

Manuscript Number:	GIGA-D-17-00225	
Full Title:	Predicting plant biomass accumulation from image-derived parameters	
Article Type:	Research	
Funding Information:	Federal Agency for Agriculture and Food (15/12-13, 530-06.01-BiKo CHN)	Dr. Christian Klukas
	Robert Bosch Stiftung (32.5.8003.0116.0)	Dr. Christian Klukas
	Bundesministerium für Bildung und Forschung (0315958A and 031A053B)	Dr. Christian Klukas
	European Plant Phenotyping Network (284443)	Dr. Christian Klukas
Abstract:	<p>Background: Image-based high-throughput phenotyping technologies have been rapidly developed in plant science recently and they provide a great potential to gain more valuable information than traditionally destructive methods. Predicting plant biomass is regarded as a key purpose for plant breeders and ecologist. However, it is a great challenge to find a suitable model to predict plant biomass in the context of high-throughput phenotyping.</p> <p>Results: In the present study, we constructed several models to examine the quantitative relationship between image-based features and plant biomass accumulation. Our methodology has been applied to three consecutive barley (<i>Hordeum vulgare</i>) experiments with control and stress treatments. The results proved that plant biomass can be accurately predicted from image-based parameters using a random forest model. The high prediction accuracy based on this model, in particular the cross-experiment performance, will contribute to relieve the phenotyping bottleneck in biomass measurement in breeding applications. The relative contribution of individual features for predicting biomass was further quantified, revealing new insights into the phenotypic determinants of plant biomass outcome. What's more, the methods could also be used to determine the most important image-based features related to plant biomass accumulation, which would be promising for subsequent genetic mapping to uncover the genetic basis of biomass.</p> <p>Conclusions: We have developed quantitative models to accurately predict plant biomass accumulation from image data. We anticipate that the analysis results will be useful to advance our views of the phenotypic determinants of plant biomass outcome, and the statistical methods can be broadly used for other plant species.</p>	
Corresponding Author:	Dijun Chen GERMANY	
Corresponding Author Secondary Information:		
Corresponding Author's Institution:		
Corresponding Author's Secondary Institution:		
First Author:	Dijun Chen	
First Author Secondary Information:		
Order of Authors:	Dijun Chen	
	Rongli Shi	
	Jean-Michel Pape	

	Kerstin Neumann
	Daniel Arend
	Andreas Graner
	Ming Chen
	Christian Klukas
Order of Authors Secondary Information:	
Opposed Reviewers:	
Additional Information:	
Question	Response
Are you submitting this manuscript to a special series or article collection?	No
<p>Experimental design and statistics</p> <p>Full details of the experimental design and statistical methods used should be given in the Methods section, as detailed in our Minimum Standards Reporting Checklist. Information essential to interpreting the data presented should be made available in the figure legends.</p> <p>Have you included all the information requested in your manuscript?</p>	Yes
<p>Resources</p> <p>A description of all resources used, including antibodies, cell lines, animals and software tools, with enough information to allow them to be uniquely identified, should be included in the Methods section. Authors are strongly encouraged to cite Research Resource Identifiers (RRIDs) for antibodies, model organisms and tools, where possible.</p> <p>Have you included the information requested as detailed in our Minimum Standards Reporting Checklist?</p>	Yes
<p>Availability of data and materials</p> <p>All datasets and code on which the conclusions of the paper rely must be either included in your submission or deposited in publicly available repositories (where available and ethically appropriate), referencing such data using a unique identifier in the references and in the “Availability of Data and Materials”</p>	Yes

section of your manuscript.

Have you have met the above requirement as detailed in our [Minimum Standards Reporting Checklist](#)?

Predicting plant biomass accumulation from image-derived parameters

Dijun Chen^{1,*}, Rongli Shi¹, Jean-Michel Pape¹, Kerstin Neumann¹, Daniel Arend¹, Andreas Graner¹, Ming Chen², and Christian Klukas¹

¹*Leibniz Institute of Plant Genetics and Crop Plant Research (IPK), Corrensstrasse 3, 06466 Gatersleben, Germany.*

²*Department of Bioinformatics, College of Life Sciences, Zhejiang University, Hangzhou 310058, China.*

*Correspondence should be addressed to D.C. (chendijun2012@gmail.com)

Abstract

Background:

Image-based high-throughput phenotyping technologies have been rapidly developed in plant science recently and they provide a great potential to gain more valuable information than traditionally destructive methods. Predicting plant biomass is regarded as a key purpose for plant breeders and ecologist. However, it is a great challenge to find a suitable model to predict plant biomass in the context of high-throughput phenotyping.

Results:

In the present study, we constructed several models to examine the quantitative relationship between image-based features and plant biomass accumulation. Our methodology has been applied to three consecutive barley (*Hordeum vulgare*) experiments with control and stress treatments. The results proved that plant biomass can be accurately predicted from image-based parameters using a random forest model. The high prediction accuracy based on this model, in particular the cross-experiment performance, will contribute to relieve the phenotyping bottleneck in biomass measurement in breeding applications. The relative contribution of individual features for predicting biomass was further quantified, revealing new insights into the phenotypic determinants of plant biomass outcome. What's more, the methods could also be used to determine the most important image-based features related to plant biomass accumulation, which would be promising for subsequent genetic mapping to uncover the genetic basis of biomass.

Conclusions:

1
2
3
4 28 We have developed quantitative models to accurately predict plant biomass accumulation from image data.

5
6 29 We anticipate that the analysis results will be useful to advance our views of the phenotypic determinants of
7
8 30 plant biomass outcome, and the statistical methods can be broadly used for other plant species.
9

10 31 **Keywords:** Barley; High-throughput phenotyping; Phenomics; Biomass; Modeling.
11

12 32

14 33 **Introduction**

15
16
17 34 Biomass accumulation is an important indicator of crop final product and plant performance. It is thus
18
19 35 considered as a key trait in plant breeding, agriculture improvement and ecological applications. The
20
21 36 conventional approach of measuring plant biomass is very time consuming and labour intensive since plants
22
23 37 need to be harvested destructively to obtain the fresh or dry weight. Moreover, the destructive method makes
24
25 38 multiple measurements of the same plant over time impossible. With the development of new technology,
26
27 39 digital image analysis has been used more broadly in many fields, as well as in plant research. It allows faster
28
29 40 and more accurate plant phenotyping and has been proposed as an alternative way to infer plant biomass.
30

31 41

32
33 42 In recent years, plant biomass has been subject to intensive investigation by using high-throughput
34
35 43 phenotyping (HTP) approaches in both controlled growth chambers [1-4] and field environments [5-9],
36
37 44 demonstrating that the ability of imaging-based methods to infer plant biomass accumulation. On the other
38
39 45 hand, to produce reliable assessments, suitable model types needs to be established and model construction
40
41 46 requires integration of many components such as efficient mathematical analysis and representative data.

42
43 47 Although there are some developed models for predicting plant biomass, most of them have certain
44
45 48 limitations. For example, Golzarian *et al.* (2011) modelled the plant biomass (dry weight) in wheat (*Triticum*
46
47 49 *aestivum* L.) as a linear function of projected area, assuming plant density was constant. However, this
48
49 50 method under-estimated dry weight of salt stressed plants and over-estimated that of control plants. Even
50
51 51 though the authors argued that the bias was largely related to plant age and the model might be improved by
52
53 52 including the factor of plant age [2], the differences in plant density between stressed and control plants may
54
55 53 be caused by different physiological properties of plants rather than plant age. In another study, Busemeyer
56
57 54 *et al.* (2013) developed a calibrated biomass determination model for triticale (*x Triticosecale* Wittmack L.)
58
59 55 under field conditions based on multiple linear regression analysis of a diverse set of parameters, considering
60
61
62
63
64
65

1
2
3
4 56 both, the volume of the plants and their density. Indeed, this model largely improved the prediction accuracy
5
6 57 of the calibration models based on a single type of parameters and can precisely predict biomass accumulation
7
8 58 across environments [8]. Another concern is that the number of traits used in these studies were quite limited
9
10 59 and perhaps not representative enough. Therefore, a more effective and powerful model is needed to
11
12 60 overcome these limitations and to allow better utilization of the image-based plant features which are
13
14 61 obtained from non-invasive phenotyping approaches.
15
16 62

17
18 63 In this study, we present a general framework for investigating the relationships between plant biomass
19
20 64 (referred to as shoot biomass hereafter) and image-derived parameters. We applied a multitude of supervised
21
22 65 and unsupervised statistical methods to investigate different aspects of biomass determinants by a list of
23
24 66 representative phenotypic traits in three consecutive experiments in barley. The results showed that image-
25
26 67 based features can accurately predict plant biomass output and collectively reflect large proportions of the
27
28 68 variation in biomass accumulation. We elucidated the relative importance of different feature categories and
29
30 69 of individual features in prediction of biomass accumulation. The differences in the contribution of the image-
31
32 70 based features for prediction of two types of biomass measurements, fresh weight and dry weight were
33
34 71 compared as well. Furthermore, our models were tested for the possibility of predicting plant biomass in
35
36 72 different experiments with different treatments. As high-throughput plant phenotyping is a technique which
37
38 73 is becoming more and more widely used for automated phenotype in plant research, especially in plant
39
40 74 breeding, we anticipate that the methodologies proposed in this work will have various potential applications.
41
42
43 75

44 45 46 76 **Results**

47 48 77 **Development of statistical models for modelling plant biomass accumulation using image-based** 49 50 78 **features**

51
52 79 In the previous studies [10, 11], we have shown that a single phenotypic trait -- the three-dimensional digital
53
54 80 volume, which is a derived feature from projected side and top areas -- can be reasonably predictive to
55
56 81 estimate plant biomass accumulation. We expect that the predictive power could be improved when multiple
57
58 82 phenotypic traits are combined in a prediction model since plant biomass is determined not only by their
59
60
61
62
63
64
65

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

83 structural features but also by their density (physiological properties). To further investigate the relationship
84 between image-derived parameters and plant biomass accumulation, deep phenotyping data which contain
85 both structural (e.g., geometric traits) and physiological traits (e.g., plant moisture content; **Fig. 1, A and B**)
86 were analysed. Pot weights of the plants were not included for the analysis although they were weighed
87 regularly. It might reflect the growth tendency of the whole plants (shoots and roots) where herein we focused
88 mainly on shoots.

89
90 Models were constructed to quantify the ability of imaging-based features to statistically predict the biomass
91 accumulation. The models were developed by using four widely used machine-learning methods (**Fig. 1C**):
92 multivariate linear regression (MLR), multivariate adaptive regression splines (MARS), random forest (RF)
93 and support vector regression (SVR), which have extensively been used in accurate prediction of gene
94 expression [12-16] and DNA methylation levels [17-20]. We combined the biomass measurements (fresh
95 weight [FW] and/or dry weight [DW]) with image-based features and then divided them into a training data
96 set and a test data set. A model was trained on the training data set and has then been applied to the test data
97 set to predict the plant biomass. The relationship between plant biomass accumulation and image-based
98 features was assessed based on the criterion of the Pearson correlation coefficient (r) between the predicted
99 values and the actual values, or the coefficient of determination (R^2 ; the percentage of variance of biomass
100 explained by the model; **Fig. 1D**).

101
102 Our methodology was applied to three consecutive experiments (**Fig. 2A; Supplemental Table S1 and Data**
103 **S1**), which were designed to investigate vegetative biomass accumulation in response to two different
104 watering regimes under semi-controlled greenhouse conditions in a core set of barley cultivars by non-
105 invasive phenotyping [11, 21]. There were 312 plants with 18 genotypes for each experiment. Plants were
106 monitored using three types of sensors (visible, fluorescence [FLUO] and near-infrared [NIR]) in a
107 LemnaTec-Scanalyzer 3D imaging system. An extensive list of phenotypic traits ranging from geometric
108 (shape descriptors) to physiological properties (i.e., colour-, FLUO- and NIR-related traits) could be extracted
109 from the image data (**Fig. 1B**) using our image processing pipeline IAP [10]. A representative list of traits for
110 each plant in the last growth day were selected to test their ability to predict plant biomass.

1
2
3
4 111

5
6 112 **Coordinated patterns of plant-image-based profiles and their relation to plant biomass**

7
8 113 We extracted a list of representative and non-redundant phenotypic traits for each plant from image datasets
9
10 114 for each experiment (see **Materials and Methods**; **Fig. 1B**). In common for these experiments, overall thirty-
11
12 115 six high-quality traits which describe plant growth status in the last growth day were obtained. As a result,
13
14 116 each dataset was assigned a matrix whose elements were the signals of different features in different plants
15
16 117 (**Fig. 1C**). Unsupervised methods, such as hierarchical clustering (HCA; **Fig. 2B**) and principal component
17
18 118 analysis (PCA; **Fig. 2C**) were applied to these datasets. We found that plants from different experiments with
19
20 119 different treatments showed clearly distinct patterns of phenotypic profiles. For instance, stressed plants and
21
22 120 control plants were separated using PCA by their first principal component (PC1) and also by the top clusters
23
24 121 obtained in HCA, while plants from different experiments were distinguished by PC2 and PC3 in PCA or
25
26 122 subordinate clusters in HCA. Accordingly, it could be observed that biomass (e.g., FW) of plants from
27
28 123 different experiments with different treatments was significantly different (two-way ANOVA, p -value $< 2e$ -
29
30 124 16; **Fig. 2D**). The relationship was reflected by a dendrogram from cluster analysis based on the means of
31
32 125 FW over genotypes (**Fig. 2E**). Furthermore, the overall phenotypic patterns of these plants were similar to
33
34 126 their biomass output (**Fig. 2, B-E**), revealing that these image-based features were potential factors reflecting
35
36 127 the accumulation of plant biomass. We thus explored the relationship between the signals of these image-
37
38 128 based features and the level of plant biomass output. We calculated the correlation coefficients for each
39
40 129 dataset. The correlation patterns were consistent for different datasets and more than half of the features
41
42 130 revealed high correlation coefficients ($r > 0.5$; **Fig. 2F**). Interestingly, both structural features (such as digital
43
44 131 volume, projected area and the length of the projected plant area border) and density-related features (such
45
46 132 as NIR and FLUO intensities) were involved in the top ranked features.

47
48
49 133

50
51 134 **Relating image-based signals to plant biomass output**

52
53 135 The above analyses suggest that plant biomass can at least be partially inferred from image-based features.
54
55 136 To examine which model has the best performance and to select an appropriate model for biomass prediction,
56
57 137 we then applied our regression models (**Fig. 1C**) to predict plant biomass using image-based features. Our
58
59 138 analyses were focused on the first experiment (i.e., Exp 1), since the phenotypic traits of the corresponding
60
61
62
63
64
65

1
2
3
4 139 dataset have been intensively investigated in our previous study [11]. In this experiment, plant biomass was
5
6 140 quantified in two forms: FW and DW. We selected a collection of 45 image-derived parameters from this
7
8 141 dataset that were non-redundant and highly representative.
9

10 142
11
12 143 We next tried to predict FW (**Fig. 3A**) and DW (**Fig. 3C**) based on this set of image-derived features using
13
14 144 four different regression models. The models were respectively tested on control plants, stressed plants and
15
16 145 the whole set of plants. The performance of these models was compared and evaluated. Although the
17
18 146 performance of these models was roughly similar, RF, SVR and MARS methods had better performance than
19
20 147 the MLR method for prediction of both FW (**Fig. 3B**) and DW (**Fig. 3D**), implying a nonlinear relationship
21
22 148 between image-based phenotypic profiles and biomass output. The RF model largely outperformed other
23
24 149 models especially in predicting biomass of control plants, accounting for the most variance ($R^2 = 0.85$ for
25
26 150 FW and $R^2 = 0.62$ for DW; **Fig. 3, B and D**, left panels) and showed the best prediction accuracy (Pearson's
27
28 151 correlation $r = 0.93$ for FW and $r = 0.80$ for DW; **Fig. 3, B and D**, middle panels). The prediction accuracy
29
30 152 of our models (the correlation coefficients between the predicted biomass and the actual biomass) was also
31
32 153 compared with the ability of individual features to predict biomass (here, the “digital volume”; **Fig. 3, B and**
33
34 154 **D**, middle panels). It was found that our models generally showed better prediction power than the single
35
36 155 digital volume-based prediction, indicating that additional features improved the predictive power. In this
37
38 156 study, we focused on the results from the RF method in the rest of analysis, although results from different
39
40 157 methods were highly consistent and led to the same conclusions.
41
42
43 158

159 **Relative importance of different image-based features for predicting plant biomass**

44
45
46
47 160 As mentioned above, the image-based features could be classified broadly into four categories: plant structure
48
49 161 properties, colour-related features, NIR signals, and FLUO-based traits (**Fig. 1B**). The last three types of
50
51 162 features reflect plant physiological properties and can be considered as plant density-related traits and are
52
53 163 thus related to their fresh or dry matter content. For each individual feature or each type of features, we
54
55 164 constructed a degenerate model for biomass prediction using the corresponding feature(s) as the predictor(s).
56
57 165 We compared the capability of each individual or type of feature for predicting biomass accumulation in the
58
59 166 first experiment (i.e., experiment 1). Geometric features showed the most predictive power among the four
60
61
62
63
64
65

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

167 categories for prediction of both FW and DW, but were slightly less predictive than all features in a full model
168 (**Fig. 4, A and B**). Strikingly, the predictability of other types of features (such as colour-related and FLUO-
169 based traits) was substantial, indicating that these traits may act as unforeseen factors in biomass prediction.
170 In addition, the NIR-based features showed higher predictive capability for FW than for DW in control and
171 stressed plants, revealing NIR signals were import factors in determining FW accumulation.

172
173 Next, we investigated the relative importance (RI) of each feature for predicting biomass using a full model
174 in the whole set of plants (i.e., “control + stressed plants”; **Fig. 4, C and D**, upper panels). In a RF model, the
175 RI of a feature is calculated as the increase of prediction error (%IncMSE) when phenotypic data for this
176 feature is permuted [22], and thus indicates the contribution of the feature after considering its intercorrelation
177 in a model. We found that the top ten most important features in the full model for predicting FW and DW
178 included both structure and density-related traits. As expected, projected area (from side or top view) and
179 digital volume were the top ranked features, which have individually been considered as proxies of shoot
180 biomass in previous studies [2, 11, 23-29].

181
182 In principle, we would expect that highly important features in the full model would be related to a high
183 predictive power in a degenerate model. Surprisingly, there was no clear correlation observed between the
184 feature importance and their predictive power (**Fig. 4, C and D**). For example, several colour-related and
185 NIR-based features which were in the top ten list of the most important features revealed insubstantial
186 predictive power in individual models. This observation implies that the relation of the underlying biomass
187 determinants is extremely complex and not a linear combinations of the investigated features.

188
189 Furthermore, we compared the relative importance of each feature in predicting FW and DW (**Fig. 4E**).
190 Although a positive correlation ($r = 0.88$) between the feature importance for FW and DW could be observed,
191 several features showed large differences in their ability to interpret FW or DW, including “nir.intensity”
192 (derived from side view images), “compactness.01” (top), “hull.pc1” (top), “leaf.count” (side),
193 “hsv.h.average” (top) and “lab.a.mean” (top). For instance, NIR intensity and plant compactness (top view)
194 may be important for predicting FW but not for DW. We also performed the above analyses by using only

1
2
3
4 195 control (**Supplemental Fig. S1**) or stressed plants (**Supplemental Fig. S2**), respectively. We found that the
5
6 196 patterns of feature importance were distinct between these two groups of plants. For example, NIR intensity
7
8 197 was ranked as the top fifth feature for predicting FW for stressed plants but was not substantially important
9
10 198 for control plants. These findings suggest that there are differences in underlying plant biomass determinants
11
12 199 in these kinds of treatment situations that are also reflected by their image-based phenotypic traits.
13

14 200

15
16 201 **Image-based features are predictive of plant biomass across experiments with similar conditions or**
17
18 202 **treatments**

19
20 203 In order to explore whether our models were generalizable across different experiments, we applied our
21
22 204 models trained in one experiment to predict biomass (herein FW) in other experiments using a common set
23
24 205 of features. Examples of such cross-experiment predictions are shown in **Figure 5A**. We tested and illustrated
25
26 206 all possibilities for cross prediction using the whole set of plants in the corresponding experiment. In general,
27
28 207 the prediction accuracy within individual experiments remained high ($r > 0.97$ and $R^2 > 0.93$ for all three
29
30 208 experiments; **Fig. 5B**), revealing that our models were effectively predicting plant biomass based on image-
31
32 209 derived feature signals among different experiments. Moreover, the prediction accuracy for cross-experiment
33
34 210 prediction was still relatively high, with $r > 0.81$ and $R^2 > 0.65$, implying that our models accurately
35
36 211 captured the relationships among the various image-based features. However, we observed that the third
37
38 212 experiment had relative weaker correlations with the other two experiments for predicting biomass, while the
39
40 213 first two experiments showed strong correlations or even nearly identical results when being compared with
41
42 214 each other (**Fig. 5A**). This might be related mainly to seasonal (temperature and illumination) differences
43
44 215 which caused different plants behaviours, namely lower biomass for both control and stressed plants, in
45
46 216 experiment 3 as explained by the authors [21]. This suggests that different plant growth conditions might
47
48 217 cause some variation for cross-experiment prediction.
49

50 218

51
52
53 219 At the same time, we tested cross predictability of our models using treatment-specific data in the experiments
54
55 220 (**Fig. 6**). Similar results were obtained as above using the whole dataset (**Fig. 5B**). The weak predictive power
56
57 221 for cross-prediction involving control plants from the third experiment was most clearly observable in the
58
59 222 low accuracy in the biomass prediction of this particular subset of plants. Generally, control and stressed
60
61
62
63
64
65

1
2
3
4 223 plants were found to have very weak predictive power when related to each other (**Fig. 6**), as also supported
5
6 224 by the distinct patterns of relative feature importance between these two plant groups (**Supplemental Figs.**
7
8 225 **S1 and S2**). For each experiment, the prediction accuracy was higher for stressed plants compared to control
9
10 226 plants. This might resulted from the imaging analysis process. Relatively small plants, stressed plants in this
11
12 227 case, would gain more clear images due to less overlapping or less area of out range. Therefore, image quality
13
14 228 would be an important variation source for our modelling and should be taking into consideration for any
15
16 229 application.
17
18 230

21 231 **Discussion**

22
23
24 232 Biomass is a complex but important trait in functional ecology and agronomy for studying plant growth, crop
25
26 233 productive potential and plant regeneration capabilities. Many different techniques, either destructive or non-
27
28 234 destructive, have been used to estimate biomass [30]. Compared with the traditional destructive methods for
29
30 235 measuring biomass, non-destructive imaging methods provide a faster, more accurate approach for plant
31
32 236 phenotyping. In recent years, more and more high-throughput plant phenotyping platforms have been set up
33
34 237 and applied worldwide. Accordingly, it becomes a current challenge to establish models utilizing the big
35
36 238 datasets gained from high-throughput imaging systems. Although previous attempts have been made to
37
38 239 estimate plant biomass from image data, most of these studies consider only a single image-based feature or
39
40 240 very few features in their models which are often linear-based, ignoring the fact that the phenotypic
41
42 241 components underlying biomass accumulation are presumably complex. Accurately predicting biomass from
43
44 242 image data requires efficient mathematical models as well as representative image-derived features.
45
46 243

47
48 244 In this study, we have presented a systematic analysis of relationships between plant biomass accumulation
49
50 245 and image-derived signals, to confirm the assumption that biomass can be accurately predicted from image-
51
52 246 based parameters. We built a random forest model of biomass accumulation using a comprehensive list of
53
54 247 representative image-based features. The comparison between a random forest model and alternative
55
56 248 regression models indicated that the RF model outperforms other models in terms of (1) better predictive
57
58 249 power – especially in comparison with the linear model, confirming the complex phenotypic architecture of
60
61
62
63
64
65

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

250 biomass, (2) better outperformance than a single-feature-prediction model – arguing the complex phenotypic
251 makeup of biomass, and (3) feasible biological interpretability – the ability to readily extract information
252 about the importance of each feature in prediction. The high prediction accuracy based on this model, in
253 particular the cross-experiment performance, is promising to relieve the phenotyping bottleneck in biomass
254 measurement in breeding applications. For example, based on an established small reference dataset which
255 is used to train a RF model, it is possible to predict biomass in several large plant populations within one
256 experiment or across several experiments using image data by taking advantage of high-throughput
257 phenotyping technologies. Alternatively, the model can be trained from a much larger reference panel of
258 plants that are grown in diverse environmental conditions which is then applied to a diverse set of experiments.
259 The first evidence for this notion is the observation that our model showed more predictive power in plants
260 with two treatments than with a single treatment (**Fig. 3, B and D**). Indeed, when applying our model to the
261 combined dataset from all the three experiments, we found the prediction accuracy remains very high ($R^2 =$
262 0.96 and $r = 0.98$, average values from ten times of ten-fold cross-validation). To keep the high prediction
263 accuracy in other application, there are some points should be take caution. Considering the environmental
264 effects on biomass accumulation, the application of our model will require the testing experiments showing
265 similar conducted conditions with that of the reference experiments. This means the plant cultivation
266 conditions should be standardized and any noise which might lower image quality should be avoided.
267 Another approach to improve applicability of models, which could not be tested in this study, would be to
268 improve the data base for the training, by acquiring data from additional environment sensors. Temperature,
269 humidity, and illumination data would certainly help to explain differences in the growth patterns among
270 experiments, performed in different growth seasons. To this end, we expect that our approach is extensible
271 by incorporating such sensor data in the data matrices. Furthermore, our results can provide suggestive hints
272 for biologists to setup phenotyping infrastructures for investigation of plant biomass. For instance, a visible
273 light imaging system would be sufficient to accurately predict fresh weight based on the observation that
274 geometric features alone show high prediction accuracy (**Fig. 4A**). However, to investigate dry weight, it
275 would be helpful to include an additional near-infrared camera system under normal growth conditions and
276 an additional fluorescence camera system under drought stress conditions (**Fig. 4B**).

277

1
2
3
4 278 In contrast to previous studies [1-4, 23-29], in which biomass was investigated using only single image-
5
6 279 derived parameter (such as projected area) or several geometric parameters, our analyses extended these
7
8 280 studies by incorporating more representative features that cover both structural and physiological-related
9
10 281 properties into a more sophistic model. Although the predictive power of our model is roughly higher than
11
12 282 that of single feature-based prediction, such as the digital volume (**Fig. 3**) [11], our model also reveals the
13
14 283 relative contribution of individual feature in prediction of biomass. The information regarding the importance
15
16 284 of each feature will offer new insights into the phenotypic determinants of plant biomass outcome.
17
18 285 Interestingly, we found that several top ranked features, such as digital volume and NIR intensity, showed
19
20 286 genetic correlations with biomass of fresh weight (**Fig. 4C**) [11], implying these top ranked features may
21
22 287 represent the main “phenotypic components” of biomass outcome and can be further used to dissect genetic
23
24 288 components underlying biomass accumulation. As image-based high-throughput phenotyping in plants
25
26 289 developed mainly in recent years and therefore few corresponding modelling studies have been performed,
27
28 290 we believe that our model could be further improved when new types of cameras and/or newly defined
29
30 291 features are available.

31
32 292
33
34 293 In summary, we have developed a quantitative model for dissecting the phenotypic components of biomass
35
36 294 accumulation based on image data. Apart from predicting biomass outcome, the methods can be used to
37
38 295 determine the most important image-based features related to plant biomass accumulation, which are
39
40 296 promising for subsequent genetic mapping to uncover the genetic basis of biomass.

41
42 297

43 44 45 298 **Potential Implications**

46
47
48 299 We anticipate that the analysis results will be useful to advance our views of the phenotypic determinants of
49
50 300 plant biomass outcome, and the statistical methods can be broadly used for other plant species and therefore
51
52 301 assist plant breeding in the context of phenomics.

53
54 302
55
56
57
58
59
60
61
62
63
64
65

303 **Materials and Methods**

304 **Germplasm and experiments**

305 Barley plant image data were obtained as described previously [11, 21]. Briefly, a core set of 16 two-rowed
306 spring barley cultivars (*Hordeum vulgare* L.) and two parental cultivars of a double haploid (DH) were
307 monitored for vegetative biomass accumulation. Three independent experiments with identical setup were
308 performed in a (semi-) controlled greenhouse at IPK by using the automated phenotyping and imaging
309 platform LemnaTec-Scanalyzer 3D. Experiments were performed consecutively from May to November
310 2011 over a period of 58 days each (**Supplemental Table S1**). The greenhouse setup enabled sowing for the
311 next experiment already 2 days before the old experiment ended. For this, new pots were placed in the middle
312 of the greenhouse, while the old experiment was still on the conveyer belts.

313

314 Each experiment consisted of two treatments: well-watered (control treatment) and water limited (drought
315 stress treatment). In each treatment, nine plants per core set cultivar as well as six plants per DH parent were
316 tested. This resulted in a total of 312 plants per experiment, corresponding to the maximal capacity of the
317 phenotyping platform. Watering and imaging were performed daily. Drought stress was imposed by
318 intercepting water supply from 27 days after sowing (DAS 27) to DAS 44. Stressed plants were re-watered
319 at DAS 45. In total, for each of the experiments about 100 GB of raw (image) data was accumulated. At the
320 end of experiments (DAS 58), plants were harvested to measure above-ground biomass in form of plant fresh
321 weight (FW; for all experiments) and/or dry weight (DW; for experiment 1).

322

323 **Image analysis**

324 Image datasets were processed by the barley analysis pipelines in the IAP software (version v1.1.2) [10].
325 Analysed results were exported in the csv file format via IAP functionalities, which can be used for further
326 data inspection. The result table includes columns for different phenotypic traits and rows as plants are
327 imaged over time. The corresponding metadata is included in the result table as well.

328

329 Each plant was characterized by a set of phenotypic traits also referred to as features, which were grouped
330 into four categories: geometric features, fluorescence-related (FLUO-related) features, colour-related

1
2
3
4 331 features and near-infrared-related (NIR-related) features. These traits were defined by considering image
5
6 332 information from different cameras (visible light, fluorescence and near infrared) and imaging views (side
7
8 333 and top views). See the IAP online documentation (<http://iapg2p.sourceforge.net/documentation.pdf>) for
9
10 334 details about trait definition.

11
12 335

13 14 336 **Feature selection**

15
16 337 Feature selection was performed with the same procedure as described in [11]. We applied the feature
17
18 338 selection technique to each dataset. Generally, we captured almost identical subset features from different
19
20 339 datasets. We manually added several representative traits due to removal by variance inflation factors. For
21
22 340 example, the digital volume and projected area are highly correlated with each other but we kept both of
23
24 341 them, because we would investigate the predictive power of both features. Moreover, the regression models
25
26 342 we used are insensitive to collinear features. We thus kept as much representative features as possible. To
27
28 343 apply the prediction models among different datasets, a common set of features supported by all the datasets
29
30 344 was used.

31
32 345

33 34 346 **Data transformation**

35
36 347 Each plant can be presented by a representative list of phenotypic traits, resulting in a matrix $X_{n \times m}$ for each
37
38 348 experiment, where n is the number of plants and m is the number of phenotypic traits. Missing values were
39
40 349 filled by mean values of other replicated plants. To make the image-derived parameters from diverse sources
41
42 350 comparable, we normalized the columns of X by dividing the values with the maximum value of each column
43
44 351 across all plants. Plants with empty values of manual measurements (FW and DW) were discarded for
45
46 352 analysis. These transformed data sets were subjected to regression models.

47
48 353

49 50 354 **Hierarchical clustering analysis and PCA**

51
52 355 Hierarchical clustering analysis (HCA) and principle component analysis (PCA) were performed on the
53
54 356 transformed data matrix $X_{n \times m}$ in the same way as described in [11]. We also performed HCA using the
55
56 357 genotype-level mean value of FW data to check the similarity of overall plant growth patterns in different
57
58 358 experiments.

59
60
61
62
63
64
65

1
2
3
4 359

5
6 360 **Models for predicting plant biomass**

7
8 361 To understand the underlying relationship between image-derived parameters and the accumulated biomass
9
10 362 (such as FW and DW), we constructed predictive models based on four different machine-learning methods:
11
12 363 multivariate linear regression (MLR), multivariate adaptive regression splines (MARS), random forest (RF)
13
14 364 and support vector regression (SVR). In these models, the normalized phenotypic profile matrices $X_{n \times m}$ for
15
16 365 a representative list of phenotypic traits were used as predictors (explanatory variables) and the measured
17
18 366 DW/FW as the response variable Y .

19
20 367

21
22 368 All these models were implemented in R (<http://www.r-project.org/>; release 2.15.2). To assess the relative
23
24 369 contribution of each phenotypic trait to predicting the biomass. We also calculated the relative feature
25
26 370 importance for each model. Specifically, for the MLR model, we used the “lm” function in the base
27
28 371 installation packages. The relative importance of predictor variables in the MLR model was estimated by a
29
30 372 heuristic method [31] which decomposes the proportionate contribution of each predictor variable to R^2 . For
31
32 373 MARS, we used the “earth” function in the *earth* R package. The “number of subsets (nsubsets)” criterion
33
34 374 (counting the number of model subsets that include the variable) was used to calculate the variables feature
35
36 375 importance, which is implemented in the “evimp” function. For the RF model, we used the *randomForest* R
37
38 376 package which implements Breiman's random forest algorithm [22]. We chose the “%IncMSE” (increase of
39
40 377 mean squared error) to represent the criteria of relative importance measure. For SVR, we utilized the *e1071*
41
42 378 R package which provides functionalities to use the *libsvm* library [32]. The absolute values of the
43
44 379 coefficients of the normal vector to the “optimal” hyperplane can be considered as the relative importance of
45
46 380 each predictor variable contributing to regression [33, 34].

47
48
49 381

50
51 382 **Evaluation of the prediction models**

52
53 383 To evaluate the performance of the predictive models, we adopted a 10-fold cross-validation strategy to check
54
55 384 the prediction power of each regression model. Specifically, each dataset was randomly divided into a training
56
57 385 set (90% of plants) and a testing set (10% of plants). We trained a model on the training data and then applied
58
59 386 it to predict biomass for the testing data. Afterwards, the predicted biomass in the testing set was compared

60
61
62
63
64
65

1
2
3
4 387 with the manually measured biomass. The predictive accuracy of the model can be measured by

- 5
6 388 1) the Pearson correlation coefficient (PCC; r) between the predicted values and the observed values;
7
8 389 2) the coefficient of determination (R^2) which equals to the fraction of variance of biomass explained
9
10 390 by the model, defined as

11
12
13 391
$$R^2 = 1 - \frac{SS_{res}}{SS_{tot}} = 1 - \frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{\sum_{i=1}^n (y_i - \bar{y})^2}$$

14
15
16 392 where SS_{res} and SS_{tot} are the sum of squares for residuals and the total sum of squares, respectively, \hat{y}_i the
17
18 393 predicted and y_i the observed biomass of the i th plant, \bar{y} is the mean value of the observed biomass; and

- 19
20 394 3) the root mean squared relative error of cross-validation, defined as

21
22
23
24 395
$$\text{RMSRE} = \sqrt{\frac{\sum_{i=1}^s \left(\frac{y_i - \hat{y}_i}{y_i} \right)^2}{s}}$$

25
26
27 396 where s denotes the sample size of the testing dataset.

28
29 397 We repeated the cross-validation procedure ten times. The mean and standard deviation of the resulting R^2
30
31 398 and RMSRE values were calculated across runs.

32
33 399
34
35 400 To illustrate the broad utility of our methods across seasons (thus different growth environments) and
36
37 401 treatments (e.g., control versus drought stress) in the same season, we applied the models in different contexts
38
39 402 with cohort validation. Specifically, we trained the biomass prediction models under one specific context and
40
41 403 predicted biomass in another different context and *vice versa*. The predictive accuracy of the model was
42
43 404 evaluated based on the measures R^2 and RMSRE as described above. Furthermore, the predictive power was
44
45 405 reflected by the bias μ between the predicted and observed values, defined as

46
47
48 406
$$\mu = \frac{1}{n} \cdot \sum_{i=1}^n \frac{\hat{y}_i - y_i}{y_i}$$

49
50
51 407 where n denotes the sample size of the dataset. This bias indicates over- ($\mu > 0$) or under-estimation ($\mu < 0$)
52
53 408 of biomass.

54
55 409

56
57 410 **Availability of source code and requirements**

- 58
59 411
 - Project name: Modeling of plant biomass accumulation with HTP data
- 60
61
62
63
64
65

- 1
2
3
4 412 • Project home page: <https://github.com/htpmod/HTPmod>
5
6 413 • Operating system(s): Windows, Linux and Mac OS.
7
8 414 • Programming language: R
9
10 415 • License: open source under GNU GPL v3.0.
11

12 416

14 417 **Availability of supporting data and materials**

16 418 The raw data sets supporting the results of this article are available in the PGP repository [35] under XXXX
17
18 419 (please use the following links for review: [https://doi.ipk-gatersleben.de/DOI/d1378e80-f25d-4f7f-99c9-](https://doi.ipk-gatersleben.de/DOI/d1378e80-f25d-4f7f-99c9-c623fb2626c9/8b3fcde6-3ea7-4a3f-affc-9bf5ac289846/2/1847940088)
20 420 [c623fb2626c9/8b3fcde6-3ea7-4a3f-affc-9bf5ac289846/2/1847940088](https://doi.ipk-gatersleben.de/DOI/d1378e80-f25d-4f7f-99c9-c623fb2626c9/8b3fcde6-3ea7-4a3f-affc-9bf5ac289846/2/1847940088), [https://doi.ipk-](https://doi.ipk-gatersleben.de/DOI/f190fef8-e009-4a1f-bce2-a830fe561d42/71943812-eff1-4c89-980b-20dc37e35b97/2/1847940088)
22 421 [gatersleben.de/DOI/f190fef8-e009-4a1f-bce2-a830fe561d42/71943812-eff1-4c89-980b-](https://doi.ipk-gatersleben.de/DOI/f190fef8-e009-4a1f-bce2-a830fe561d42/71943812-eff1-4c89-980b-20dc37e35b97/2/1847940088)
24 422 [20dc37e35b97/2/1847940088](https://doi.ipk-gatersleben.de/DOI/f190fef8-e009-4a1f-bce2-a830fe561d42/71943812-eff1-4c89-980b-20dc37e35b97/2/1847940088), and
26 423 [https://doi.ipk-gatersleben.de/DOI/d23537d7-8d78-4f00-8bc2-018eca2d83d3/23ede2c4-a44e-44b7-9009-](https://doi.ipk-gatersleben.de/DOI/d23537d7-8d78-4f00-8bc2-018eca2d83d3/23ede2c4-a44e-44b7-9009-1f6d7514aa5c/2/1847940088)
28 424 [1f6d7514aa5c/2/1847940088](https://doi.ipk-gatersleben.de/DOI/d23537d7-8d78-4f00-8bc2-018eca2d83d3/23ede2c4-a44e-44b7-9009-1f6d7514aa5c/2/1847940088)), according to the ISA-Tab format and the recommendations of the MIAPPE
30 425 (Minimum Information About a Plant Phenotyping Experiment) standard [36]. The selected data for
32 426 modelling are available in the **Supplemental Data S1**.
34
35 427

38 428 **Declarations**

40 429 **List of abbreviations**

42 430 DAS: Days After Sowing
44 431 DW: Dry Weight
46 432 FLUO: Fluorescence
48 433 FW: Fresh Weight
50 434 HCA: Hierarchical Clustering Analysis
52 435 HTP: High-Throughput Phenotyping
54 436 MLR: Multivariate Linear Regression
56 437 MARS: Multivariate Adaptive Regression Splines
58 438 NIR: Near-Infrared

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

439 PCA: Principal Component Analysis

440 PCC: Pearson Correlation Coefficient

441 RF: Random Forest

442 RMSRE: Root Mean Squared Relative Error

443 SVR: Support Vector Regression

444

445 **Consent for publication**

446 Not applicable.

447

448 **Funding**

449 This work was supported by the Leibniz Institute of Plant Genetics and Crop Plant Research (IPK), the Robert
450 Bosch Stiftung (32.5.8003.0116.0) and the Federal Agency for Agriculture and Food (BEL, 15/12-13, 530-
451 06.01-BiKo CHN) and the Federal Ministry of Education and Research (BMBF, 0315958A and 031A053B).

452 This research was furthermore enabled with support of the European Plant Phenotyping Network (EPPN,
453 grant agreement no. 284443) funded by the FP7 Research Infrastructures Programme of the European Union.

454

455 **Competing interests**

456 The authors declare that they have no competing interests.

457

458 **Author contributions**

459 D.C. designed the research. C.K. and M.C. supervised the project. K.N. and G.A. performed the LemnaTec
460 experiments. D.A. created the ISA-Tab formatted description and uploaded data records in the PGP repository.

461 J.M.P. and C.K. analyzed image data. D.C. implemented the methods, analyzed data, interpreted the results,
462 and wrote the manuscript with contribution from R.S.. All authors read and approved the final version of the
463 article.

464

465 **Acknowledgements**

466 We would like to thank Ingo Mücke for his management of the LemnaTec system operations. We thank

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

467 Michael Ulrich for performing software tests and helping in data analysis.

468

469 **Figure Legends**

470 **Figure1.** Modeling pipeline for predicting plant biomass accumulation based on image-derived parameters.

471 (A) Input data, including high-throughput image data and manually measured biomass data. Plants were
472 phenotyped using various cameras such as visible (or color), fluorescence (FLUO) and near-infrared (NIR)
473 sensors. Image analysis was performed with IAP software [10] for feature extraction. The same plants were
474 harvested and measured at the end of growth. Generally, two types of biomass were measured: fresh weight
475 (FW) and dry weight (DW). (B) Trait processing. All the phenotypic traits were grouped into four categories:
476 geometric, color-related, FLUO-related and NIR-related traits. Phenotypic data were subjected to quality
477 check to remove low-quality data. (C) Each plant was described by a list of traits, resulting in a predictor
478 matrix whose rows represent plants and columns represent image-based traits. This matrix was used to
479 predicted plant biomass accumulation by MLR (multivariate linear regression), MARS (multivariate adaptive
480 regression splines), RF (random forest) and SVR (support vector regression) models. The right panel
481 represents the schema of model validation. In the first schema, a dataset (Dataset 1) was divided into training
482 set and testing set in a ten-fold cross-validation manner. In the second schema, the whole of one dataset
483 (Dataset 1) was used for training and another dataset (Dataset 2) was used for testing. (D) Model selection,
484 evaluation and result interpretation. The correlation of the predicted values and measured values was used to
485 assess the overall performance of the model.

486
487 **Figure 2.** Predictability of image-based traits to plant biomass.

488 (A) Schema depicting three consecutive high-throughput phenotyping experiments in barley. Plants in each
489 experiment were harvested for biomass measurements: fresh weight (FW; for all experiments) and dry weight
490 (DW; only for experiment 1). (B) Heatmap of Pearson's correlations between plants. Pearson's correlation
491 coefficient (PCC) was calculated based on image-derived traits. Cluster dendrograms for experiments (left)
492 and treatments (top) are shown. (C) Scatter plots showing projections of the top four Principal components
493 (PCs) based on PCA of image-based data. The component scores (shown in points) are colored and shaped
494 according to the experiments (as legend listed in the box). The component loading vectors (represented in
495 lines) of all traits (as colored according to their categories) were superimposed proportionally to their
496 contribution. (D) Boxplot showing the distribution of FW across different experiments. (E) A dendrogram

1
2
3
4 497 from cluster analysis based on the means of FW data over genotypes. **(F)** Pearson's correlation (mean values
5
6 498 in the three datasets) between image-based traits and FW. Traits with the largest mean correlations values are
7
8 499 labeled: 1 -- sum of leaf length (side view), 2 -- sum of FLUO intensity (side), 3 -- plant area border length
9
10 500 (side), 4 -- sum of NIR intensity (top), 5 -- sum of FLUO intensity (top), 6 -- projected area (top), 7 --
11
12 501 projected area (side) and 8 -- digital volume.

13
14 502
15
16 503 **Figure 3.** Quantitative relationship between image-based features and plant biomass.

17
18 504 **(A)** and **(C)** Scatter plots of manually measured plant biomass (fresh weight [FW] and dry weight [DW])
19
20 505 versus predicted biomass values using four prediction models: multivariate linear regression (MLR),
21
22 506 multivariate adaptive regression splines (MARS), random forest (RF) and support vector regression (SVR).
23
24 507 The red line indicates the expected prediction ($y = x$). The quantitative relationship between image-based
25
26 508 features and biomass was evaluated by Pearson's correlation coefficient (PCC r and its corresponding p -
27
28 509 value), RMSRE (root mean squared relative error) and the percentage of variance explained by the models
29
30 510 (the coefficient of determination R^2). **(B)** and **(D)** Summary of the predictive power of each regression model.
31
32 511 The results were based on ten-fold cross-validation with ten trials. Models were evaluated based on control
33
34 512 plants, stressed plants and the whole set of plants. The solid lines in the middle panel represent PCC between
35
36 513 digital volume and biomass for specific datasets.

37
38
39 514
40
41 515 **Figure 4.** The relative importance of image-based features in prediction of plant biomass.

42
43 516 The capabilities of different types of image-based features to predict plant biomass based on evaluation of
44
45 517 either fresh weight (FW) **(A)** or dry weight (DW) **(B)**. The overall predictive accuracies of each type of
46
47 518 features are indicated. Grey bar denote the predictive accuracy using all features. The relative importance of
48
49 519 each feature in the Random Forest model (upper panel) and the predictive accuracy of each individual feature
50
51 520 as the single predictor (lower panel) based on investigation of either FW **(C)** or DW **(D)**. The calculation was
52
53 521 based on the whole set of plants (control and stressed plants). Note that feature labels are shared in the upper
54
55 522 and lower panels. Features are shown in numbers as ordered by their names. The three features highlighted
56
57 523 in the red dash box are digital volume, projected side area and projected top area. **(E)** Comparison of the
58
59 524 relative importance of features in prediction of FW and DW. The top six most different features are

1
2
3
4 525 highlighted and labeled.

5
6 526

7
8 527 **Figure 5.** Comparison of prediction accuracy across different experiments.

9
10 528 (A) Application of the model learned from one experiment to other experiments. (B) Boxplots of coefficient
11 529 determination (R^2 , left) Pearson's correlation coefficients (r , middle) and the root mean squared relative error
12 530 (RMSRE, right) for different comparisons. “Within” denotes a model trained and tested on data from the same
13 531 dataset with specific treatments (control, stress or both), and “Cross” represents a model trained on one
14 532 dataset and tested on another dataset. “Control → stress” denotes a model trained on data with control
15 533 treatment and tested on data with stress treatment, and vice versa for “stress → control”.

16 534

17 535 **Figure 6.** Comparison of prediction accuracy across different treatments. Refer to **Figure 5A** for legend.

18 536

19 537 **Supplemental Data**

20 538 The following supplemental materials are available.

21 539 **Supplemental Figure S1.** The relative importance of image-based features in prediction of biomass in
22 540 control plants. Refer to **Figure 4** for legend. The calculation was based on control plants.

23 541 **Supplemental Figure S2.** The relative importance of image-based features in prediction of biomass in
24 542 stressed plants. Refer to **Figure 4** for legend. The calculation was based on stressed plants.

25 543

26 544 **Supplemental Table S1.** Overview of three high-throughput phenotyping experiments in barley.

27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46

Experiment	#plants/#genotypes ¹	Date of sowing	Date of harvesting	Biomass ²
Exp. 1 (1121KN)	310/18	27.05.2011	24.07.2011	FW & DW
Exp. 2 (1130KN)	310/18	22.07.2011	18.09.2011	FW
Exp. 3 (1137KN)	309/18	16.09.2011	13.11.2011	FW & DW

47 545 ¹Number of plants or genotypes used in analysis (filtered data).

48 546 ²Types of biomass measurement. FW: fresh weight; DW: dry weight.

49 547
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

1
2
3
4 548 **Supplemental Data S1.** Manual data and image-derived data in the three experiments.
5
6

7 549
8
9

10 550 **References**
11

- 12 551 1. Tackenberg O: **A new method for non-destructive measurement of biomass, growth rates,**
13 552 **vertical biomass distribution and dry matter content based on digital image analysis.** *Annals of*
14 553 *botany* 2007, **99**(4):777-783.
15 554 2. Golzarian MR, Frick RA, Rajendran K, Berger B, Roy S, Tester M, Lun DS: **Accurate inference of shoot**
16 555 **biomass from high-throughput images of cereal plants.** *Plant methods* 2011, **7**:2-2.
17 556 3. Feng H, Jiang N, Huang C, Fang W, Yang W, Chen G, Xiong L, Liu Q: **A hyperspectral imaging system**
18 557 **for an accurate prediction of the above-ground biomass of individual rice plants.** *Review of*
19 558 *Scientific Instruments* 2013, **84**(9):095107-095107.
20 559 4. Neumann K, Zhao Y, Chu J, Keilwagen J, Reif JC, Kilian B, Graner A: **Genetic architecture and**
21 560 **temporal patterns of biomass accumulation in spring barley revealed by image analysis.** *BMC*
22 561 *plant biology* 2017, **17**(1):137.
23 562 5. Ehlert D, Horn H-J, Adamek R: **Measuring crop biomass density by laser triangulation.** *Computers*
24 563 *and electronics in agriculture* 2008, **61**(2):117-125.
25 564 6. Ehlert D, Heisig M, Adamek R: **Suitability of a laser rangefinder to characterize winter wheat.**
26 565 *Precision Agric* 2010, **11**(6):650-663.
27 566 7. Erdle K, Mistele B, Schmidhalter U: **Comparison of active and passive spectral sensors in**
28 567 **discriminating biomass parameters and nitrogen status in wheat cultivars.** *Field Crops Research*
29 568 2011, **124**(1):74-84.
30 569 8. Busemeyer L, Ruckelshausen A, Moller K, Melchinger AE, Alheit KV, Maurer HP, Hahn V, Weissmann
31 570 EA, Reif JC, Wurschum T: **Precision phenotyping of biomass accumulation in triticale reveals**
32 571 **temporal genetic patterns of regulation.** *Scientific reports* 2013, **3**:2442-2442.
33 572 9. Cao Q, Miao Y, Wang H, Huang S, Cheng S, Khosla R, Jiang R: **Non-destructive estimation of rice**
34 573 **plant nitrogen status with Crop Circle multispectral active canopy sensor.** *Field Crops Research*
35 574 2013, **154**:133-144.
36 575 10. Klukas C, Chen D, Pape JM: **Integrated Analysis Platform: An Open-Source Information System for**
37 576 **High-Throughput Plant Phenotyping.** *Plant Physiol* 2014, **165**(2):506-518.
38 577 11. Chen D, Neumann K, Friedel S, Kilian B, Chen M, Altmann T, Klukas C: **Dissecting the phenotypic**
39 578 **components of crop plant growth and drought responses based on high-throughput image**
40 579 **analysis.** *Plant Cell* 2014, **26**:4636-4655.
41 580 12. Cheng C, Alexander R, Min R, Leng J, Yip KY, Rozowsky J, Yan K-K, Dong X, Djebali S, Ruan Y *et al*:
42 581 **Understanding transcriptional regulation by integrative analysis of transcription factor binding**
43 582 **data.** *Genome research* 2012, **22**(9):1658-1667.
44 583 13. Cheng C, Gerstein M: **Modeling the relative relationship of transcription factor binding and**
45 584 **histone modifications to gene expression levels in mouse embryonic stem cells.** *Nucleic Acids*
46 585 *Research* 2012, **40**(2):553-568.
47 586 14. Cheng C, Yan K-K, Yip KY, Rozowsky J, Alexander R, Shou C, Gerstein M, others: **A statistical**
48 587 **framework for modeling gene expression using chromatin features and application to**
49 588 **modENCODE datasets.** *Genome Biol* 2011, **12**(2):R15-R15.
50 589 15. Dong X, Greven MC, Kundaje A, Djebali S, Brown JB, Cheng C, Gingeras TR, Gerstein M, Guigó R,
51 590 Birney E *et al*: **Modeling gene expression using chromatin features in various cellular contexts.**
52 591 *Genome Biol* 2012, **13**(9):R53-R53.
53 592 16. Karličić R, Chung H-R, Lasserre J, Vlahoviček K, Vingron M: **Histone modification levels are predictive**
54 593 **for gene expression.** *Proceedings of the National Academy of Sciences* 2010, **107**(7):2926-2931.
55 594 17. Ma B, Wilker EH, Willis-Owen SAG, Byun H-M, Wong KCC, Motta V, Baccarelli AA, Schwartz J,
56 595 Cookson WOCM, Khabbaz K *et al*: **Predicting DNA methylation level across human tissues.** *Nucleic*
57
58
59
60
61
62
63
64
65

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

596 *acids research* 2014, **42**(6):3515-3528.

597 18. Zhang W, Spector T, Deloukas P, Bell J, Engelhardt B: **Predicting genome-wide DNA methylation**
598 **using methylation marks, genomic position, and DNA regulatory elements.** *Genome Biology* 2015,
599 **16**(1):14-14.

600 19. Das R, Dimitrova N, Xuan Z, Rollins RA, Haghghi F, Edwards JR, Ju J, Bestor TH, Zhang MQ:
601 **Computational prediction of methylation status in human genomic sequences.** *Proceedings of*
602 *the National Academy of Sciences* 2006, **103**(28):10713-10716.

603 20. Zheng H, Wu H, Li J, Jiang S-W: **CpGIMethPred: computational model for predicting methylation**
604 **status of CpG islands in human genome.** *BMC medical genomics* 2013, **6**(Suppl 1):S13-S13.

605 21. Neumann K, Klukas C, Friedel S, Rischbeck P, Chen D, Entzian A, Stein N, Graner A, Kilian B:
606 **Dissecting spatiotemporal biomass accumulation in barley under different water regimes using**
607 **high-throughput image analysis.** *Plant, cell & environment* 2015.

608 22. Breiman L: **Random forests.** *Machine learning* 2001, **45**(1):5-32.

609 23. Dietz H, Steinlein T: **Determination of plant species cover by means of image analysis.** *Journal of*
610 *Vegetation Science* 1996, **7**(1):131-136.

611 24. Leister D, Varotto C, Pesaresi P, Niwergall A, Salamini F: **Large-scale evaluation of plant growth in**
612 **Arabidopsis thaliana by non-invasive image analysis.** *Plant Physiology and Biochemistry* 1999,
613 **37**(9):671-678.

614 25. Paruelo JM, Lauenroth WK, Roset PA: **Estimating aboveground plant biomass using a photographic**
615 **technique.** *Journal of Range Management* 2000:190-193.

616 26. Walter A, Scharr H, Gilmer F, Zierer R, Nagel KA, Ernst M, Wiese A, Virnich O, Christ MM, Uhlig B *et*
617 *al*: **Dynamics of seedling growth acclimation towards altered light conditions can be quantified**
618 **via GROWSCREEN: a setup and procedure designed for rapid optical phenotyping of different**
619 **plant species.** *New Phytol* 2007, **174**(2):447-455.

620 27. Arvidsson S, Perez-Rodriguez P, Mueller-Roeber B: **A growth phenotyping pipeline for Arabidopsis**
621 **thaliana integrating image analysis and rosette area modeling for robust quantification of**
622 **genotype effects.** *New Phytol* 2011, **191**(3):895-907.

623 28. Hairmansis A, Berger B, Tester M, Roy SJ: **Image-based phenotyping for non-destructive screening**
624 **of different salinity tolerance traits in rice.** *Rice* 2014, **7**(1):16-16.

625 29. Neilson EH, Edwards AM, Blomstedt CK, Berger B, Møller BL, Gleadow RM: **Utilization of a high-**
626 **throughput shoot imaging system to examine the dynamic phenotypic responses of a C4 cereal**
627 **crop plant to nitrogen and water deficiency over time.** *Journal of experimental botany* 2015.

628 30. Catchpole WR, Wheeler CJ: **Estimating plant biomass: a review of techniques.** *Australian Journal*
629 *of Ecology* 1992, **17**(2):121-131.

630 31. Johnson JW: **A Heuristic Method for Estimating the Relative Weight of Predictor Variables in**
631 **Multiple Regression.** *Multivariate Behavioral Research* 2000, **35**(1):1-19.

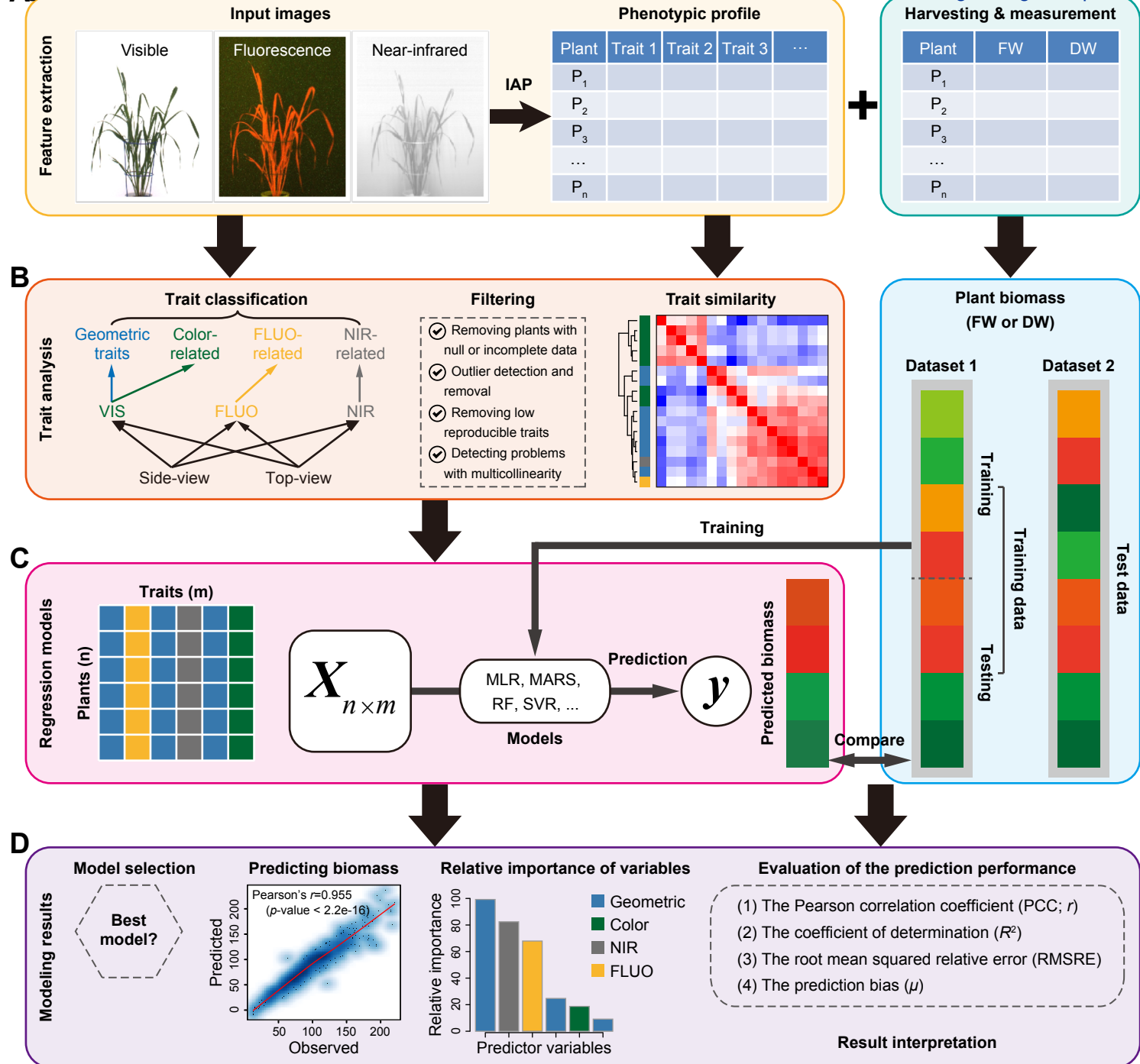
632 32. Chang C-C, Lin C-J: **LIBSVM: a library for support vector machines.** *ACM Transactions on Intelligent*
633 *Systems and Technology (TIST)* 2011, **2**(3):27.

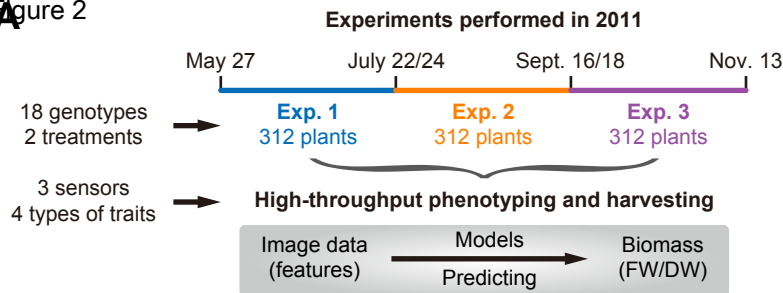
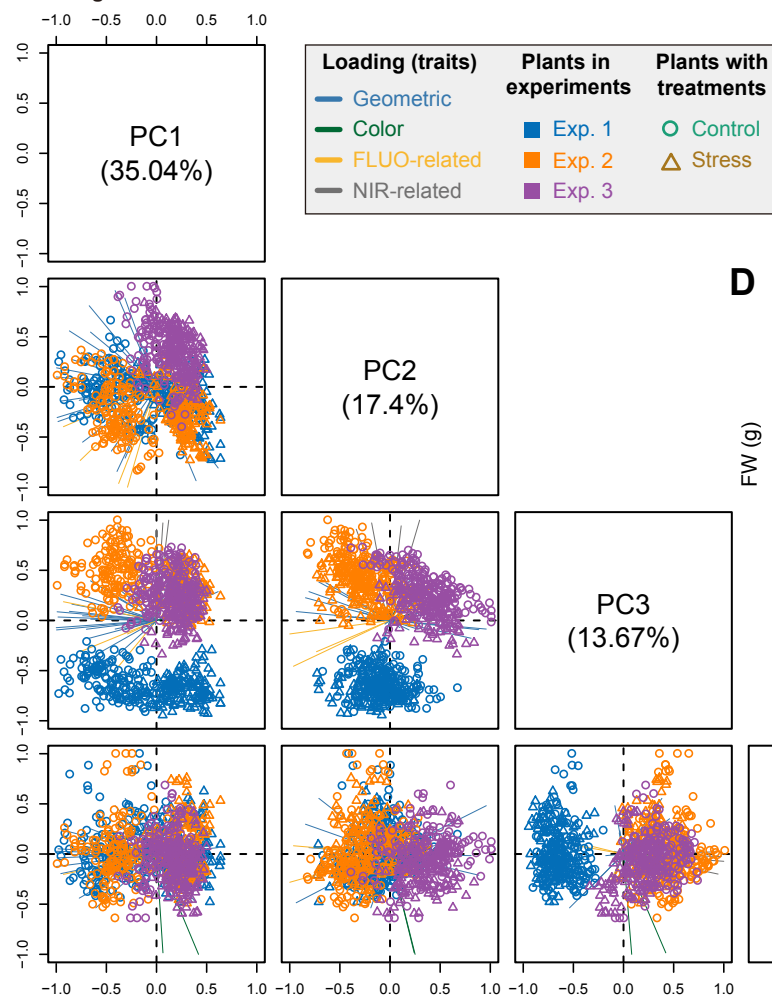
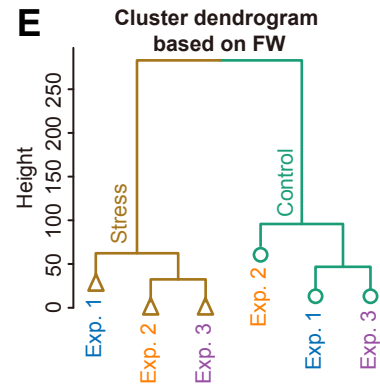
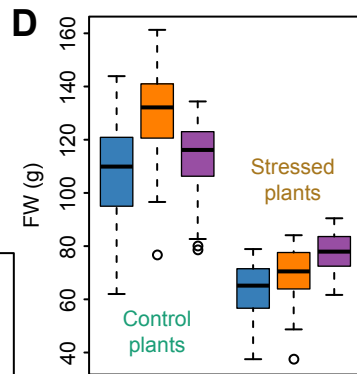
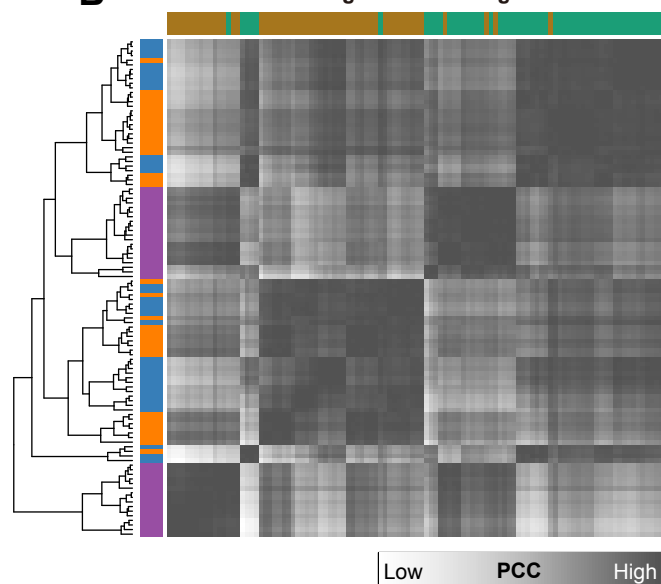
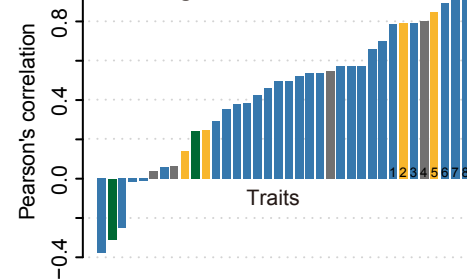
634 33. Loo LH, Wu LF, Altschuler SJ: **Image-based multivariate profiling of drug responses from single**
635 **cells.** *Nature methods* 2007, **4**(5):445-453.

636 34. Iyer-Pascuzzi AS, Symonova O, Mileyko Y, Hao Y, Belcher H, Harer J, Weitz JS, Benfey PN: **Imaging**
637 **and analysis platform for automatic phenotyping and trait ranking of plant root systems.** *Plant*
638 *physiology* 2010, **152**(3):1148-1157.

639 35. Arend D, Junker A, Scholz U, Schuler D, Wylie J, Lange M: **PGP repository: a plant phenomics and**
640 **genomics data publication infrastructure.** *Database : the journal of biological databases and*
641 *curation* 2016, **2016**.

642 36. Cwiek-Kupczynska H, Altmann T, Arend D, Arnaud E, Chen D, Cornut G, Fiorani F, Frohmberg W,
643 Junker A, Klukas C *et al*: **Measures for interoperability of phenotypic data: minimum information**
644 **requirements and formatting.** *Plant methods* 2016, **12**:44.



A Figure 2**C** PCA based on image-derived traits**B** Hierarchical clustering based on image-derived traits**F** Correlations between image-based traits and FW

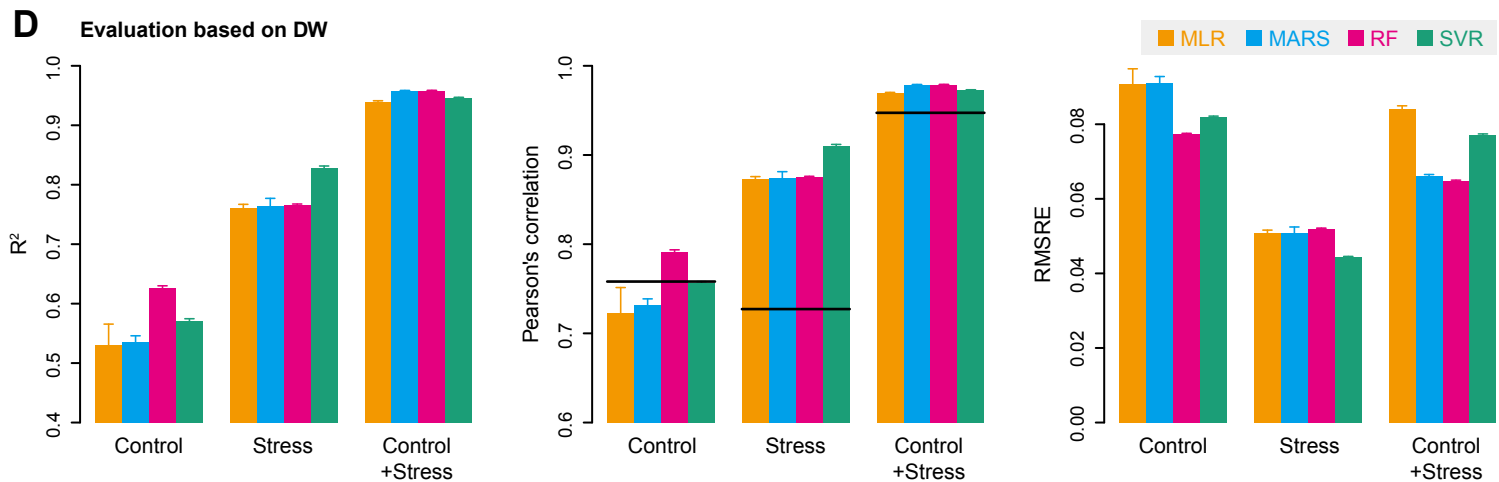
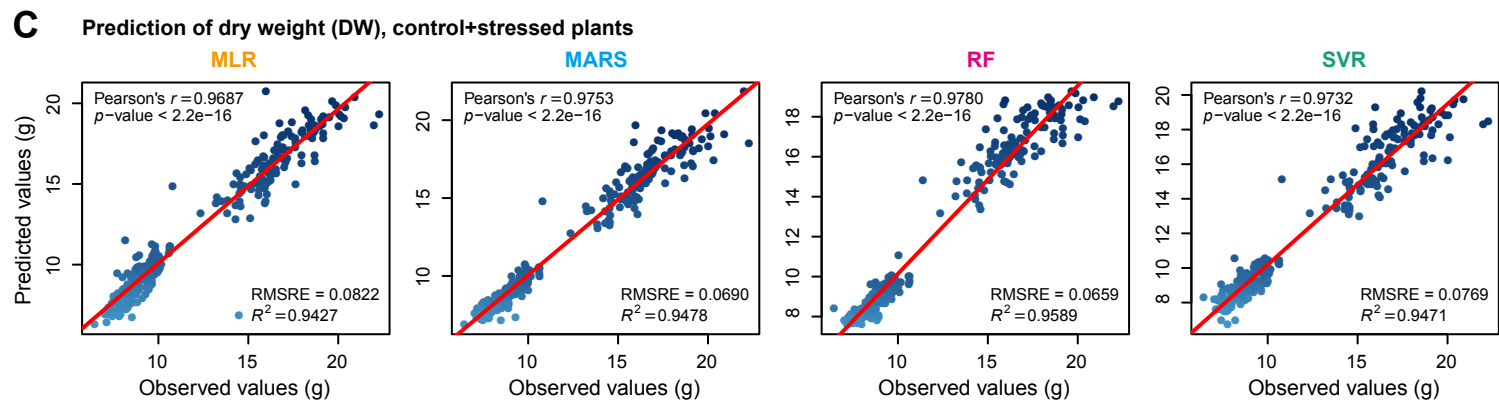
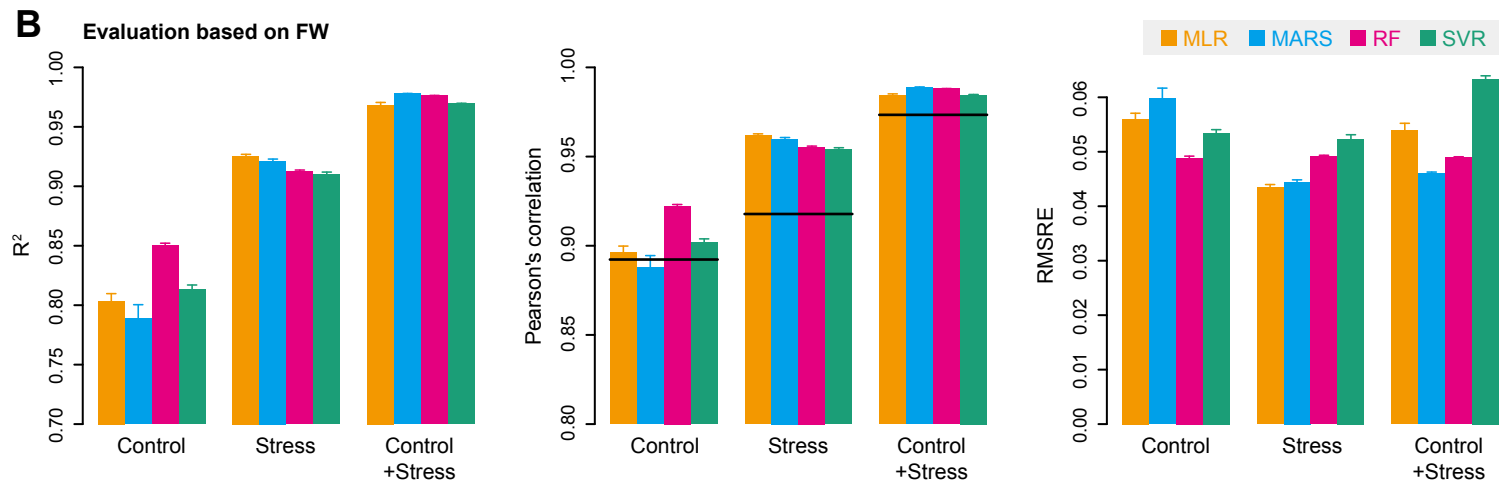
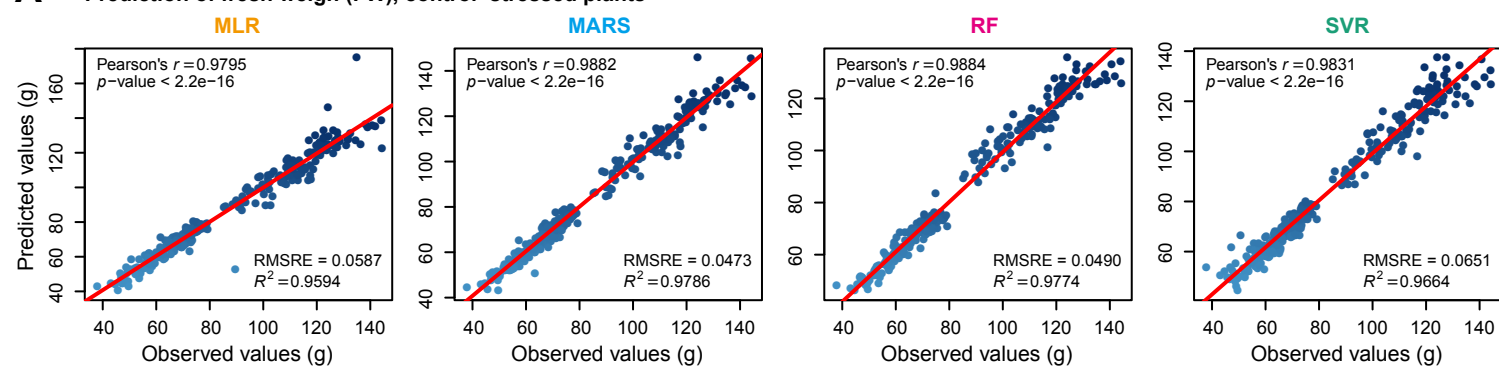
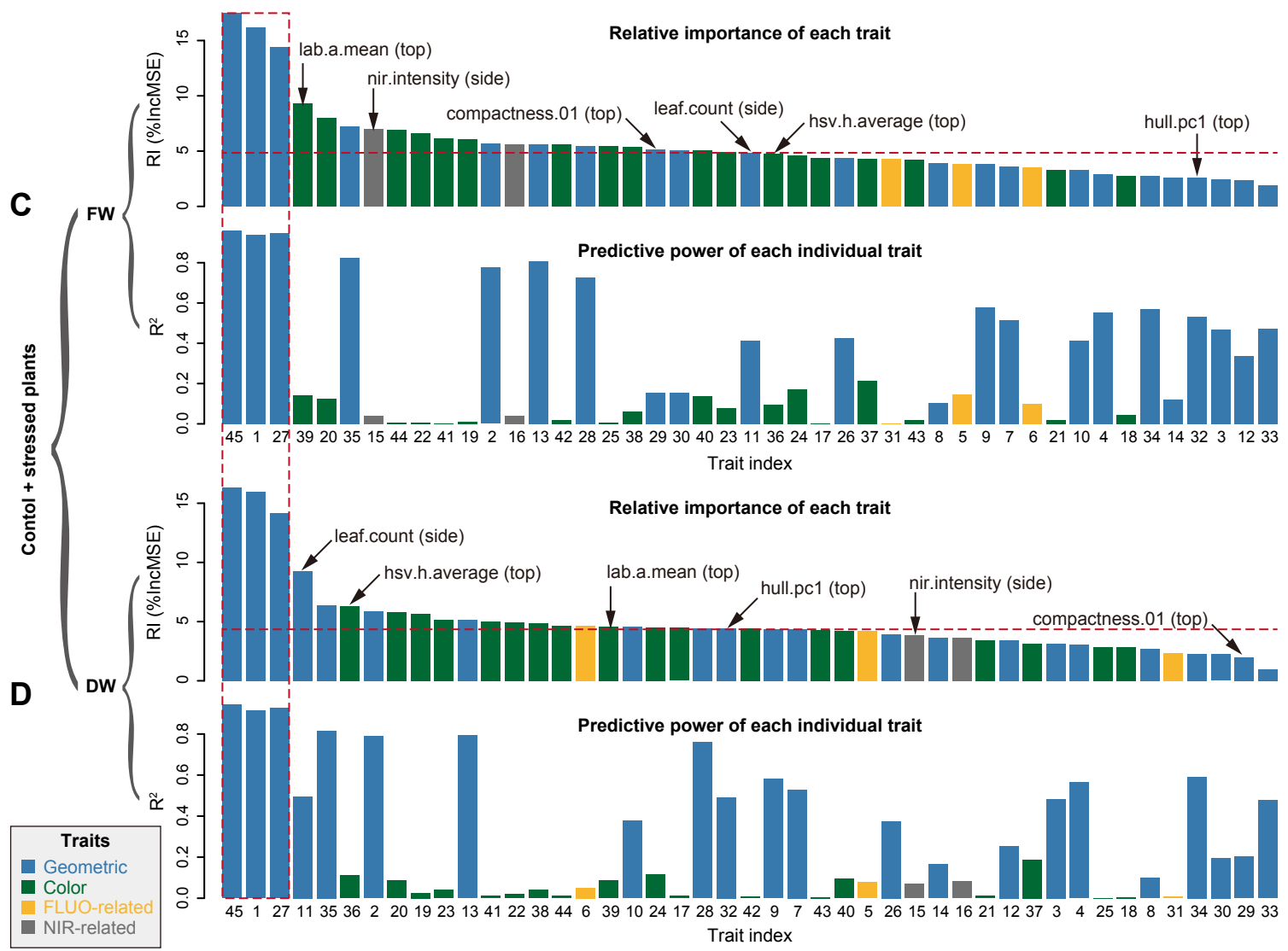
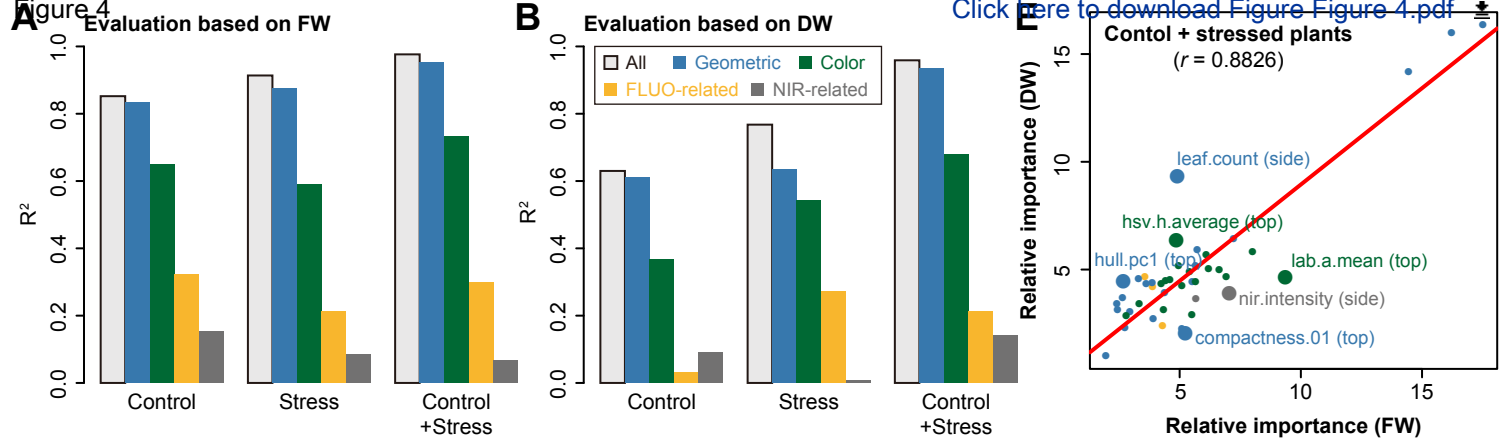
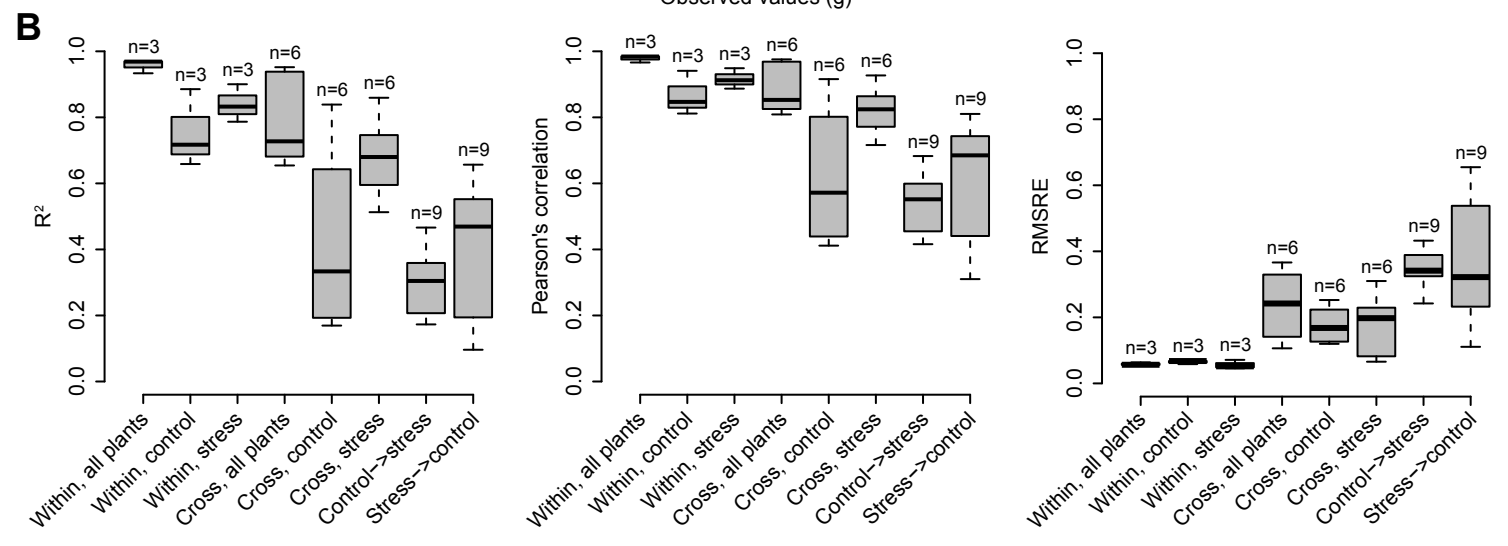
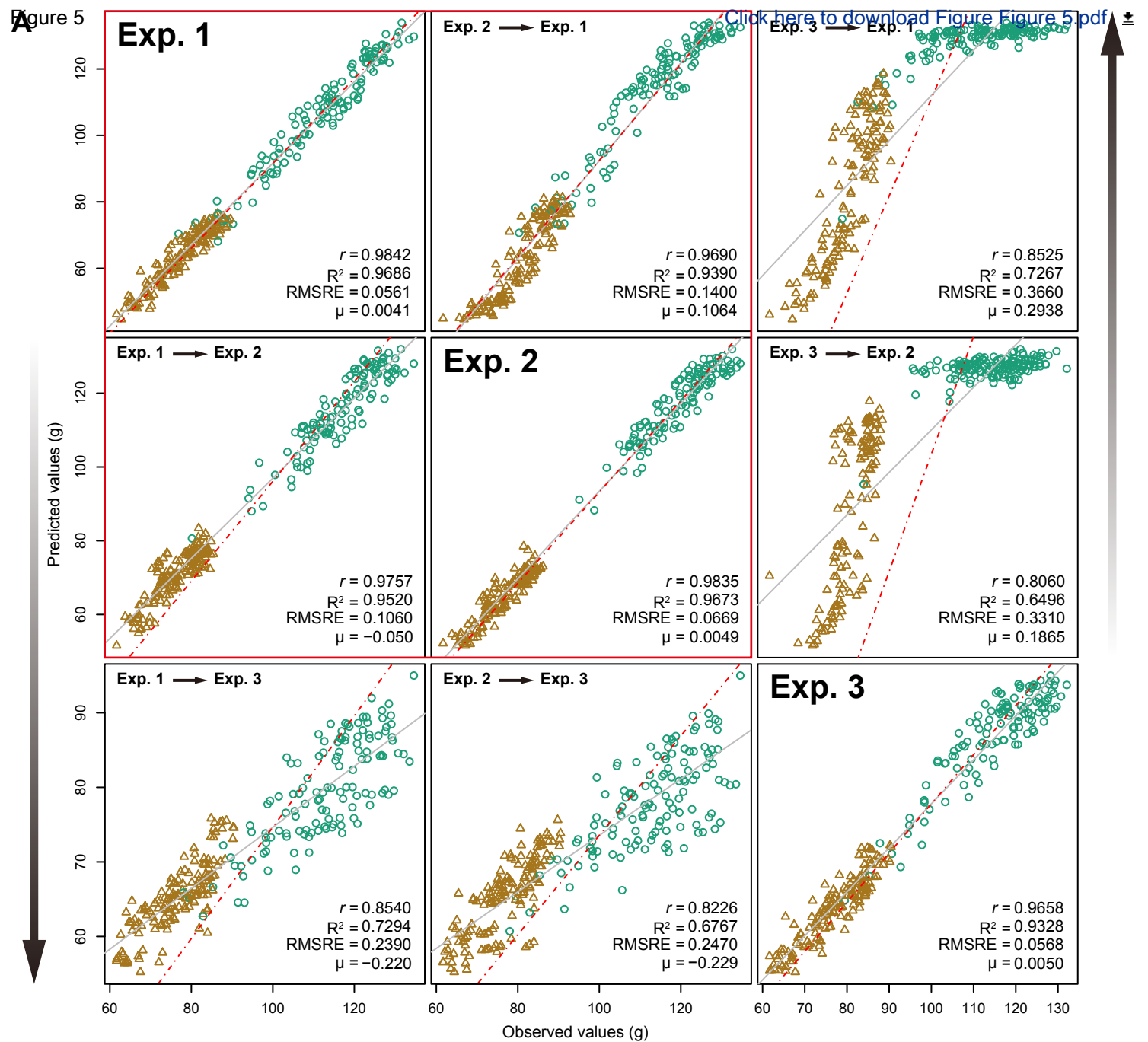
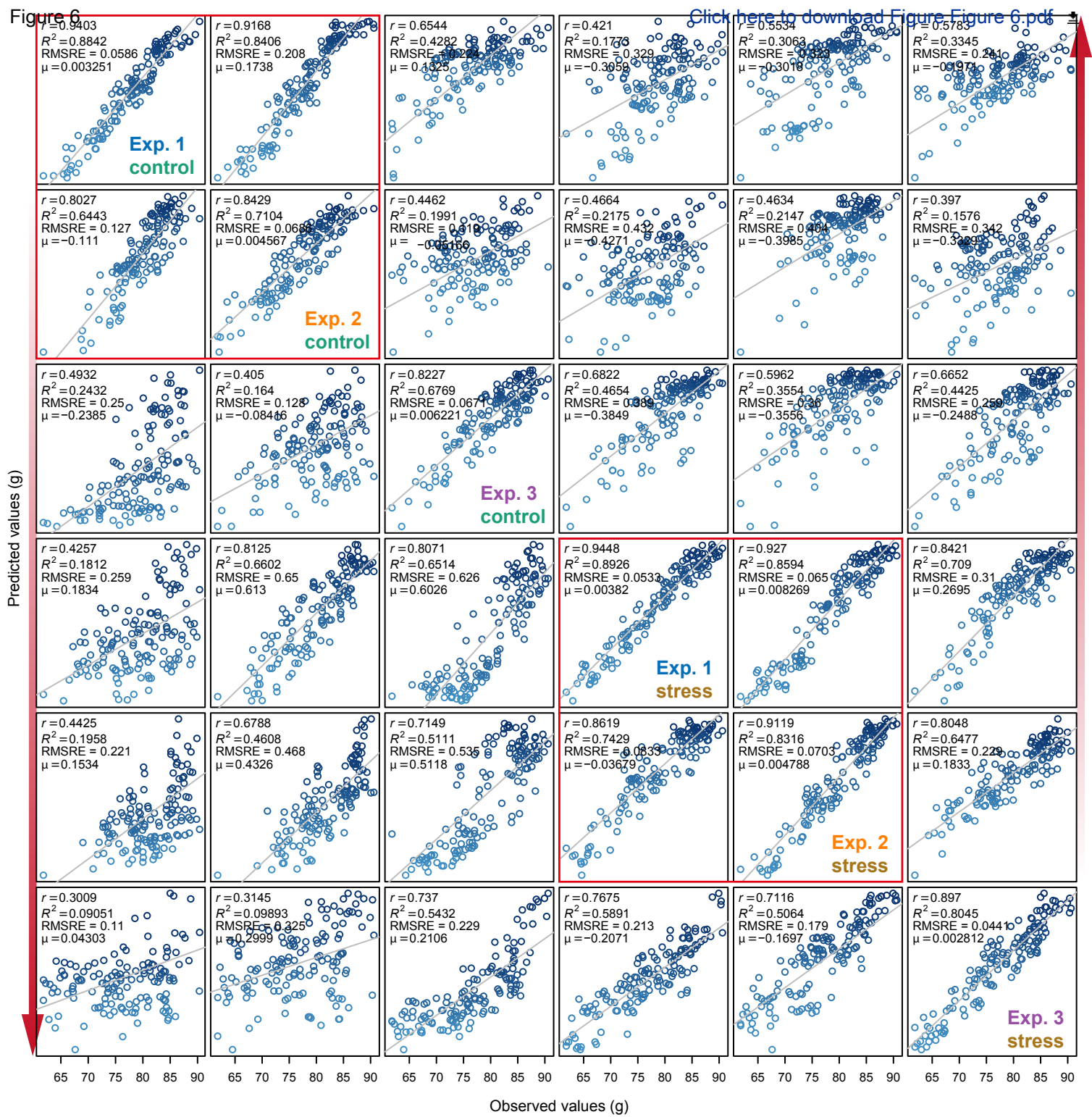


Figure 4











Click here to access/download
Supplementary Material
Supplemental Data S1.xlsx

