sub-national crop calendars and cropping intensities (in part because the

445    goal of MIRCA is to produce monthly and not annual crop distribution maps). In contrast to

446    MIRCA and SPAM, GAEZ relies on its own data for information about cultivation intensity of

447    irrigated crops.

448

449    *Ancillary Data*

450    SPAM and GAEZ incorporate datasets beyond those used by M3 and MIRCA as a means of

451    differentiating between production levels within cropping systems. The SPAM approach requires

452    additional sets of data because it attempts further disaggregation of its rainfed production

453    statistics amongst commercial and subsistence categories, and bases its approach to distribution

454    of individual crops within the cropland extent on agronomic, economic and demographic

455    principles and assumptions. These  include crop area and production shares amongst irrigated

456    production and large-scale/commercial and smallholder rainfed production, the spatial

457    differences in the biophysical suitability of individual crops for irrigated and rainfed (commercial

458    and subsistence) production, and estimates of the spatial patterns of population density as well as

459    crop prices. GAEZ similarly divides the land into land use types to reflect variable management

460    and input conditions. Data used to differentiate among land use types reflect the specific

461    requirements of each and include road infrastructure, livestock density, population density and

462    distance to market. Table 2 reflects the overlapping and separate ancillary datasets used by

463    SPAM and GAEZ. In addition to available ancillary datasets, SPAM leverages the international

464    network and field presence of CGIAR to undergo a systemic validation process. The feedback

465  from this validation is used to inform future model simulations. This process is unique to the

466  SPAM model.

467

468  **Appendix S3: Gaussian Filter Sensitivity Analysis**

469

470  The 2-dimensional Gaussian filter for pixel $i$ may be expressed as:

471  $$g(x,y)_i = \left(\frac{1}{2\pi\sigma^2}\right)\left(e^{-\frac{x^2+y^2}{2\sigma^2}}\right) \quad\quad\quad (S1)$$

472  Where $x$ is the distance from the horizontal axis, $y$ is the distance from the vertical axis and $\sigma$ is the

473  standard deviation of the Gaussian distribution, used to control the kernel density as illustrated below

474

475  Figure S1: Differences in Cropland extent with Gaussian filters having kernel densities of 0 (pixel-level

476  comparison), 1 (4 pixel radius), 2 (8 pixel radius), 3 (12 pixel radius) and 4 (16 pixel radius).



| Pixel Size<br>No Filter<br>(1 pixel; 0.5 deg.) | Gaussian Filter<br>1σ kernel den.<br>(4 pixels; 2 deg.) | Gaussian Filter<br>2 σ kernel den.<br>(8 pixels; 4 deg.) | Gaussian Filter<br>3 σ kernel den.<br>(12 pixels; 6 deg.) | Gaussian Filter<br>4σ kernel den.<br>(16 pixels; 8 deg.) |

477

478

479

480   Figure S2: Differences in Cropland extent with Gaussian filters having kernel densities of 0 (pixel-level

481   comparison), 1 (4 pixel radius), 2 (8 pixel radius), 3 (12 pixel radius) and 4 (16 pixel radius).



482

483

**Appendix S4: Cropland extent supplementary material**

486   Figure S1: Cropland extent of A) GAEZ and B) Ramankutty et al., (2008)



487

488     Figure S2: Pixel-wise cropland extent differences



489
490

491

492

493     *Figure S3: GAEZ and Ramankutty cropland Extent by Latitude*

494 **Appendix S5: Supplementary Figures for Wheat**

495 Figure S1: Model agreement on magnitude of harvested area of wheat by threshold. A) Harvested Area >
496 0% of cell, B) Harvested Area > 1% of cell, C) Harvested Area > 10% of cell, D) Harvested Area > 25%
497 of cell



498
499

**Appendix S6: Supplementary Analysis for Rice and Maize**

*S6.1 Rice harvested areas and yields*

Although there are minor departures from model consensus at higher thresholds, the majority

of the discrepancies in the spatial distribution of rice are at the lowest threshold (see

supplementary Figure S1 panel A). As discussed with wheat, these dissimilarities arise due to

differences in model methodology. But in contrast to wheat, the models in disagreement are

primarily the SPAM and MIRCA models (see supplementary Fig. S5 in Appendix S6.2), not

only those that spread harvested area across plausible cropland, implying that the inconsistencies

are in part attributable to the collection of supplementary subnational statistics as described in

Section 4.1. At higher thresholds for harvested area the disagreement between products is

minimal and will be explored further in the following analyses.

The harvested areas of rice appear to be less influenced by discrepancies in the cropland

extent products, and instead differ as a function of downscaling method or input data. Evaluating

the harvested areas of rice by latitude reveals that MIRCA predicts significantly more harvested

area for rice north of 30N than do the other products (see Fig. S2). This divergence may stem

from the fact that the MIRCA method of crop distribution does not consider biophysical

limitations, or it may reflect differences in the input statistical crop yields at a sub-national level

given that the above average estimation by MIRCA appears to be concentrated in eastern China

(see Fig. S3).

With the exception of MIRCAs large harvested area in east China, the relation between

GAEZ and each of the products that use Ramankutty cropland extent (M3, MIRCA and SPAM)

is nearly identical (see Fig. S3). This may indicate that all three use similar sub-national rice data

in China. M3, MIRCA and SPAM differ from GAEZ, for example, in their distribution of rice

523    within India: GAEZ distributes more rice area to the southwest while other products distribute

524    more rice to the northeast.

525         Maps of rice yields match even less well than did the wheat yields, differing significantly

526    over all latitudes north of the equator (see Fig. S2). M3 predicts yields higher than SPAM and

527    GAEZ over most latitudes but particularly in Asia (see Figs. S2 and S4). This can be seen in the

528    skew of the distributions of the histograms in each pair-wise yield comparison. The differences

529    between GAEZ and SPAM are most pronounced in China. SPAM predicts consistently lower

530    yields than both M3 and GAEZ over nearly all of China.

531         Figure S5 illustrates the model-dependent differences, and resulting uncertainty, in

532    calculating the yield gap (panel A) using both an absolute measure (panel C) and with regard to

533    existing yields (panel B). Areas in which the yield gap uncertainty ratio approaches 1 signify

534    areas in which uncertainty dominates the estimate of the yield gap. However, it is equally

535    important to contextualize these uncertainties relative to existing yields and using an absolute

536    measure of model difference. Areas displaying large values in all three panels indicate areas in

537    which the model-estimated yield gaps disagree (panel C), where this disagreement is a

538    significant proportion of the estimated yield gap (panel A) and in which the differences are

539    important in the context of existing food production systems (panel B).

540         The model dependent uncertainty exceeds the estimated yield gap in significant parts of

541    every rice-growing continent, meaning that uncertainty in the estimation of yields dominates the

542    yield gap calculation in these regions. As with wheat, the differences relative to existing yields

543    are decreased in some major producing areas, but remain significant in others. Of particular

544    interest are the Ganges basin in India, Peru, parts of China and Indonesia, which display high

545    values across all three indices.

546

547

548  *S6.2 Maize harvested area and yield*

549  Model comparisons for the harvested area of maize should be interpreted with some care as

550  not all models measure identical quantities. M3, GAEZ and SPAM all measure maize as it is

551  grown for grain only, while MIRCA measures maize as the sum of maize grown for silage and

552  for grain. While there was insufficient model data to separate (or aggregate) the models

553  quantities to align identically, a comparison between models is still useful for broadly identifying

554  inter-model differences. For the yield analysis all models (M3, SPAM and GAEZ) did measure

555  the same quantity.

556  The consensus analysis reveals that as compared with rice and wheat, the agreement between

557  models on the extent of maize harvested area is generally good with a few notable exceptions. At

558  the lowest threshold of harvested area, MIRCA is alone in designating growing area in much of

559  Canada, Japan, the UK and Ireland while SPAM is the only model that indicates additional area

560  in central Asia (see supplementary Fig. S6 panel A and supplementary Fig. S11 in Appendix

561  S7.3). The SPAM differences are likely due to collection of additional national or sub-national

562  statistics, while the MIRCA differences may arise due to a combination of the inclusion of silage

563  and additional data collection. At the threshold of > 1% (panel B), model methodology

564  dominates the differences: MIRCA indicates higher proportions of maize harvested area due to a

565  more even distribution across statistical reporting units. In southeast China all models agree on

566  the general extent of maize harvested area but disagree trivially on the spatial placement as

567  evidenced by comparison with the results of the Gaussian analysis in this region.  All models

568 show a relative consensus on spatial extent of maize harvested area that covers >10% and > 25%

569 of the cell (see panels C and D of S6).

570  As with the harvested area of wheat, the broad patterns of disagreement in the harvested

571 areas of maize reflect differences in the cropland extent products in many regions. This relation

572 is particularly apparent in South America, West Africa and North America (see Fig. S8). By

573 latitude, MIRCA displays a larger harvested area than any other product at 50-60N (see Fig. S7).

574 This may have to do with the lack of biophysical constraints in the distribution method or may be

575 due to the inclusion of maize used for silage, as described earlier.

576  The maize yields show consistently different patterns both spatially and by lattitude. GAEZ,

577 for example, displays consistently higher yields in the tropics (see Fig. S7), while M3 predicts

578 significantly lower yields than either SPAM or GAEZ in Eastern United States (see Fig. S9). As

579 evidenced by the skew in the histogram insets of the pair-wise comparisons, SPAM distributes

580 maize yields to be significantly larger in a smaller number of cells (insets, Fig. S9).

581  Figure S10 illustrates the uncertainty ratio (panel A), the model-dependent differences in

582 calculating the yield gap (panel C), and those differences relative to existing yields (panel C).

583 Differences in the estimated yield gap exceed 1 tonne / ha over a majority of maize producing

584 areas of the globe. As with both rice and wheat, uncertainty dominates the calculation of the

585 yield gap in significant portions of every continent (see Fig. S10). Areas displaying large values

586 in all three indices include parts of East and Southern Africa, Brazil, Mexico and Pakistan. These

587 differences highlight the importance of continued efforts towards improving model estimates of

588 harvested area and yield.

589

590    Figure S1: Model agreement on magnitude of harvested area of rice by threshold. A) Harvested Area >

591    0% of cell, B) Harvested Area > 1% of cell, C) Harvested Area > 10% of cell, D) Harvested Area > 25%

592    of cell



593

594

595

596

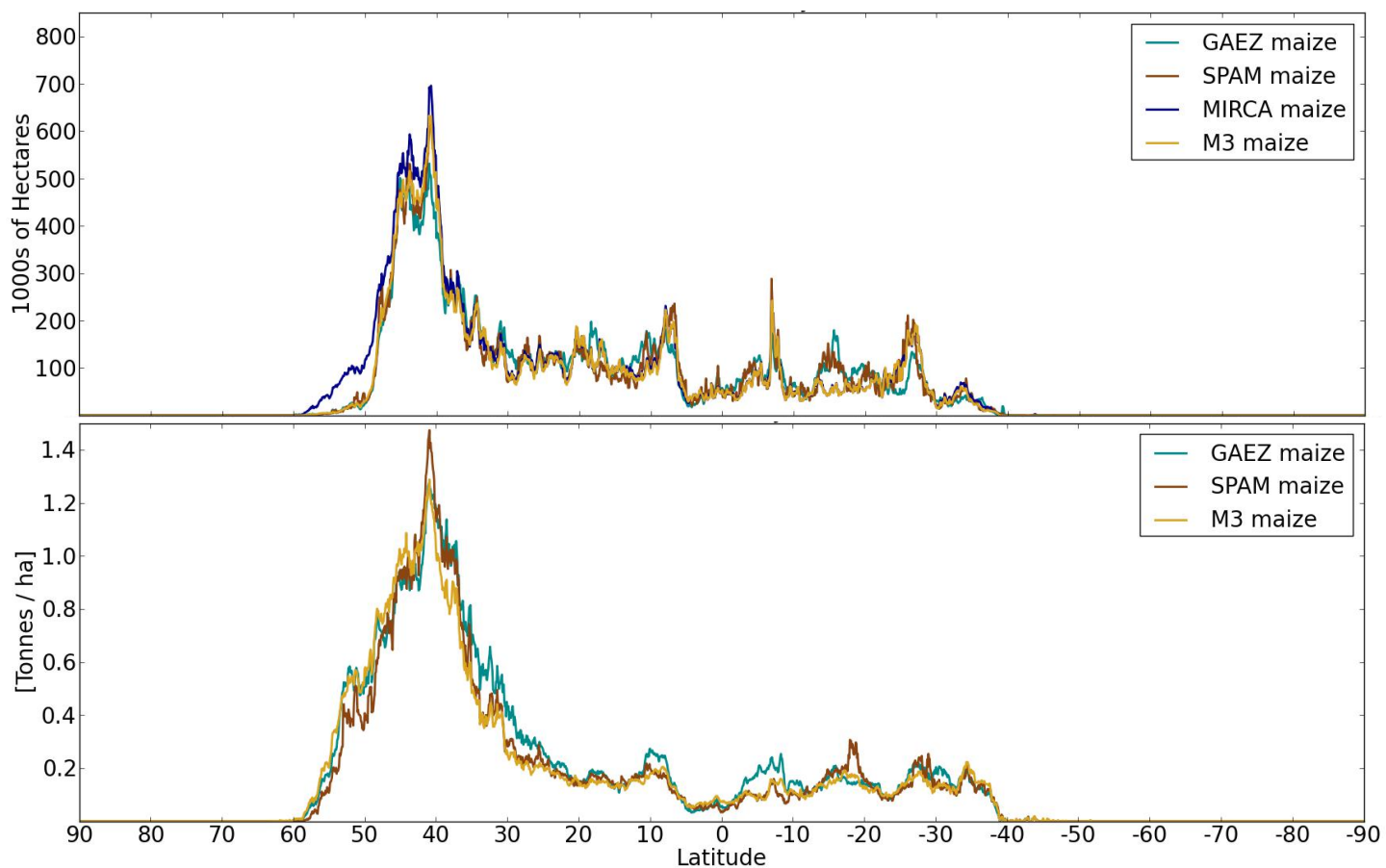*Figure S2: Rice harvested area and yield by latitude*

601  *Figure S3: Comparison of rice harvested area by model following a Gaussian filter of three sigma kernel*
602      *density. Histograms in each panel display the normalized percent of pixels as a function of*
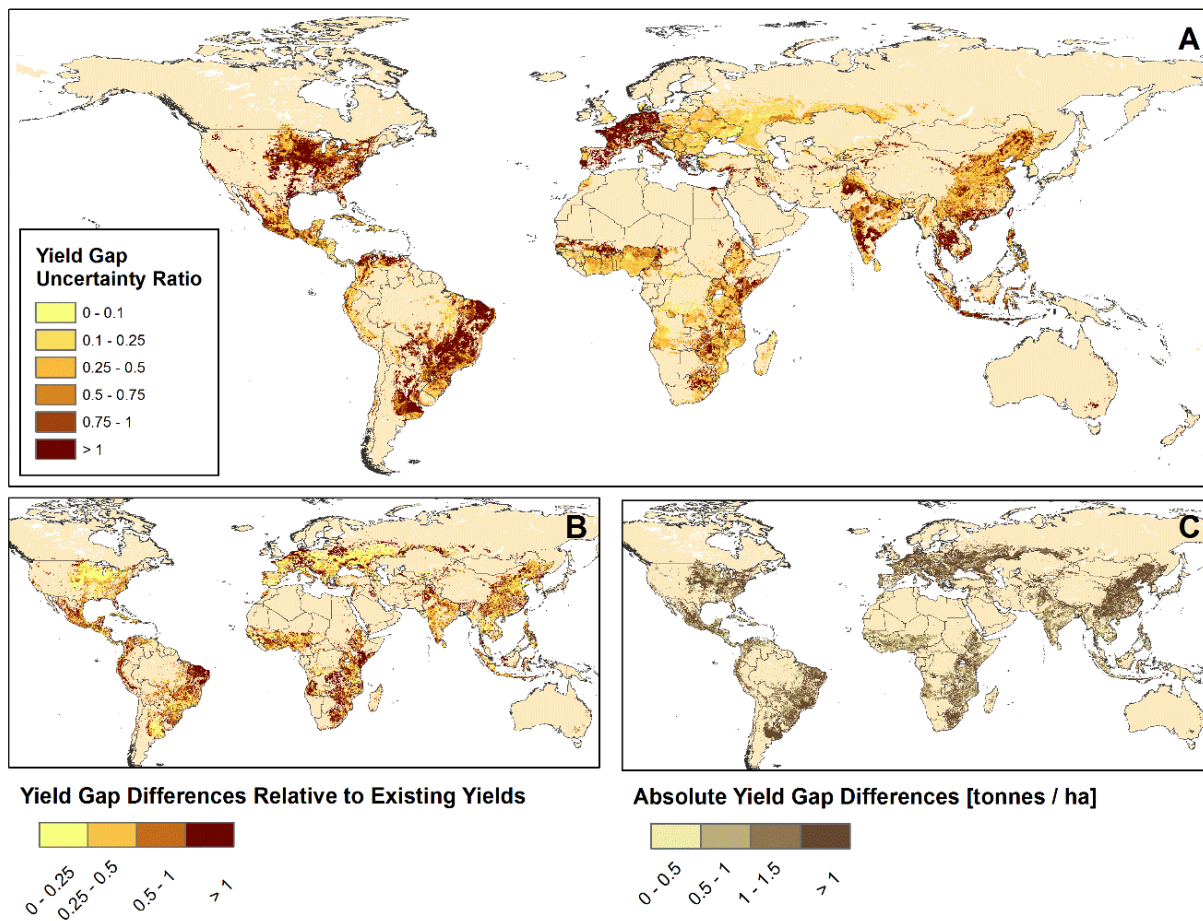603      *harvested area, y-axis limits[0, 50%], x-axis limits[-5000, 5000] ha.*
604



605
606

607    *Figure S4: Comparison of rice yield by model following a Gaussian filter of three sigma kernel density.*

608    *Histograms in each panel display the normalized percent of pixels as a function of yield, y-axis*

609    *limits[0, 35%], x-axis limits[-20, 20] tonnes/ha.*

615    *Figure S5: Implications of model differences for estimated rice yield gaps.  A) Yield Gap Uncertainty*

616    *Ratio: average model difference divided by average estimated yield gap B) average difference in*

617    *estimated yield gap divided by existing yield C) average difference in estimated yield gap*

618



619

620     Figure S6: Model agreement on magnitude of harvested area of maize by threshold. A) Harvested Area >

621     0% of cell, B) Harvested Area > 1% of cell, C) Harvested Area > 10% of cell, D) Harvested Area > 25%

622     of cell



623

624

625

626

627

*Figure S7: Maize harvested area and yield by latitude*

632 *Figure S8: Comparison of maize harvested area by model following a Gaussian filter of three sigma*

633 *kernel density. Histograms in each panel display the normalized percent of pixels as a function of*

634 *harvested area, y-axis limits[0, 35%], x-axis limits[-5000, 5000] ha.*



635

636

637

638　　*Figure S9: Comparison of maize yield by model following a Gaussian filter of three sigma kernel density.*

639　　*Histograms in each panel display the normalized percent of pixels as a function of yield, y-axis limits[0,*

640　　*35%], x-axis limits[-20, 20] tonnes/ha.*

641



642

*Figure S10: Implications of model differences for estimated maize yield gaps.  A) Yield Gap Uncertainty Ratio: average model difference divided by average estimated yield gap B) average difference in estimated yield gap divided by existing yield C) average difference in estimated yield gap*

**Appendix S7: Supplementary pixel-wise figures for wheat, rice and maize**


S7.1 Pixel-wise figures for wheat:


Figure S1: Wheat harvested area for A) M3, B) GAEZ, C) MIRCA and D) SPAM

Figure S2: pixel-wise comparison of the wheat harvested area by model

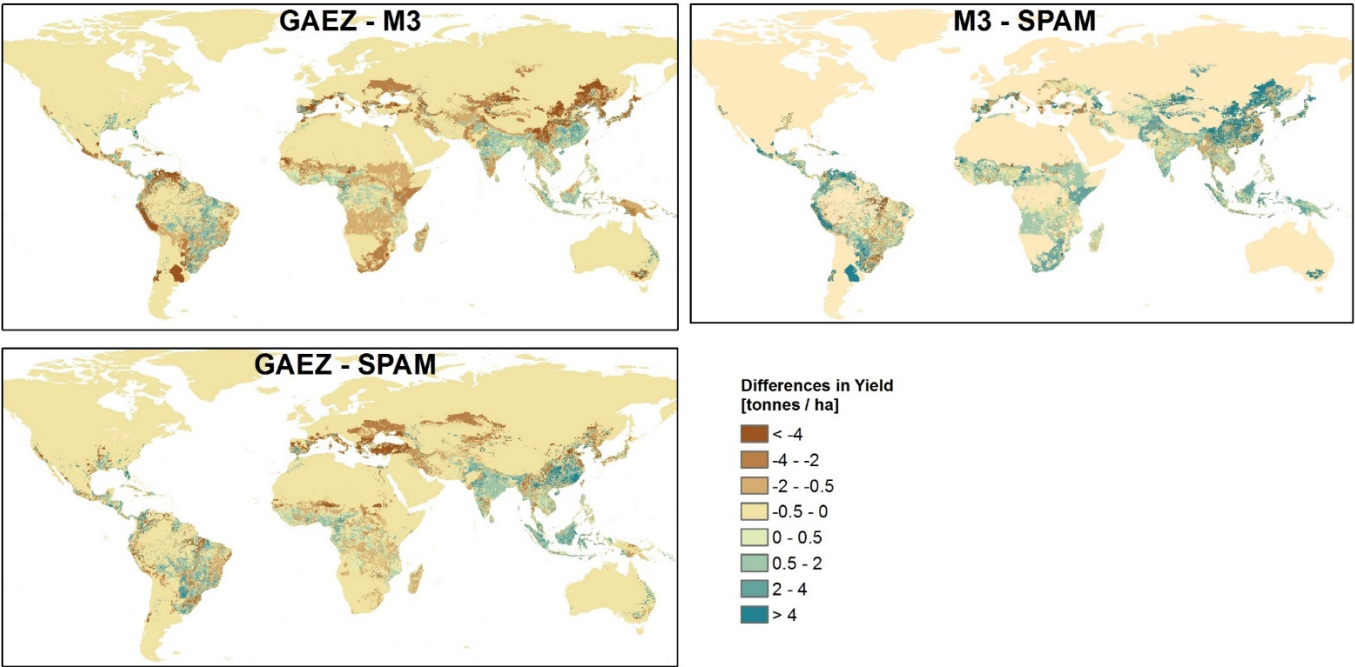Figure S3: Wheat yield for A) M3, B) GAEZ, and C) SPAM



Figure S4: pixel-wise comparison of the wheat yield by model

*S6.2 Pixel-wise figures for rice*

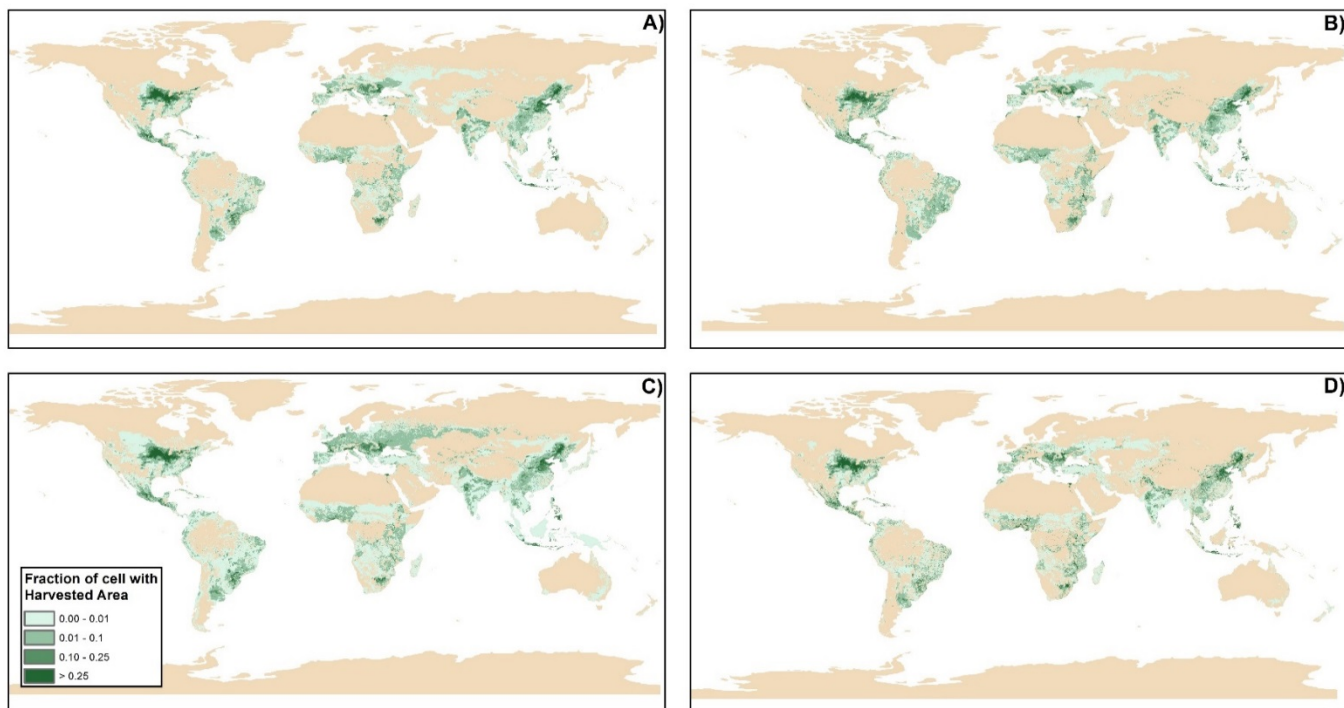Figure S5: Rice harvested area for A) M3, B) GAEZ, C) MIRCA and D) SPAM

Figure S6: pixel-wise comparison of the rice harvested area by model

Figure S7: Rice yield for A) M3, B) GAEZ, and C) SPAM



Figure S8: pixel-wise comparison of the rice yield by model

S7.3 Pixel-wise figures for maize

Figure S9: Maize harvested area for A) M3, B) GAEZ, C) MIRCA and D) SPAM

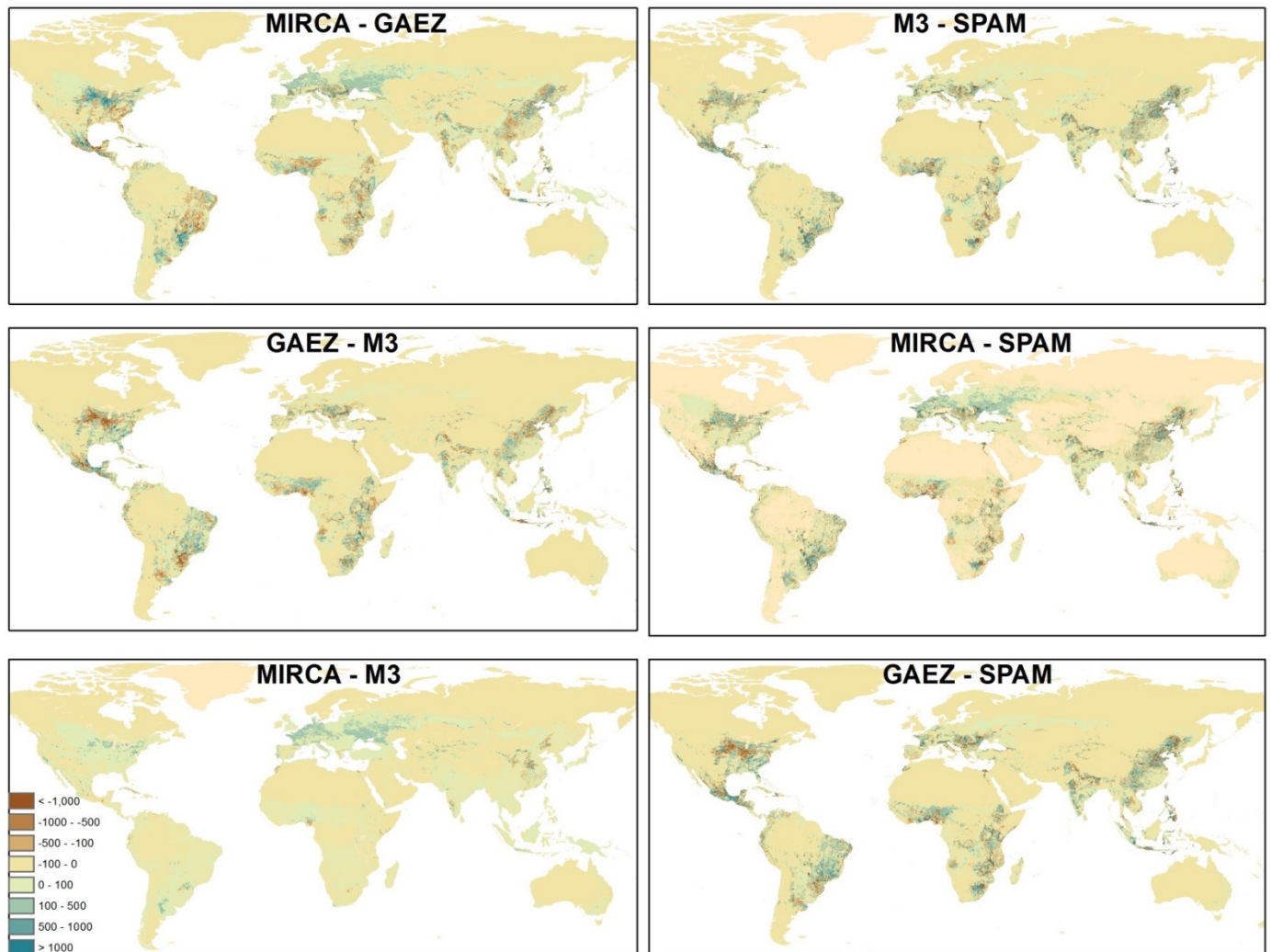Figure S10: pixel-wise comparison of the maize harvested area by model

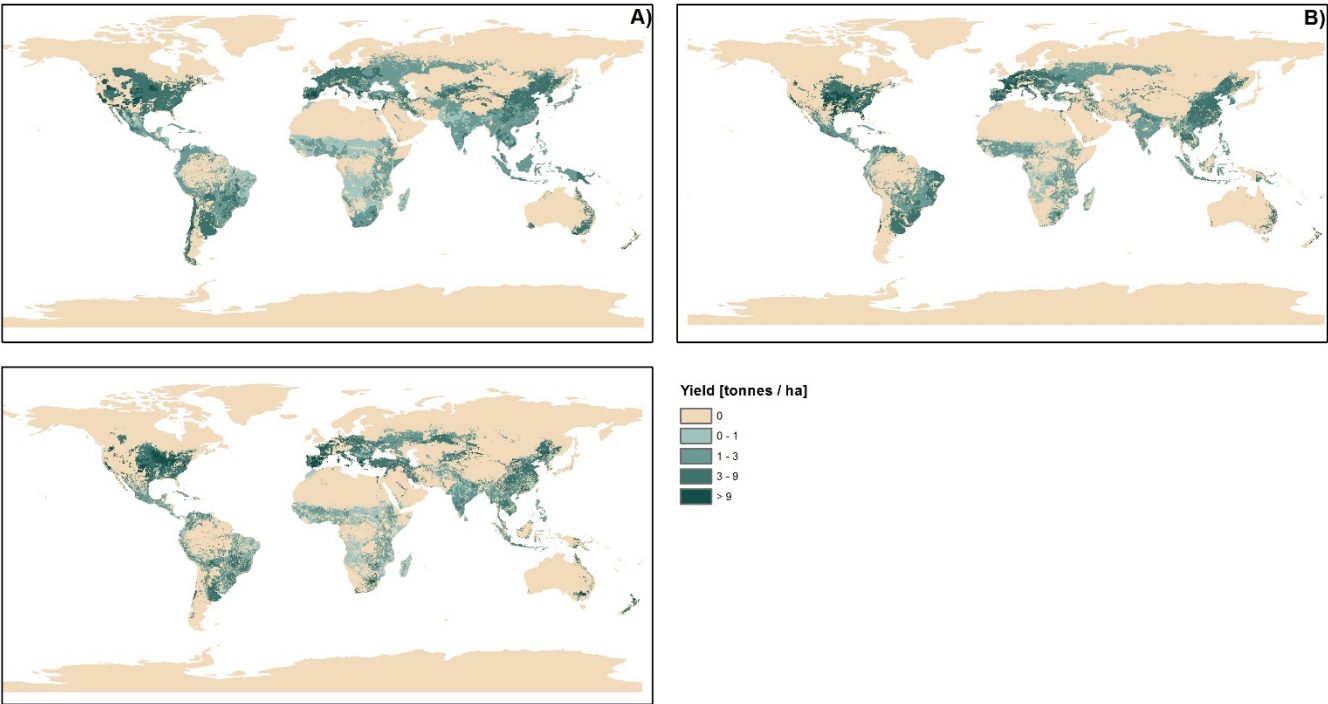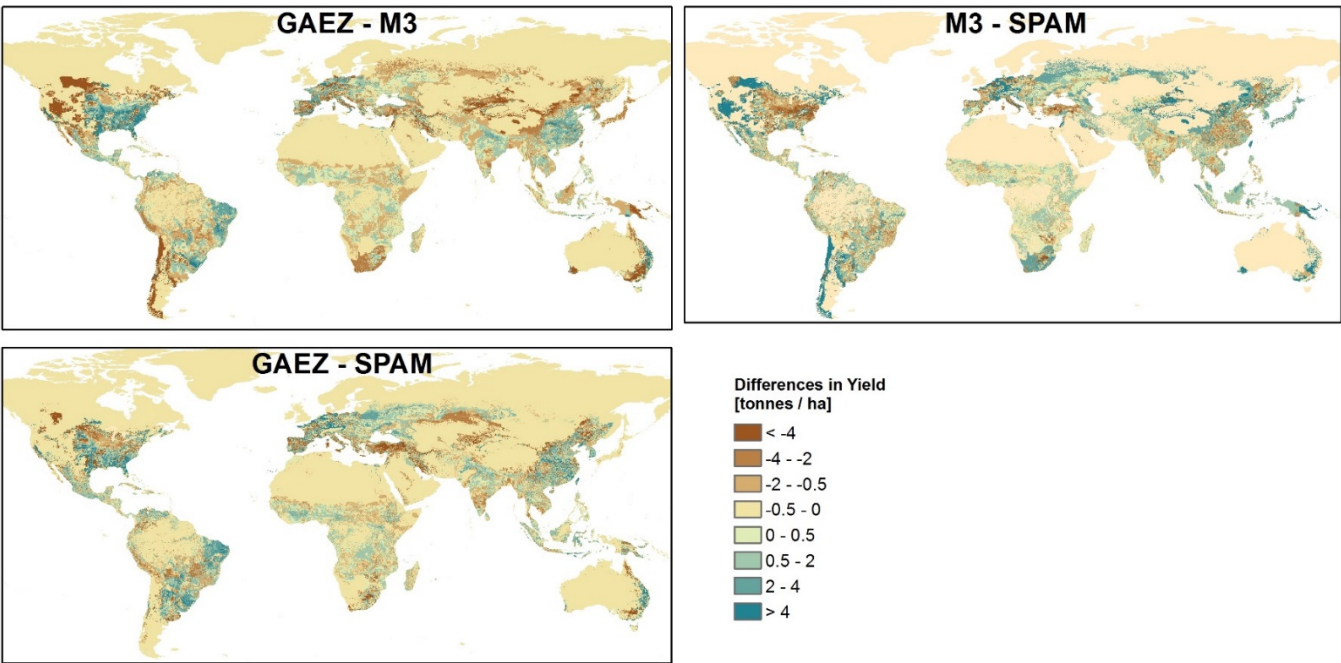Figure S11: Maize yield for A) M3, B) GAEZ, and C) SPAM



Figure S12: pixel-wise comparison of the maize yield by model

# WORKS CITED

[1]   Bartholome, E., & Belward, A. S. (2005), GLC2000: A new approach to global land cover mapping from Earth Observation data, *International Journal of Remote Sensing*, **26**, 1959– 1977.

[2]   Bhaduri, B.B., Bright, E., Coleman, P. & Dobson, J. (2002) LandScan. Geoinformatics, 5, 34–37.

[3]   Dobson, J.E., Brlght, E.A., Coleman, P.R. & Worley, B.A. (2000) LandScan: A Global Population Database for Estimating Populations at Risk. 66, 849–857.

[4]   Fischer, G., Ermolieva, T., Ermoliev, Y. & Velthuizen, H.T. Van (2006) Spatial Recovering of Agricultural Values from Aggregate Information : Sequential Downscaling Methods. International Journal of Knowledge and Systems Sciences (ISKSS), 3.

[5]   Fischer, G.,  F.O. Nachtergaele, S.  Prieler, E.Teixeira, G.Tóth, H. van Velthuizen, L. Verelst, D. Wiberg (2013). Global Agro-Ecological Zones (GAEZ v3.0), http://gaez.fao.org/Main.html#. Accessed October 2013.

[6]   Golan, A., Judge, G. & Miller, D., 1996. Maximum Entropy Econometrics: Robust Estimation with Limited Data, New York: John Wiley & Sons

[7]   Jaynes, E.T. (1957) Information Theory and Statistical Mechanics. 320–330.

[8]   Paris, Q. (1998) An analysis of ill-posed production problems using maximum entropy. American Journal of Agricultural Economics, 80, 124–138.

[9]   Portmann, F., Siebert, S., Bauer, C. & Döll, P. (2008) Global dataset of monthly growing areas of 26 irrigated crops. Frankfurt Hydrology Paper 06, Institute of Physical Geography.University of Frankfurt, Frankfurt am Main, Germany

[10]  Portmann, F.T., Siebert, S. & Döll, P. (2010) MIRCA2000-Global monthly irrigated and rainfed crop areas around the year 2000: A new high-resolution data set for agricultural and hydrological modeling. Global Biogeochemical Cycles, 24.

[11]  Shannon, C.E. (1948) A Mathematical Theory of Communication. Bell System Technology Journal, 27, 379–423.