

1

2

Genetic image-processing using regularized selection indices

3

by

4

5 Marco Lopez-Cruz¹, Eric Olson¹, Gabriel Rovere^{2,3,4}, Jose Crossa⁶, Susanne Dreisigacker⁶,

6 Sushismita Mondal⁶, Ravi Singh⁶, and Gustavo de los Campos^{3,4,5,*}

7

8 ¹ Department of Plant, Soil and Microbial Sciences, Michigan State University, USA

9 ² Department of Animal Science, Michigan State University, USA

10 ³ Department of Epidemiology and Biostatistics, Michigan State University, USA

11 ⁴ Institute for Quantitative Health Science and Engineering, Michigan State University, USA

12 ⁵ Department of Statistics and Probability, Michigan State University, USA

13 ⁶ International Maize and Wheat Improvement Center (CIMMYT), Mexico

14

15 * Corresponding author. E-mail: gustavoc@msu.edu. (GDLC)

16

17

18 **Abstract**

19 High-throughput phenotyping (HTP) technologies can produce data on thousands of phenotypes
20 per unit being monitored. These data can be used to breed for economically and environmentally
21 relevant traits (e.g., drought tolerance); however, incorporating high-dimensional phenotypes in
22 genetic analyses and in breeding schemes poses important statistical and computational
23 challenges. To address this problem, we developed regularized selection indices; the
24 methodology integrates techniques commonly used in high-dimensional phenotypic regressions
25 (including penalization and rank-reduction approaches) into the selection index (SI) framework.
26 Using extensive data from CIMMYT's (International Maize and Wheat Improvement Center)
27 wheat breeding program we show that image-based regularized SIs offer consistently higher
28 accuracy for grain yield than those achieved by canonical SIs and by vegetation indices commonly
29 used to predict agronomic traits. Regularized SIs offer an effective approach to leverage HTP data
30 that is routinely generated in agriculture; the methodology can also be used to conduct genetic
31 studies using high-dimensional phenotypes that are often collected in humans and model
32 organisms including body images and whole-genome gene expression profiles.

33

34 **Author summary**

35 A more intensive use of High-throughput phenotyping (HTP) in breeding programs can increase
36 selection gains and can enable breeding for traits that are otherwise difficult to measure and to
37 breed for (e.g., drought resistance in plants). Most of the phenotypes generated by HTP platforms
38 are high-dimensional, making the use of these data for breeding decisions challenging. We

39 propose to address this problem by using regularized selection indices (SIs). The methodology
40 combines ideas from quantitative genetics with methods used in high-dimensional regressions.
41 Using wheat data from CIMMYT's wheat breeding program we show that regularized SIs deliver
42 more accurate selection decisions than that of canonical SIs.

43

44 **Introduction**

45 High-throughput phenotyping (HTP) technologies have been adopted at a fast pace in agriculture;
46 applications range from the use of HTP in highly controlled environments (e.g., growth chambers
47 [1]) to extensive HTP using sensing devices mounted on aerial (e.g., hyper-spectral cameras
48 mounted on aerial vehicles [2]) and terrestrial equipment such as tractors and combine
49 harvesters [3]. Modern agricultural production systems use HTP data to optimize management
50 practices [4], forecast agricultural outputs [5] and to assess the quality (e.g., protein content) of
51 agricultural commodities [6]. HTP data can also be a valuable input for breeding programs. For
52 instance, extensive HTP may enable an expansion of genetic testing that can lead to higher
53 intensity of selection and faster genetic progress. Moreover, HTP data may offer opportunities
54 to improve traits such as drought tolerance that are otherwise difficult to measure and breed for.

55 Sensors can generate data on hundreds or thousands of phenotypes per unit being
56 monitored. An extensive body of research deals with the use HTP data to predict phenotypes
57 such as grain yield [5,7–9], dry matter [3], oil and protein content [10,11]. However, there has
58 been much less research on how to integrate HTP data in genetic studies and in breeding
59 schemes. In genetics, the problem of predicting the genetic merit of a target trait given a set of

60 correlated phenotypes was first addressed by Smith [12] and Hazel [13] who introduced the
61 concept of selection index (SI) in plant and animal breeding, respectively.

62 A SI seeks to improve a target trait y_i (e.g., grain yield) using information from another
63 set of measured traits (e.g., hyper-spectral image data). A linear SI is a weighted sum of the
64 measured phenotypes with weights derived to maximize the correlation between the genetic
65 merit for the selection target and the SI. Thus, the SI methodology offers a natural framework for
66 integrating HTP data into breeding decisions. However, when the measured phenotype is high-
67 dimensional, the naïve application of the SI can lead to overfitting and sub-optimal accuracy of
68 indirect selection.

69 To address this problem we developed regularized selection indices (including penalized
70 and reduced-rank methods) that are tailored to achieve accurate prediction of genetic values
71 using high-dimensional phenotypes. The proposed methodology integrates into the SI framework
72 methods often used to prevent overfitting in high-dimensional phenotypic regressions [14].
73 Using extensive multi-environment crop imaging data from CIMMYT's wheat breeding program
74 we show that regularized SIs offer improved accuracy of indirect selection in both optimal and
75 stress environments.

76

77 **Results**

78 A **selection index** is a linear combination of p measured phenotypes, $\mathbf{x}_i = (x_{i1}, \dots, x_{ip})'$, of the
79 form $I_i = \mathbf{x}_i' \boldsymbol{\beta}$, where $\boldsymbol{\beta} = (\beta_1, \dots, \beta_p)'$ is a vector of regression coefficients whose entries define
80 the weights of each of the measured phenotypes in the SI. In a canonical SI those weights are

81 derived by minimizing the squared deviation between the genetic merit for the selection target
82 (g_{y_i} , e.g., the genetic merit for grain yield of the i^{th} genotype) and the SI, that is:

$$83 \quad \hat{\boldsymbol{\beta}} = \arg \min_{\boldsymbol{\beta}} \frac{1}{2} \mathbb{E}(g_{y_i} - \mathbf{x}'_i \boldsymbol{\beta})^2. \quad (1)$$

84 The solution to this optimization problem is (see *Methods* section):

$$85 \quad \hat{\boldsymbol{\beta}} = \mathbf{P}_x^{-1} \mathbf{G}_{x,y}, \quad (2)$$

86 where $\mathbf{G}_{x,y} = \mathbb{E}(\mathbf{x}_i g_{y_i}) = (G_{x_1,y}, \dots, G_{x_p,y})'$ is a vector containing the genetic covariances
87 between the selection objective (y_i) and each of the measured traits (\mathbf{x}_i), and \mathbf{P}_x is the
88 (population) phenotypic variance-covariance matrix of the measured phenotypes, that is, $\mathbf{P}_x =$
89 $\mathbb{E}(\mathbf{x}_i \mathbf{x}'_i) = \text{Cov}(\mathbf{x}_i, \mathbf{x}'_i)$. Thus, a canonical SI takes the form $I_i = \mathbf{x}'_i \mathbf{P}_x^{-1} \mathbf{G}_{x,y}$. The theory
90 underlying the derivation of SIs and response to indirect selection is well established [15,16].

91 The SI is by construction the best predictor (in the mean-squared error sense) of the
92 genetic merit for the selection target; this property holds when $\mathbf{G}_{x,y}$ and \mathbf{P}_x are known. However,
93 when the number of measured phenotypes is large errors in the estimation of \mathbf{P}_x and $\mathbf{G}_{x,y}$ may
94 lead to overfitting and sub-optimal accuracy of indirect selection.

95 **Regularized selection indices**

96 Reduced-rank (e.g., principal components methods) and penalized regression [14] are two
97 approaches commonly used to confront overfitting in high-dimensional regression problems.
98 These methodologies were developed for regression problems involving an observable
99 phenotype (y_i). In the SI, the response (g_{y_i}) is unobservable; however, the same principles that

100 are applied in phenotypic reduced-rank and penalized regressions can be integrated into the SI
101 framework.

102 **Reduced-rank selection indices.** In principal components (PC) regression, the response is
103 regressed on a reduced number ($q < p$) of PCs extracted from a set of predictors (x_i); the same
104 concept can be used to derive a reduced-rank SI. For instance, one can extract a reduced number
105 of PCs from a crop image and the resulting PCs can be used as ‘measured traits’ in expression (1).
106 The solution of expression (1) will render estimates of the regression coefficients for the PCs,
107 which can be transformed back to coefficients applicable to the measured traits (see *Methods*).
108 Thus, a reduced-rank SI (referred to as PC-SI) can be derived following these steps: (i) extract,
109 using the singular value decomposition, q PCs from the matrix containing the measured
110 phenotypes, (ii) estimate the genetic covariances between the first q PCs and the selection
111 objective, (iii) use these estimated (co)variances to derive coefficients associated with the top q
112 PCs; finally, (iv) transform these coefficients into coefficients for the measured phenotypes. This
113 process can be done using $q = 1, 2, \dots, p$ PCs ($q = p$ renders the canonical SI). For the sequence
114 of estimates $(\hat{\beta}^{(1)}, \hat{\beta}^{(2)}, \dots, \hat{\beta}^{(p)})$, one can evaluate the accuracy of indirect selection of the
115 resulting SI and an *optimal rank* for the PC-SI can be chosen to maximize the accuracy of indirect
116 selection.

117 **Penalized selection indices.** In a penalized regression, regularization is achieved by
118 including in the objective function a penalty on model complexity. In the context of a SI, we have

119
$$\hat{\beta} = \arg \min_{\beta} \left\{ \frac{1}{2} \mathbb{E}(g_{y_i} - x'_i \beta)^2 + \lambda J(\beta) \right\}, \quad (3)$$

120 where λ is a penalty parameter ($\lambda = 0$ yields the coefficients for the canonical SI) and $J(\boldsymbol{\beta})$ is a
121 penalty function. Commonly used penalties include the L2 ($\|\boldsymbol{\beta}\|_2^2 = \sum_{j=1}^p \beta_j^2$) and L1 ($\|\boldsymbol{\beta}\|_1 =$
122 $\sum_{j=1}^p |\beta_j|$) norms [17], or a weighted sum of the two [18].

123 Using $J(\boldsymbol{\beta}) = \frac{1}{2} \sum_{j=1}^p \beta_j^2$ in expression (3) renders a **Ridge-regression-type PSI** (RR-PSI, see
124 *Methods*):

$$125 \quad \hat{\boldsymbol{\beta}}^{L2} = (\mathbf{P}_x + \lambda \mathbf{I})^{-1} \mathbf{G}_{x,y},$$

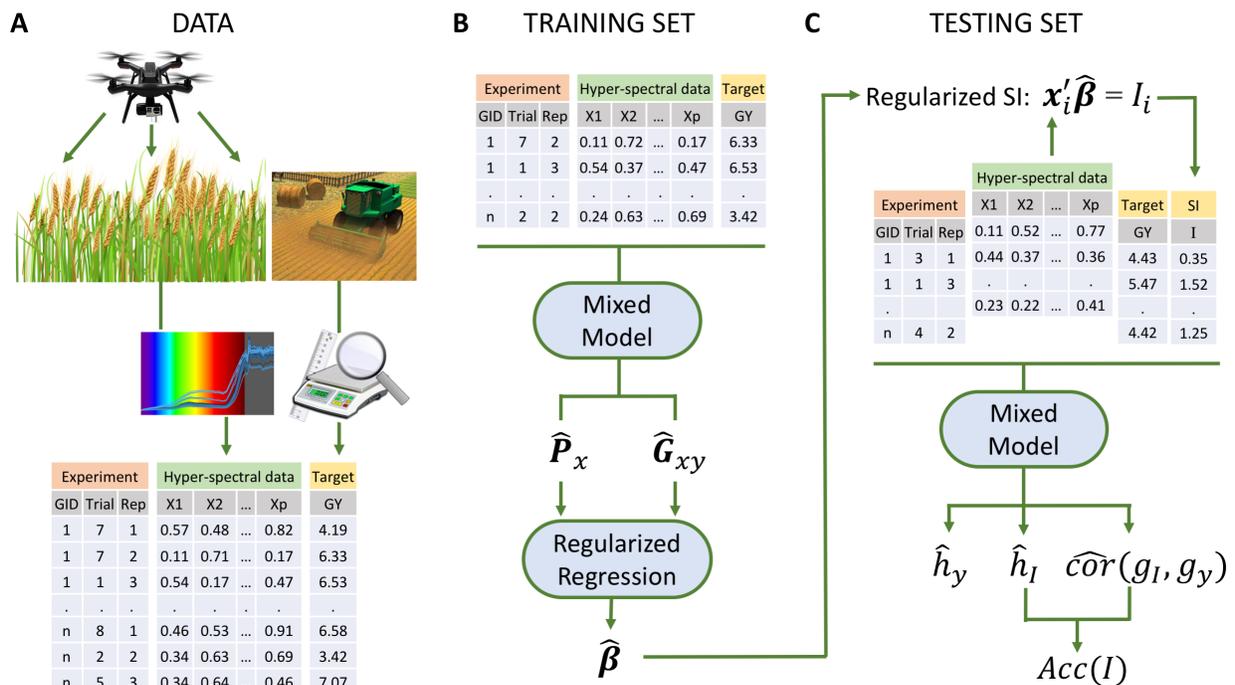
126 where \mathbf{I} is a $p \times p$ identity matrix. The RR-PSI (referred to as the L2-PSI) yields shrunken estimates
127 of the regression coefficients.

128 In many applications, variable selection (i.e., a SI that is a function of a subset of the
129 measured phenotypes) may be desirable. This property can be obtained using penalties involving
130 the L1-norm, either alone, $J(\boldsymbol{\beta}) = \sum_{j=1}^p |\beta_j|$ (LASSO [19]), or in combination with the L2-norm,
131 $J(\boldsymbol{\beta}) = \frac{1}{2} (1 - \alpha) \sum_{j=1}^p \beta_j^2 + \alpha \sum_{j=1}^p |\beta_j|$ (elastic-net [18]). Unlike the L2-PSI, the LASSO and
132 elastic-net SIs (hereinafter referred to as L1-PSI and EN-PSI, respectively) do not have a closed-
133 form solution. However, solutions for PSIs involving an L1-penalty can be obtained using iterative
134 procedures such as the coordinate descent [20] and the least angle regression [21] (LARS)
135 algorithms (see *Methods*). As with the PC-SI, an optimal PSI can be obtained by choosing the
136 values of the regularizing parameters (λ, α) that maximize the accuracy of indirect selection.

137 **Accuracy of indirect selection**

138 Indirect selection accuracy is defined as the correlation between the index used to rank
 139 genotypes and the genetic merit of the selection objective, that is, $Acc(I) = cor(I_i, g_{y_i})$. This
 140 parameter is equal to the product of the square root of the heritability of the SI (h_I) times the
 141 genetic correlation between the SI and the selection target, $cor(g_{I_i}, g_{y_i})$ [16]. To avoid
 142 estimation bias $Acc(I)$ must be estimated using data that was not used to derive the coefficients
 143 of the index (Fig 1); therefore, in the application presented below we: (i) partitioned the data into
 144 training and testing sets, (ii) derived the coefficients of the SI in the training set, (iii) applied these
 145 coefficients to image data of the testing set ($I_i = \mathbf{x}'_i \hat{\boldsymbol{\beta}}$), and (iv) estimated h_I , $cor(g_{I_i}, g_{y_i})$, and
 146 $Acc(I)$ in the testing set. Furthermore, we quantified the efficiency of indirect selection relative
 147 to mass phenotypic selection (RE) using $RE = \frac{h_I}{h_y} cor(g_{I_i}, g_{y_i})$ [16].

148



149

150 **Fig 1. Prediction of the genetic merit for grain yield using hyper-spectral crop image data. (A)**

151 Data consists of hyper-spectral reflectance data (x_i) and phenotypic measurements of the target
152 trait (y_i , e.g., grain yield). (B) A subset of the data (the training set) is used to derive the
153 coefficients (β) of a selection index. (C) These coefficients are then applied to image data of
154 individuals in the testing set to derive the index (I_i) for each individual. The predictive ability of
155 the index is assessed by calculating the accuracy of indirect selection ($Acc(I)$) in the testing set.

156

157 **Regularized selection indices for wheat grain yield using hyper-spectral image data**

158 We applied the methodology described in the previous section to data ($n=3,276$) from the
159 CIMMYT Global Wheat Program consisting of grain yield (ton ha^{-1}) and hyper-spectral image
160 data. The data were collected at CIMMYT's experimental station in Ciudad Obregon, Sonora,
161 Mexico ($27^{\circ}20'$ N, $109^{\circ}54'$ W, 38 masl) from 39 yield trials in which a total of 1,092 genotypes
162 were tested. Rainfall in Obregon is very limited; therefore, four different environments were
163 generated representing a combination of planting methods (*Flat* or *Bed*), controlled irrigation
164 (minimal, 2 or 5 irrigations), and planting dates (optimum or early-heat). As expected, average
165 yield decreased as drought stress intensity increased (see Table 1 and S1 Fig for boxplots of yield
166 by environment).

167

168

169

170 **Table 1 . Average grain yield and heritability by environmental condition**

Planting conditions		Number of irrigations	Abbreviation	Average (SD) Yield	Heritability (SD)
Date	System				
	Flat	Minimal	Flat-Drought	2.06 (0.58)	0.83 (0.016)
Optimum		2	Bed-2IR	3.67 (0.43)	0.66 (0.032)
	Bed	5	Bed-5IR	6.11 (0.61)	0.43 (0.025)
Early		5	Bed-EHeat	6.43 (0.73)	0.61 (0.018)

171 SD: standard deviation.

172

173 Image data was collected using an infrared and an hyper-spectral camera and consisted
 174 of reflectance of electromagnetic power at 250 wavelengths (or bands) within the visible and
 175 near-infrared spectrums (392-850 nm). Separate images were collected at 9 time-points covering
 176 vegetative (VEG), grain filling (GF), and maturity (MAT) stages of the crop (see S2 Fig). Grain yield
 177 and image data were pre-adjusted using mixed-effects model that accounted for genotype, trial,
 178 replicate, and sub-block (see *Methods* section).

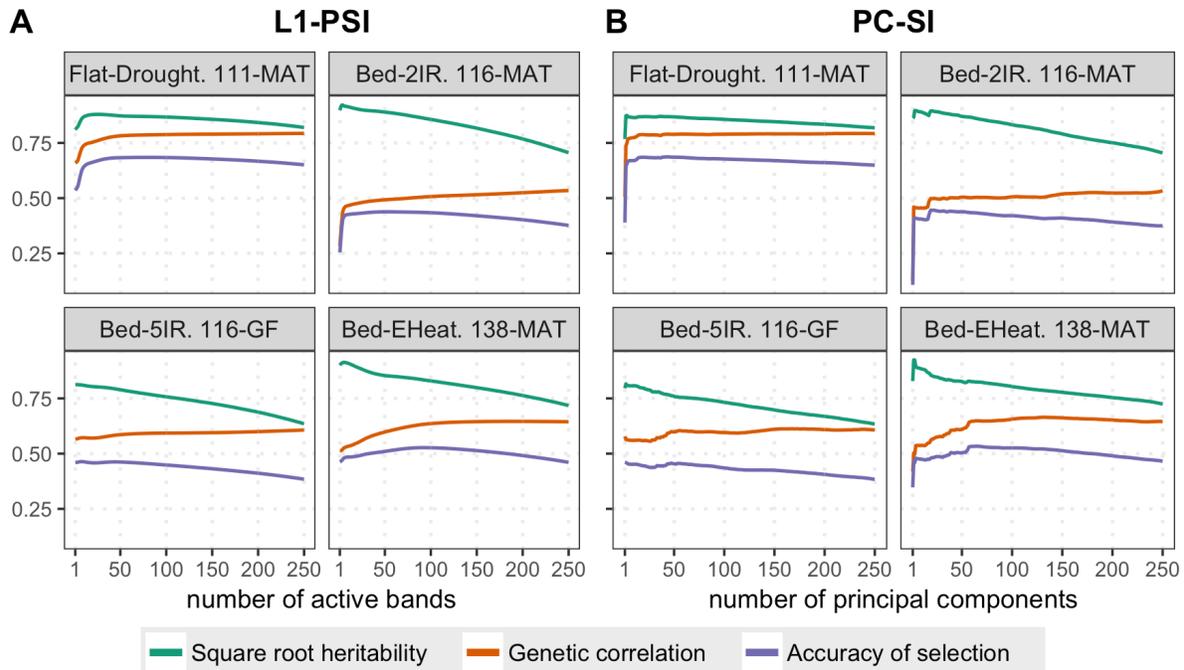
179 **Regularization improves the heritability and the accuracy of the index**

180 To assess the effect of regularization on the accuracy of indirect selection we fitted an L1-PSI over
 181 a grid of values of the regularization parameter ($\lambda^{(1)} > \lambda^{(2)} > \dots > 0$ in expression (3), using $\lambda =$
 182 0 renders a canonical SI). For each of the solutions ($\hat{\beta}(\lambda^{(1)}), \hat{\beta}(\lambda^{(2)}), \dots$) we estimated the
 183 heritability of the resulting index and the genetic correlation between the index and the selection

184 target, and from those estimates we derived the accuracy of indirect selection. The same
185 approach was used to evaluate the accuracy of indirect selection of PC-SIs with a varying number
186 (1, 2, ...) of PCs.

187 We first fitted PSIs and PC-SIs using data from a single time-point; the results from the
188 latest time-point (corresponding to MAT or late GF stages depending on the environment) are
189 presented in Fig 2 (see S3-S5 Figs for other time-points). The heritability of the L1-PSI (Fig 2A)
190 decreased as more bands became active in the index. Likewise, the heritability of PC-SI (Fig 2B)
191 decreased with the number of PCs used. However, the genetic correlation increased as either
192 more bands become active in the L1-PSI or more PCs were used in the PC-SI. Consequently, the
193 maximum accuracy of indirect selection was achieved with a SI of intermediate complexity (with
194 anywhere between 20 and 60 of the 250 bands being active in the L1-PSI, and between 20-60 PCs
195 in the PC-SI). Results for other time-points and environments (S3-S5 Figs) exhibited similar
196 patterns with some differences between environments. The accuracy of indirect selection of the
197 optimal L1-PSI was always close to that of the optimal PC-SI and that of the optimal L2-PSI (S1
198 Table). Importantly, in all cases the accuracy of indirect selection of the optimal regularized SIs
199 was considerably higher than that of the canonical SI, which is the one corresponding to 250
200 active bands or 250 PCs (i.e., the right-most results in the plots in Fig 2).

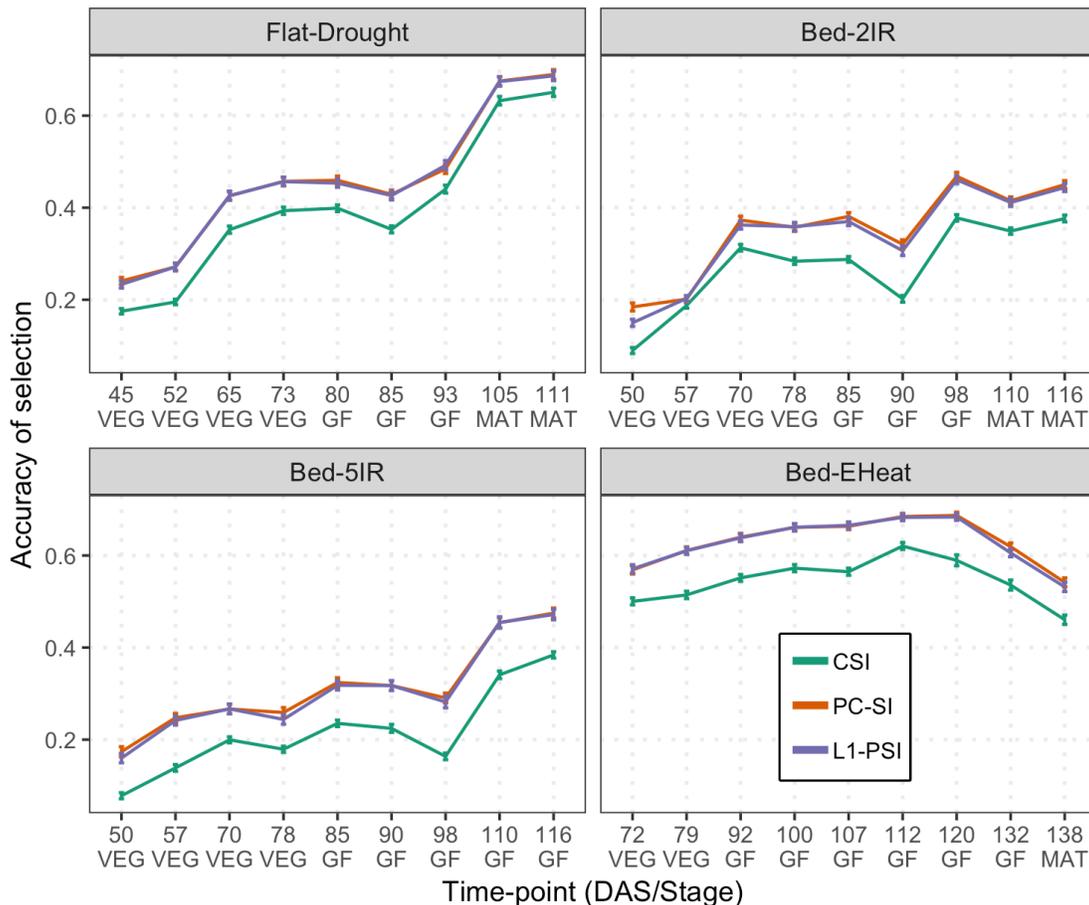
201
202



203
204 **Fig 2. Accuracy of indirect selection of regularized SIs and its components.** Square root
205 heritability (green), genetic correlation (orange) and accuracy of indirect selection (purple), all
206 averaged over 100 training-testing partitions, versus the number of predictors used to build the
207 index: (A) number of active bands in the case of the L1-PSI, or (B) number of PCs in the PC-SI.
208 Each panels represents one environment (latest time-point).

209
210 Fig 3 displays the accuracy of indirect selection across time-points for the optimal (i.e.,
211 the one with the highest accuracy of indirect selection) L1-PSI and PC-SI. For comparison we also
212 display in the plot the accuracy of indirect selection achieved by a canonical SI (in green).
213 Regularization increased the selection accuracy across time-points and environments.
214 Regularized SIs (either PC-SI or L1-PSI) had an accuracy of indirect selection that was 10-30%
215 higher than the accuracy achieved by a canonical SI. Interestingly, there were no sizable
216 differences between the accuracy of indirect selection achieved with the optimal L1-PSI and that

217 of the optimal PC-SI. Compared with the canonical SI, regularized SIs had higher heritability (S6
 218 Fig); this was achieved without compromising the genetic correlation (S7 Fig), thus leading to
 219 a001 higher accuracy of indirect selection achieved by either penalization or reduced-rank
 220 strategies.
 221



222
 223 **Fig 3. Accuracy of indirect selection achieved by a canonical (CSI) and by regularized (PC-SI and**
 224 **L1-PSI) selection indices.** The lines provide the average accuracy over 100 training-testing
 225 partitions. Horizontal lines represent the 95% CI. The horizontal axis give the time-point at which
 226 images were collected and are expressed in both days after sowing (DAS) and stages
 227 (VEG=vegetative, GF=grain filling, MAT=maturity).

228 **Using data from multiple time-points further improves selection accuracy**

229 The results presented above were based on data from a single time-point. We also generated
 230 selection indices using data from multiple time-points (in this case, x_i was a vector containing
 231 2,250 traits, corresponding to 250 wavelengths measured at each of 9 time-points). Integrating
 232 data from multiple time-points further increased the accuracy of L1-PSI by a margin that ranged
 233 from 1 to 8 points on the correlation scale (Table 2). The gains in selection accuracy obtained
 234 using data from multiple time-points were more evident in environments with lower accuracy;
 235 similar results were obtained for the PC-SI and L2-PSI (S1 Table).

236

237 **Table 2. Accuracy and relative efficiency of indirect selection of an L1-penalized SI using data**
 238 **from one and nine time-points.**

Environment	Accuracy (SD)		Relative Efficiency (SD)	
	Best single time-point*	Nine time-points combined	Best single time-point*	Nine time-points combined
Flat-Drought	0.69 (0.05)	0.70 (0.05)	0.74 (0.05)	0.75 (0.05)
Bed-2IR	0.46 (0.04)	0.54 (0.03)	0.57 (0.05)	0.67 (0.04)
Bed-5IR	0.47 (0.06)	0.55 (0.05)	0.72 (0.08)	0.83 (0.08)
Bed-EHeat	0.68 (0.04)	0.71 (0.04)	0.88 (0.05)	0.91 (0.04)

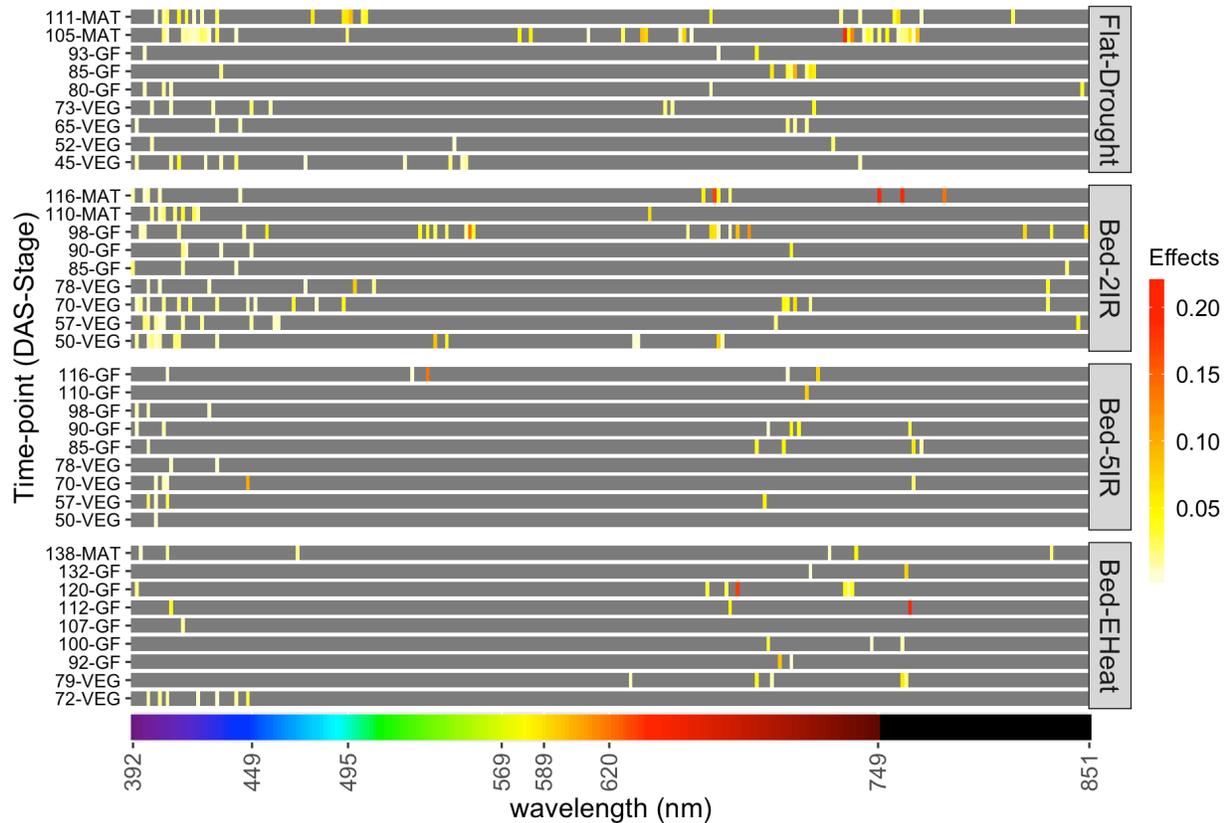
239 Values are presented as an average across 100 training-testing partitions. SD: standard deviation.

240 *For each environment we include the time-point that gave the highest accuracy of selection (see
 241 Fig 3 for other time-points).

242 **L1-penalization leads to sparse selection indices**

243 Fig 4 shows a heatmap for the solutions of the optimal L1-PSI that integrated data from the 9
244 time-points. Each panel represents an environment, horizontal bands represent time-points.
245 Within each time-point bands not entering in the solution are in grey and non-zero coefficients
246 are represented in a yellow-red scale (red indicates large absolute-value coefficients). The well-
247 irrigated environments (*Bed-5IR* and *Bed-EHeat*) had considerably sparser indices with only a
248 reduced number of wavelengths in the solutions; these were mostly located in the violet, blue
249 and red regions of the spectrum. In stressed environments (*Flat-Drought* and *Bed-2IR*) there were
250 also a few wavelengths in the green and infrared regions that were active. In all the indices, there
251 were wavelengths from several time-points that were active in the optimal solution, suggesting
252 that data from both early and late phenological stages are informative about the genetic merit
253 for grain yield.

254



255
 256 **Fig 4. Heatmap of regression coefficients for L1-penalized selection indices.** Separate indices
 257 were derived for each environment using multi time-point data. DAS=days after sowing, VEG, GF,
 258 MAT represent vegetative, grain-filling and maturity stages, respectively. The bottom color-bar
 259 shows the light color associated with each wavelength in the visible spectrum (≤ 750 nm); black
 260 was used to represent the near-infrared spectrum (wavelength > 750 nm).

261

262 Comparison with phenotypic prediction

263 We compared the accuracy of indirect selection of the PSI and PC-SI with vegetation indices and
 264 penalized phenotypic prediction. *Vegetation indices* are often used to predict yield (e.g., [22]),
 265 biomass, and chlorophyll content (e.g., [23,24]). We considered two vegetation indices: the Red

266 and Green Normalized Difference Vegetation Indices (RNDVI [25] and GNDVI [26] respectively).
267 For each of these indices we estimated the genetic correlation with grain yield, as well as their
268 heritability and accuracy of indirect selection (S1 Table). Overall, the accuracy of indirect
269 selection of the GNDVI and RNDVI was lower than the one achieved with a PSI (the average
270 difference in accuracy between RNDVI and the L1-PSI varied by environment from 0.02 to 0.14
271 points in correlation, S1 Table, in favor of the L1-PSI). The heritability of the GNDVI and RNDVI
272 was similar and superior in some cases to that of the L1-PSI (S6 Fig); however, the genetic
273 correlation between the vegetation indices and grain yield was (in most time-points and
274 environments) lower than the genetic correlation between the L1-PSI and grain yield (S7 Fig).
275 Thus, the main driver of the difference in accuracy between the L1-PSI and the vegetation indices
276 was the difference in genetic correlation.

277 We also fitted L1-penalized phenotypic prediction (L1-Phen) and compared the accuracy
278 of indirect selection of these phenotypic prediction methods with that of penalized SIs. Overall,
279 the L1-Phen achieved an accuracy of indirect selection very close to that of the L1-PSI (S1 Table);
280 however, in a few environments at some time-points, the L1-PSI achieved a higher accuracy of
281 indirect selection than the phenotypic prediction.

282

283 **Discussion**

284 High-throughput phenotyping has been extensively adopted in agricultural research and
285 commercial production. Extracting interpretable information from HTP data poses important
286 statistical challenges. The clear majority of research in this area has focused on calibrating

287 equations to predict phenotypes (e.g., total biomass, grain yield) using HTP data as inputs. This
288 approach is well-suited for phenotypic prediction; however, the same approach can be sub-
289 optimal for selection because the best predictor of a phenotype is not always the best predictor
290 of the genetic merit of the same trait.

291 The best phenotypic predictor is the sum of the best predictor of the genetic merit (g_y)
292 plus the best predictor of the environmental term (ε_y), that is, $\mathbb{E}(y|\mathbf{x}) = \mathbb{E}(g_y|\mathbf{x}) + \mathbb{E}(\varepsilon_y|\mathbf{x})$.
293 The first term, $\mathbb{E}(g_y|\mathbf{x})$, is the SI and it is, by construction, maximally correlated with the genetic
294 merit. The second term, $\mathbb{E}(\varepsilon_y|\mathbf{x})$, is relevant for phenotypic prediction but represents noise when
295 the problem is that of selecting the best genotypes.

296 Selection indices exploit genetic covariances, while phenotypic prediction relies on
297 phenotypic covariances between the selection target and the measured phenotype (e.g., crop
298 imaging). Thus, the two methods yield different results whenever the patterns of phenotypic
299 correlations are sufficiently different from the patterns of genetic correlations. In our data set,
300 environmental conditions were highly controlled, with relatively low un-controlled within-trial
301 variability in environmental conditions. Consequently, the patterns of phenotypic and genetic
302 correlations were very similar (see S8 Fig). This was true for many time-points and environments
303 but not in others (e.g., 80, 85 and 93 DAS in *Flat-Drought*, and 90 and 98 DAS in *Bed-2IR*); it was
304 exactly in those time-points and environments that the L1-PSI achieved higher accuracy of
305 indirect selection than the L1-Phen method (S1 Table).

306 A canonical SI (expression (1)) is, by construction, maximally correlated with the genetic
307 merit of the selection objective. This optimality property holds when the genetic and phenotypic

308 (co)variance matrices that are needed to derive the coefficients of the SI (see expression (2)) are
309 known without error. However, when the measured phenotype is high-dimensional, estimation
310 errors in the phenotypic (co)variance matrix (\mathbf{P}_x), as well as in the genetic covariances ($\mathbf{G}_{x,y}$), can
311 make the canonical SI sub-optimal. Our empirical results confirm this: canonical SIs over-fitted
312 the data, this leads to a SI with low heritability and low accuracy of indirect selection.

313 To prevent overfitting, we considered integrating ideas commonly used in high-
314 dimensional regression into the SI methodology. Our empirical results show that regularization
315 consistently improves the accuracy of indirect selection relative to canonical SIs. We verified this
316 for various environmental conditions and for crop imaging data collected at 9 different time-
317 points. The optimal PSI and the optimal PC-SI achieved almost the same accuracy of indirect
318 selection for all the environments and time-points, suggesting that either type of regularization
319 can be effective.

320 **Reduced-rank selection indices** are appealing because after dimension reduction the
321 problem of deriving a SI is trivial and can be dealt with methods commonly used to derive
322 canonical SIs. Moreover, after HTP has been reduced to a few derived-traits (say the top 10 PCs),
323 these traits can be integrated into genetic evaluations (either pedigree-based [27] or genomic-
324 enabled [28]) using standard multi-trait models.

325 Principal components-based methods have been considered before in the analysis of
326 Fourier-transformed infrared (FTIR) spectra derived from milk samples. For instance, Soyeurt,
327 Misztal & Gengler [29] used a reduced number of FTIR-derived PCs to estimate variance
328 components for selection objectives (e.g., fat or protein content in milk). Building upon this idea,

329 Dagnachew, Meuwissen & Ådnøy [30] suggested predicting the genetic merit for milk fatty acids
330 using FTIR-derived PCs as ‘traits’ in a genetic evaluation. However, when mapping from genetic
331 predictions of PC-lodgings onto genetic predictions for the selection objective the authors used
332 coefficients derived from a phenotypic (partial least squares) regression. This does not guarantee
333 that the resulting index is maximally correlated with the genetic merit of the selection target. The
334 penalized and PC-SI presented in this study address that problem by using coefficients that are
335 derived using genetic (and not phenotypic) covariances.

336 A disadvantage of the PC-SI is that the methodology does not naturally provide variable
337 selection, a feature that may be desirable when the measured phenotype is high-dimensional.

338 **Penalized selection indices** can perform variable selection based on genetic covariances.
339 While the derivation of a PSI is a bit more challenging than that of the PC-SI, the computational
340 burden involved in the derivation of a PSI is not extremely high.

341 **Integration of PSI and PC-SI into genetic evaluations.** The SIs considered here predict
342 genetic merit for a selection target from a set of traits measured on an individual ($I_i = \mathbf{x}'_i\boldsymbol{\beta}$); such
343 indices exploit borrowing of information between traits within an individual. Borrowing of
344 information between individuals increases selection accuracy; we envision two ways in which
345 regularized SIs can be integrated into pedigree or genomic-based genetic evaluations.

346 One possibility is to use **two-steps** whereas in the first step a PSI or a PC-SI is used to
347 predict the genetic merit using within-individual information. This step can be considered as a
348 task where patterns attributable to genetic covariances are extracted and those attributable to

349 environmental covariances are smoothed-out. Then, in a second step, the resulting index-data
350 $\{I_1, \dots, I_n\}$ could be used as a trait in a genetic evaluation.

351 A **one-step** approach is also conceptually possible: the optimization problem of
352 expression (3) can be modified by replacing \mathbf{x}_i , the vector with the measured phenotypes on the
353 i^{th} individual, with a vector $\mathbf{x} = (\mathbf{x}'_1, \mathbf{x}'_2, \dots, \mathbf{x}'_n)'$ that contains all the available HTP data
354 (measured on n individuals); after expanding the squared error loss and taking expectations we
355 get

$$356 \quad \hat{\boldsymbol{\beta}}_i = \arg \min_{\boldsymbol{\beta}_i} \left\{ \frac{1}{2} \mathbb{E}(g_{y_i}^2) - \boldsymbol{\beta}'_i \mathbf{G}_{gx} + \frac{1}{2} \boldsymbol{\beta}'_i \mathbf{P}_x \boldsymbol{\beta}_i + \lambda J(\boldsymbol{\beta}_i) \right\},$$

357 where \mathbf{G}_{gx} is a $pn \times 1$ vector of genetic covariances including between-traits-within-individual
358 (co)variances and between-subjects covariances. In standard genetic models, \mathbf{G}_{gx} takes a
359 Kronecker form $\mathbf{G}_{gx} = \mathbf{A}_i \circ \mathbf{G}_{x,y}$, where \mathbf{A}_i are genetic (either DNA- or pedigree-derived)
360 relationships between the candidate for selection and each of the individuals in the training set,
361 and $\mathbf{G}_{x,y}$ is, as before, a vector of genetic covariances between the selection objective and the
362 measured traits (\mathbf{x}). Likewise, \mathbf{P}_x is a $pn \times pn$ phenotypic (co)variance matrix. Estimating \mathbf{P}_x
363 would require estimating all the genetic and environmental covariances among the measured
364 traits. Therefore, while a one-step approach is conceptually feasible, the implementation can be
365 computationally challenging.

366 **Regularized selection indices can also be a valuable tool in genetic research.** High-
367 dimensional phenotypes are also becoming increasingly available in genetic studies involving
368 human subjects and model organisms. Performing genetic studies (e.g., genome-wide association
369 analyses) on high-dimensional phenotypes is challenging and the burden of multiple testing

370 across hundreds or thousands of phenotypes (e.g., RNA-abundance across thousands of genes)
371 may critically compromise power. The PSI and PC-SI presented in this study could be used to
372 extract genetic patterns from high dimensional phenotype data such as brain imaging or whole-
373 genome gene expression profiles and these patterns can then be used as traits in genetic studies.

374 **Conclusion:** we proposed two novel methods for predicting the genetic merit for selection
375 objectives from high-dimensional phenotypes. These phenotypes are becoming increasingly
376 available as the adoption of HTP in crop and animal production increases. The proposed methods
377 integrate regularization procedures commonly used in high-dimensional regressions into the SI
378 methodology. Regularization prevents overfitting and increases the accuracy of index. The
379 methods proposed here can be used to extract genetic patterns from almost any kind of high-
380 dimensional phenotype, including not only HTP data emerging in agriculture but also high-
381 dimensional phenotypes that emerge in genetic studies involving human subjects and model
382 organisms.

383

384 **Methods**

385 **Canonical selection index**

386 The weights on a SI are derived as the solution to the optimization problem of expression (1):

$$387 \quad \hat{\boldsymbol{\beta}} = \arg \min_{\boldsymbol{\beta}} \frac{1}{2} \mathbb{E}(g_{y_i} - \mathbf{x}'_i \boldsymbol{\beta})^2.$$

388 The right-hand side can be expressed as $\mathbb{E}(g_{y_i} - \mathbf{x}'_i \boldsymbol{\beta})^2 = \mathbb{E}(g_{y_i}^2) - 2 \mathbb{E}(g_{y_i} \mathbf{x}_i)' \boldsymbol{\beta} +$
389 $\boldsymbol{\beta}' \mathbb{E}(\mathbf{x}_i \mathbf{x}'_i) \boldsymbol{\beta}$. The first term, $\mathbb{E}(g_{y_i}^2)$, does not involve $\boldsymbol{\beta}$; therefore, it can be dropped from the

390 objective function. Furthermore, if \mathbf{x}_i has null mean, and assuming that the environmental
 391 effects on \mathbf{x}_i are orthogonal to g_{y_i} , then $\mathbb{E}(g_{y_i}\mathbf{x}_i) = \mathbf{G}_{x,y}$ is a vector containing the genetic
 392 covariances between the selection target and each of the measured phenotypes. Likewise,
 393 $\mathbb{E}(\mathbf{x}_i\mathbf{x}_i') = \mathbf{P}_x$ is the phenotypic (co)variance matrix of \mathbf{x}_i . Therefore, the problem in expression
 394 (1) can be written as

$$395 \quad \hat{\boldsymbol{\beta}} = \arg \min_{\boldsymbol{\beta}} \left\{ -\mathbf{G}'_{x,y}\boldsymbol{\beta} + \frac{1}{2}\boldsymbol{\beta}'\mathbf{P}_x\boldsymbol{\beta} \right\}.$$

396 Differentiating the right-hand side with respect to vector $\boldsymbol{\beta}$ and setting the derivatives
 397 equal to zero leads to the first order conditions: $\mathbf{P}_x\hat{\boldsymbol{\beta}} = \mathbf{G}_{x,y}$; therefore,

$$398 \quad \hat{\boldsymbol{\beta}} = \mathbf{P}_x^{-1}\mathbf{G}_{x,y}.$$

399 **Reduced-rank selection index**

400 Recall that the singular value decomposition of a real-valued matrix, $\mathbf{X} = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n]'$
 401 (individuals in rows, phenotypes in columns) takes the form $\mathbf{X} = \mathbf{U}\mathbf{D}\mathbf{V}'$, where $\mathbf{U} = [\mathbf{u}_1, \dots, \mathbf{u}_p]$
 402 is the matrix containing the left-singular vectors that span the row-space of \mathbf{X} , $\mathbf{V} = [\mathbf{v}_1, \dots, \mathbf{v}_p]$ is
 403 the matrix with the right-singular vectors, and $\mathbf{D} = \text{diag}(d_1, \dots, d_p)$ is a diagonal matrix with
 404 positive or zero elements. The PCs $\mathbf{W} = \mathbf{X}\mathbf{V} = \mathbf{U}\mathbf{D}$ are linear combinations of the measured
 405 phenotypes. A reduced-rank regression uses the first q PCs ($\tilde{\mathbf{W}} = [\mathbf{w}_1, \dots, \mathbf{w}_q]$, $q \leq p$) as
 406 ‘measured phenotypes’ in the SI:

$$407 \quad \hat{\boldsymbol{\gamma}}^{(q)} = \arg \min_{\boldsymbol{\gamma}} \frac{1}{2} \mathbb{E}(g_{y_i} - \tilde{\mathbf{w}}_i'\boldsymbol{\gamma}^{(q)})^2,$$

408 where $\tilde{\mathbf{w}}_i$ is a vector containing the scores for the i^{th} observation on the first q PCs. The solution
 409 to the optimization problem takes the form $\hat{\boldsymbol{\gamma}}^{(q)} = \mathbf{P}_{\tilde{\mathbf{w}}}^{-1} \mathbf{G}_{\tilde{\mathbf{w}},y}$, where $\mathbf{P}_{\tilde{\mathbf{w}}}$ is the phenotypic
 410 (co)variance matrix of the first q PCs and $\mathbf{G}_{\tilde{\mathbf{w}},y}$ is a vector containing the genetic covariances
 411 between each of the top q PCs and the selection objective. Since the left-singular vectors are
 412 orthonormal (i.e., $\mathbf{u}'_j \mathbf{u}_j = 1$ and $\mathbf{u}'_j \mathbf{u}_k = 0$, for $j \neq k$), then $\mathbf{W}'\mathbf{W} = \mathbf{D}^2 = \text{diag}(d_1^2, \dots, d_p^2)$.
 413 Hence, a method-of-moments estimate of the phenotypic (co)variance matrix of $\tilde{\mathbf{W}}$ contains only
 414 the first q elements $\tilde{\mathbf{D}}^2 = \text{diag}(d_1^2, \dots, d_q^2)$; this is

$$415 \quad \hat{\mathbf{P}}_{\tilde{\mathbf{w}}} = \frac{1}{n-1} \tilde{\mathbf{D}}^2.$$

416 Using $\hat{\mathbf{P}}_{\tilde{\mathbf{w}}}$ makes the coefficients of the PCs proportional to the genetic covariance
 417 between each of the PCs and the selection objective: $\hat{\boldsymbol{\gamma}}^{(q)} = (n-1)(\tilde{\mathbf{D}}^2)^{-1} \mathbf{G}_{\tilde{\mathbf{w}},y}$. This solution
 418 can be mapped to coefficients for the measured traits using $\hat{\boldsymbol{\beta}}^{(q)} = (n-1)\tilde{\mathbf{V}}(\tilde{\mathbf{D}}^2)^{-1} \mathbf{G}_{\tilde{\mathbf{w}},y}$, where
 419 $\tilde{\mathbf{V}}$ is the matrix containing only the first q right-singular vectors.

420 Penalized selection indices

421 The objective function of the penalized SI is given by expression (3). Here we considered PSIs
 422 using either L1 or L2-norms or a combination of the two.

423 **L2-PSI:** Using an L2-norm as penalty, $J(\boldsymbol{\beta}) = \frac{1}{2} \sum_{j=1}^p \beta_j^2 = \frac{1}{2} \boldsymbol{\beta}'\boldsymbol{\beta}$, in expression (3) leads to
 424 the following optimization problem:

$$425 \quad \hat{\boldsymbol{\beta}}^{L2} = \arg \min_{\boldsymbol{\beta}} \left\{ \frac{1}{2} \mathbb{E}(g_{y_i} - \mathbf{x}'_i \boldsymbol{\beta})^2 + \lambda \frac{1}{2} \boldsymbol{\beta}'\boldsymbol{\beta} \right\}.$$

426 Therefore:

$$427 \quad \widehat{\boldsymbol{\beta}}^{L2} = \arg \min_{\boldsymbol{\beta}} \left\{ -\mathbf{G}'_{x,y} \boldsymbol{\beta} + \frac{1}{2} \boldsymbol{\beta}' \mathbf{P}_x \boldsymbol{\beta} + \lambda \frac{1}{2} \boldsymbol{\beta}' \boldsymbol{\beta} \right\}.$$

428 The second and third right-hand side terms can be combined to obtain:

$$429 \quad \widehat{\boldsymbol{\beta}}^{L2} = \arg \min_{\boldsymbol{\beta}} \left\{ -\mathbf{G}'_{x,y} \boldsymbol{\beta} + \frac{1}{2} \boldsymbol{\beta}' (\mathbf{P}_x + \lambda \mathbf{I}) \boldsymbol{\beta} \right\},$$

430 where \mathbf{I} is a $p \times p$ identity matrix. Differentiating with respect to $\boldsymbol{\beta}$ and setting the derivatives
431 equal to zero, we obtain the first-order conditions: $(\mathbf{P}_x + \lambda \mathbf{I}) \widehat{\boldsymbol{\beta}}^{L2} = \mathbf{G}_{x,y}$; therefore:

$$432 \quad \widehat{\boldsymbol{\beta}}^{L2} = (\mathbf{P}_x + \lambda \mathbf{I})^{-1} \mathbf{G}_{x,y}.$$

433 **EN-PSI.** The coefficients for the elastic-net family are obtained by considering an objective
434 function as in expression (3), with $J(\boldsymbol{\beta}) = \frac{1}{2} (1 - \alpha) \sum_{j=1}^p \beta_j^2 + \alpha \sum_{j=1}^p |\beta_j|$; therefore,

$$435 \quad \widehat{\boldsymbol{\beta}}^{EN} = \arg \min_{\boldsymbol{\beta}} \left\{ -\mathbf{G}'_{x,y} \boldsymbol{\beta} + \frac{1}{2} \boldsymbol{\beta}' \mathbf{P}_x \boldsymbol{\beta} + \lambda \frac{1}{2} (1 - \alpha) \sum_{j=1}^p \beta_j^2 + \lambda \alpha \sum_{j=1}^p |\beta_j| \right\}.$$

436 The L1-PSI and L2-PSI are particular cases corresponding to $\alpha = 1$ and $\alpha = 0$,
437 respectively. When $\alpha = 0$ the solution has a closed form (see L2-PSI above). If $\alpha > 0$, no closed-
438 form solution exists; however, a solution can be obtained using the same iterative algorithms
439 that are used to solve elastic-net regressions (e.g., LARS and coordinate descent [14]). These
440 algorithms can be implemented either by ‘partial residuals’ or using ‘covariance updates’ [31]. In
441 our case, the objective function is entirely based on (co)variance terms. The objects \mathbf{P}_x and $\mathbf{G}_{x,y}$
442 enter in the objective function of the PSI in the same way that $\mathbf{X}'\mathbf{X}$ and $\mathbf{X}'\mathbf{y}$ enter in a standard
443 elastic-net regression. Therefore, to obtain solutions, we implemented the standard LARS
444 algorithm (e.g., Hastie *et al.* [14]) entirely based on covariance updates.

445 **Data**

446 The data set consists of 1,092 inbred wheat lines grouped into 39 trials and grown during the
447 2013-2014 season at the Norman Borlaug experimental research station in Ciudad Obregon,
448 Sonora, Mexico. Each trial consisted of 28 breeding lines that were arranged in an alpha-lattice
449 design with three replicates and six sub-blocks. The trials were grown in four different
450 environments: *Flat-Drought* (sowing in flat with irrigation of 180 mm through drip system), *Bed-*
451 *2IR* (sowing in bed with 2 irrigations approximately 250 mm), *Bed-EHeat* (bed sowing 30 days
452 before optimal planting date with 5 normal irrigations approximately 500 mm), and *Bed-5IR* (bed
453 sowing with 5 normal irrigations). In 2013, all the trials were planted by mid-November (optimal
454 planting date), on the 21st (*Bed-2IR* and *Bed-5IR*) and on the 26th for *Flat-Drought*. Trials for *Bed-*
455 *EHeat* were planted on October 30th. Grain yield (ton ha⁻¹, total plot yield after maturity) was
456 recorded. Reflectance phenotypic data were collected from the fields using both infrared (A600
457 series Infrared camera, FLIR, Wilsonville, OR) and hyper-spectral (A-series, Micro-Hyperspec,
458 VNIR Headwall Photonics, Fitchburg, MA) cameras mounted on a Piper PA-16 Clipper aircraft on
459 9 different dates (time-points) between January 10 and March 27th, 2014. During each flight,
460 data from $p = 250$ wavelengths ranging from 392 to 850 nm were collected for each pixel in the
461 pictures. The average reflectance of all the pixels for each wavelength was calculated from each
462 of the geo-referenced trial plots and reported as each line reflectance. Days to heading were
463 recorded as the number of days from the date of sowing/first irrigation until 50% of spike
464 emergence in each plot. Heading of about 50-80% of the total number of plots was used as
465 criterion to distinguish between vegetative (VEG) and grain filling (GF) stages. The crop was
466 considered to be at maturity (MAT) stage when the average RNDVI decreased to ~0.4.

467 **Phenotype pre-processing**

468 Within each environment, grain yield phenotypic records were pre-adjusted by fitting the
469 following mixed model,

$$470 \quad y_{jklm} = \mu + g_j + t_k + r_{l(k)} + b_{m(kl)} + \varepsilon_{jklm},$$

471 where y_{jklm} is the grain yield phenotype value for the j^{th} genotype, k^{th} trial, l^{th} replicate (within
472 trial), m^{th} sub-block (within trial and replicate), μ is the overall mean and g_j , t_k , $r_{l(k)}$, and
473 $b_{m(kl)}$ are the genotype, trial, replicate, and sub-block effects, respectively (all assumed to be
474 random) and ε_{jklm} is an error term. Random effects were assumed to be independently and
475 identically distributed (*iid*) normal with null mean and effect-specific variances. Likewise, the
476 error terms were assumed to be *iid* with null mean and common error variance.

477 Grain yield data were pre-adjusted by subtracting from the phenotypic record (y_{jklm}) the
478 mean ($\hat{\mu}$) plus BLUPs of trial, replicate, and sub-block effects; this is

$$479 \quad y_{jklm}^* = y_{jklm} - \hat{\mu} - \hat{t}_k - \hat{r}_{l(k)} - \hat{b}_{m(kl)} = \hat{g}_j + \hat{\varepsilon}_{jklm}. \quad (4)$$

480 Reflectance data was pre-adjusted by fitting the above model, using reflectance at
481 individual bands as phenotype expanded with the inclusion of a time-point effect. Separate
482 models were fitted to each of the wavelengths. As with grain yield, reflectance data were pre-
483 adjusted by subtracting from the measured reflectance the estimated mean and predicted time-
484 point, trial, replicate, and sub-block effects.

485 For quality control, pre-adjusted grain yield and reflectance phenotypes were removed
486 for those grain yield scores lying beyond 3 times the inter-quantile region from the 0.25 and 0.75
487 quantiles.

488 After pre-adjusting, all phenotypes were standardized (to have unit variance); for ease of

489 exposition, hereinafter we refer to the adjusted-scaled phenotypes (including grain yield and the
490 image data) simply as phenotypes.

491 **Heritability estimation**

492 After pre-adjusting standardization, we analyzed phenotypes using a mixed model of the form

$$493 \quad y_{ij} = g_j + \varepsilon_{ij}, \quad (5)$$

494 where y_{ij} is the phenotype for the i^{th} observation (i here is a single index for indices k , l , and m

495 in expression (4)) of the j^{th} genotype; the genetic values are $g_j \stackrel{iid}{\sim} N(0, \sigma_{g_y}^2)$, where $\sigma_{g_y}^2$ is the

496 genetic variance; and the environmental terms are $\varepsilon_{ij} \stackrel{iid}{\sim} N(0, \sigma_{\varepsilon_y}^2)$. Plot-basis heritability was

497 calculated from variance components estimates using

$$498 \quad h_y^2 = \frac{\sigma_{g_y}^2}{\sigma_{g_y}^2 + \sigma_{\varepsilon_y}^2}.$$

499 **Training-testing partitions**

500 The data set contains information from 39 trials with 84 observations each. To assess the

501 accuracy of indirect selection, we randomly assigned trials to training or testing sets. Twenty-six

502 trials ($n_{trn} \approx 2,184$ observations) were randomly assigned to the training set, and the remaining

503 13 trials ($n_{tst} \approx 1,092$) were used as the testing set. The regression coefficients of the indices

504 (the β 's for the canonical SI, PSI, and PC-SI) were calculated using grain yield and reflectance data

505 of the training set. Estimates of the coefficients and reflectance data were then used to calculate

506 the SI $I_{ij} = \mathbf{x}'_{ij} \hat{\boldsymbol{\beta}}$, for each observation i in the testing set ($i = 1, \dots, n_{tst}$). The heritability of the

507 index and the genetic correlation between the index and the selection goal were estimated in

508 the testing set.

509 The training-testing procedure was repeated 100 times by randomly assigning trials to
510 training and testing sets. From these analyses, we reported the mean of heritability, genetic
511 correlation, and accuracy; and their standard deviation across training-testing partitions.

512 Estimation of phenotypic and genetic parameters

513 The population phenotypic (co)variance matrix \mathbf{P}_x was estimated within the training set using
514 the unbiased sample (co)variance matrix given by $\hat{\mathbf{P}}_x = \frac{1}{n-1} \sum_{i=1}^{n_{trn}} (\mathbf{x}_i - \bar{\mathbf{x}})(\mathbf{x}_i - \bar{\mathbf{x}})'$, where $\bar{\mathbf{x}}$ is
515 the vector containing the sample mean of each wavelength. Since reflectance data are centered
516 and standardized, this reduces to $\hat{\mathbf{P}}_x = \frac{1}{n-1} \mathbf{X}'\mathbf{X}$, where $\mathbf{X} = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n]'$ is the matrix
517 containing all measured traits in the training set.

518 The genetic covariance ($G_{x_j, y}$) between grain yield and the j^{th} measured trait ($j = 1, \dots, p$)
519 was estimated using a sequence of univariate genetic models as in expression (5). We fitted that
520 model with grain yield phenotypes as response, then for each of the reflectance bands and then
521 for the sum of grain yield and each of the bands. The genetic covariances between the bands and
522 grain yield were then estimated using

$$523 \quad \hat{\mathbf{G}}_{y, x_j} = \frac{1}{2} \left(\hat{\sigma}_{g_{y+x_j}}^2 - \hat{\sigma}_{g_y}^2 - \hat{\sigma}_{g_{x_j}}^2 \right)',$$

524 where $\hat{\sigma}_{g_y}^2$, $\hat{\sigma}_{g_{x_j}}^2$ and $\hat{\sigma}_{g_{y+x_j}}^2$ are the estimated genetic variances for grain yield, the j^{th} band, and
525 the sum of grain yield and the j^{th} band, respectively.

526

527

528 Estimation of the accuracy of indirect selection

529 To assess the accuracy of indirect selection we applied the regression coefficients derived in the
530 training set to image data from the testing set to derive $I_{ij} = \mathbf{x}'_{ij} \widehat{\boldsymbol{\beta}}$. Then, using a mixed model
531 analysis like that described in the previous section we estimated the heritability of the SI (h_I^2),
532 the heritability of grain yield (h_y^2), and the genetic correlation between the SI and grain yield
533 ($cor(g_{I_i}, g_{y_i})$). From these estimates, we derived the accuracy of indirect selection, $Acc(I) =$
534 $h_I cor(g_{I_i}, g_{y_i})$, and the relative efficiency, $RE = \frac{h_I}{h_y} cor(g_{I_i}, g_{y_i})$.

535 Software

536 All the aforementioned analyses were implemented in the R software environment [32], version
537 3.5.1. Linear mixed models were implemented using the 'lmer' function from the LME4 [33] R-
538 package. Code that implements LARS in the context of SIs was programed based on the LARS [34]
539 R-package.

540 Materials and data availability

541 The data used in this study are publicly available by CIMMYT (<https://www.cimmyt.org/>) who
542 owns all rights in the data. Data sets and R-scripts to perform all the analyses are publicly
543 available upon request to the corresponding author.

544

545

546

547 **Acknowledgments**

548 We acknowledge CIMMYT's Global Wheat Program that provided both experimental field and
549 HTP data used in this work. MLC was supported by the Monsanto's Beachell-Borlaug
550 International Scholarship Program (MBBISP).

551 **Author contributions**

552 RS, SD, JC and SM were involved in the design of the field experiments and data collection. SM
553 performed the HTP data correction and georeferencing. MLC and GDLC conceived the idea,
554 performed the analyses and produced a first draft, all the authors contributed to the final
555 manuscript.

556

557 **References**

- 558 1. Nagel KA, Putz A, Gilmer F, Heinz K, Fischbach A, Pfeifer J, et al. GROWSCREEN-Rhizo is a
559 novel phenotyping robot enabling simultaneous measurements of root and shoot growth
560 for plants grown in soil-filled rhizotrons. *Funct Plant Biol.* 2012;39(11):891–904.
- 561 2. Araus L, Cairns JE. Field high-throughput phenotyping: the new crop breeding frontier.
562 *Trends Plant Sci.* 2014;19(1):52–61.
- 563 3. Montes JM, Utz HF, Schipprack W, Kusterer B, Muminovic J, Paul C, et al. Near-infrared
564 spectroscopy on combine harvesters to measure maize grain dry matter content and
565 quality parameters. *Plant Breed.* 2006;125:591–5.
- 566 4. White JW, Andrade-Sanchez P, Gore MA, Bronson KF, Coffelt TA, Conley MM, et al. Field
567 Crops Research Field-based phenomics for plant genetics research. *F Crop Res.*

- 568 2012;133:101–12.
- 569 5. Ferrio JP, Villegas D, Zarco J, Aparicio N, Araus JL, Royo C. Assessment of durum wheat
570 yield using visible and near-infrared reflectance spectra of canopies. *F Crop Res.*
571 2005;94:126–48.
- 572 6. Spielbauer G, Armstrong P, Baier JW, Allen WB, Richardson K, Shen B, et al. High-
573 throughput near-infrared reflectance spectroscopy for predicting quantitative and
574 qualitative composition phenotypes of individual maize kernels. *Cereal Chem.*
575 2009;86(5):556–64.
- 576 7. Garriga M, Romero-Bravo S, Estrada F, Escobar A, Matus IA, del Pozo A, et al. Assessing
577 wheat traits by spectral reflectance: do we really need to focus on predicted trait-values
578 or directly identify the elite genotypes group? *Front Plant Sci.* 2017;8(280):1–12.
- 579 8. Weber VS, Araus JL, Cairns JE, Sanchez C, Melchinger AE, Orsini E. Prediction of grain yield
580 using reflectance spectra of canopy and leaves in maize plants grown under different
581 water regimes. *F Crop Res.* 2012;128:82–90.
- 582 9. Aguate FM, Trachsel S, González-Pérez L, Burgueño J, Crossa J, Balzarini M, et al. Use of
583 hyperspectral image data outperforms vegetation indices in prediction of maize yield. *Crop*
584 *Sci.* 2017;57:2517–24.
- 585 10. Garnsworthy PC, Wiseman J, Fegeros K. Prediction of chemical, nutritive and agronomic
586 characteristics of wheat by near infrared spectroscopy. *J Agric Sci.* 2000;135:409–17.
- 587 11. Oblath EA, Isbell TA, Berhow MA, Allen B, Archer D, Brown J, et al. Development of near-
588 infrared spectroscopy calibrations to measure quality characteristics in intact Brassicaceae

- 589 germplasm. *Ind Crop Prod.* 2016;89:52–8.
- 590 12. Smith HF. A discriminant function for plant selection. *Ann Eugen.* 1936;7:240–50.
- 591 13. Hazel LN. The genetic basis for constructing selection indexes. *Genetics.* 1943;28(6):476–
592 90.
- 593 14. Hastie T, Tibshirani R, Friedman JH. *The elements of statistical learning: data mining,*
594 *inference, and prediction.* 2nd ed. New York, USA: Springer; 2009. 745 p.
- 595 15. Bulmer MG. *The mathematical theory of quantitative genetics.* New York, USA: Oxford
596 University Press; 1985.
- 597 16. Falconer DS, Mackay TFC. *Introduction to quantitative genetics.* 4th ed. Essex, UK: Prentice
598 Hall; 1996.
- 599 17. Fu WJ. Penalized regressions: the Bridge versus the LASSO. *J Comput Graph Stat.*
600 1998;7(3):397–416.
- 601 18. Zou H, Hastie T. Regularization and variable selection via the elastic net. *J R Stat Soc B.*
602 2005;67(2):301–20.
- 603 19. Tibshirani R. Regression shrinkage and selection via the LASSO. *J R Stat Soc B.*
604 1996;58(1):267–88.
- 605 20. Friedman J, Hastie T, Höfling H, Tibshirani R. Pathwise coordinate optimization. *Ann Appl*
606 *Stat.* 2007;1(2):302–32.
- 607 21. Efron B, Hastie T, Johnstone I, Tibshirani R. Least angle regression. *Ann Stat.*
608 2004;32(2):407–99.

- 609 22. Tattaris M, Reynolds MP, Chapman SC. A direct comparison of remote sensing approaches
610 for high-throughput phenotyping in plant breeding. *Front Plant Sci.* 2016;7:1–9.
- 611 23. Babar MA, Reynolds MP, Ginkel M Van, Klatt AR, Raun WR, Stone ML. Spectral reflectance
612 to estimate genetic variation for in-season biomass, leaf chlorophyll, and canopy
613 temperature in wheat. *Crop Sci.* 2006;46:1046–57.
- 614 24. Haboudane D, Miller JR, Tremblay N, Zarco-Tejada PJ, Dextraze L. Integrated narrow-band
615 vegetation indices for prediction of crop chlorophyll content for application to precision
616 agriculture. *Remote Sens Environ.* 2002;81:416–26.
- 617 25. Tucker CJ. Red and photographic infrared linear combinations for monitoring vegetation.
618 *Remote Sens Environ.* 1979;8(2):127–50.
- 619 26. Gitelson AA, Kaufman YJ, Merzlyak MN. Use of a green channel in remote sensing of global
620 vegetation from EOS-MODIS. *Remote Sens Environ.* 1996;58(3):289–98.
- 621 27. Henderson CR, Quaas RL. Multiple trait evaluation using relatives' records. *J Anim Sci.*
622 1976;43(6):1188–97.
- 623 28. Meuwissen THE, Hayes BJ, Goddard ME. Prediction of total genetic value using genome-
624 wide dense marker maps. *Genetics.* 2001;157(4):1819–29.
- 625 29. Soyeurt H, Misztal I, Gengler N. Genetic variability of milk components based on mid-
626 infrared spectral data. *J Dairy Sci.* 2010;93(4):1722–8.
- 627 30. Dagnachew BS, Meuwissen THE, Ådnøy T. Genetic components of milk Fourier-transform
628 infrared spectra used to predict breeding values for milk composition and quality traits in

- 629 dairy goats. J Dairy Sci. 2013;96(9):5933–42.
- 630 31. Friedman J, Hastie T, Tibshirani R. Regularization paths for generalized linear models via
631 coordinate descent. J Stat Softw. 2010;33(1):1–22.
- 632 32. R Core Team. R: A Language and Environment for Statistical Computing. Vienna, Austria: R
633 Foundation for Statistical Computing; 2018.
- 634 33. Bates D, Mächler M, Bolker BM, Walker SC. Fitting linear mixed-effects models using lme4.
635 J Stat Softw. 2015;67(1):1–48.
- 636 34. Hastie T, Efron B. lars: least angle regression, Lasso and forward stagewise [Internet]. 2013.
637 Available from: <https://cran.r-project.org/package=lars>

638

639 **Supporting information**

640 **S1 Fig. Box-plot of grain yield phenotypic records by environmental condition.** $n \approx 3200$. SD:
641 standard deviation.

642 **S2 Fig. Light reflectance patterns as function of the wavelength.** Each line represents the mean
643 (across $n \approx 3200$ observations) for each wavelength, within time-point (flight date). Within each
644 environment, means were scaled to lie within 0 and 1 by dividing them by the maximum average.

645 **S3 Fig. Accuracy of indirect selection of L1-PSI and its components.** Square root heritability,
646 genetic correlation and accuracy of indirect selection, all averaged over 100 training-testing
647 partitions versus the number of bands entering in the index; by time-point (DAS=days after
648 sowing, Stage: VEG=vegetative, GF=grain filling, or MAT=maturity) within environment.

649 **S4 Fig. Accuracy of indirect selection of L2-PSI and its components.** Square root heritability,
650 genetic correlation and accuracy of indirect selection, all averaged over 100 training-testing
651 partitions versus the penalization parameter (λ , logarithm scale) used to build the index; by time-
652 point (DAS=days after sowing, Stage: VEG=vegetative, GF=grain filling, or MAT=maturity) within
653 environment.

654 **S5 Fig. Accuracy of indirect selection of PC-SI and its components.** Square root heritability,
655 genetic correlation and accuracy of indirect selection, all averaged over 100 training-testing
656 partitions versus the number of principal components used to build the index; by time-point
657 (DAS=days after sowing, Stage: VEG=vegetative, GF=grain filling, or MAT=maturity) within
658 environment.

659 **S6 Fig. Square root of heritability of the canonical (CSI), of the regularized (PC-SI and L1-PSI)**
660 **selection indices, and of the RNDVI.** The lines provide the average square root heritability over
661 100 training-testing partitions. Horizontal lines represent the 95% CI. The horizontal axis give the
662 time-point at which images were collected and are expressed in both days after sowing (DAS)
663 and stages (VEG=vegetative, GF=grain filling, MAT=maturity).

664 **S7 Fig. Genetic correlation between grain yield and all: the canonical (CSI), the regularized (PC-**
665 **SI and L1-PSI) selection indices, and the RNDVI.** The lines provide the average genetic correlation
666 over 100 training-testing partitions. Horizontal lines represent the 95% CI. The horizontal axis
667 give the time-point at which images were collected and are expressed in both days after sowing
668 (DAS) and stages (VEG=vegetative, GF=grain filling, MAT=maturity).

669 **S8 Fig. Phenotypic, genetic, and environmental covariances between wavelengths and grain**
670 **yield.** 'D': discrepancy between phenotypic and genetic covariances as measured by the sum of
671 the absolute differences; by time-point (DAS: days after sowing, Stage: VEG=vegetative, GF=grain
672 filling, MAT=maturity) within environment.

673 **S1 Table. Accuracy of indirect selection (average over 100 training-testing partitions) for best**
674 **phenotypic prediction (principal components (PCR), L1-penalized prediction (L1-Phen), RNDVI,**
675 **and GNDVI) and for best genotypic prediction (canonical SI (CSI), optimal PC-SI, L1-PSI, and L2-**
676 **PSI).** Each row contains results for each environment and time-point (DAS: days after sowing,
677 Stage: VEG= vegetative, GF=grain filling, MAT=maturity). Models with the same letter (within
678 each row) are not significantly different from each other ($\alpha=0.05$, ANOVA followed by Tuckey
679 test).

680

681