

Phenomics data processing: A plot-level model for repeated measurements to extract the timing of key stages and quantities at defined time points

Lukas Roth^{a,*}, María Xosé Rodríguez-Álvarez^{b,c}, Fred van Eeuwijk^d, Hans-Peter Piepho^e, Andreas Hund^a

^aETH Zurich, Institute of Agricultural Sciences, Universitätstrasse 2, 8092 Zurich, Switzerland

^bBCAM – Basque Center for Applied Mathematics, Alameda de Mazarredo, 14. E-48009 Bilbao, Basque Country, Spain

^cIKERBASQUE, Basque Foundation for Science, Bilbao, Spain

^dWageningen University and Research, Biometris, P.O. Box 16, 6700 AA Wageningen, The Netherlands

^eUniversity of Hohenheim, Institute for Crop Science, Biostatistics Unit, Fruwirthstrasse 23, 70593 Stuttgart, Germany

Abstract

Decision-making in breeding increasingly depends on the ability to capture and predict crop responses to changing environmental factors. Advances in crop modeling as well as high-throughput field phenotyping (HTFP) hold promise to provide such insights. Processing HTFP data is an interdisciplinary task that requires broad knowledge on experimental design, measurement techniques, feature extraction, dynamic trait modeling, and prediction of genotypic values using statistical models. To get an overview of sources of variations in HTFP, we develop a general plot-level model for repeated measurements. Based on this model, we propose a seamless stage-wise process that allows to carry on estimated means and variances from stage to stage and approximates the gold standard of a single-stage analysis. The process builds on the extraction of three intermediate trait categories; (1) timing of key stages, (2) quantities at defined time points or periods, and (3) dose-response curves. In a first stage, these intermediate traits are extracted from low-level traits' time series (e.g., canopy height) using P-splines and the quarter of maximum elongation rate method (QMER), as well as final height percentiles. In a second and third stage, extracted traits are further processed using a stage-wise linear mixed model analysis. Using a wheat canopy growth simulation to generate canopy height time series, we demonstrate the suitability of the stage-wise process for traits of the first two above-mentioned categories. Results indicate that, for the first stage, the P-spline/QMER method was more robust than the percentile method. In the subsequent two-stage linear mixed model processing, weighting the second and third stage with error variance estimates from the previous stages improved the root mean squared error. We conclude that processing phenomics data in stages represents a feasible approach if using appropriate weighting through all stages. P-splines in combination with the QMER method are suitable tools to extract timing of key stages and quantities at defined time points from HTFP data.

Keywords: High-throughput field phenotyping, dynamic modeling, stage-wise processing, canopy height

*Corresponding author: lukas.roth@usys.ethz.ch

30 Highlights

- 31 • General plot-level model for repeated high-throughput field phenotyping measurements
- 32 • Three main intermediate trait categories for dynamic modeling
- 33 • Seamless stage-wise process that allows to carry on estimated means and variances
- 34 • Phenomics data processing cheatsheet

35 1. Introduction

36 Advances in high-throughput field phenotyping (HTFP) allow capturing large data sets with high temporal
37 and spatial resolution (Rebetzke et al., 2019). Summarizing these spatio-temporal data in a meaningful way is
38 essential to support selection and decision-making in breeding. In HTFP the primary data often consists of images,
39 point measurements, orthophotos, or point clouds from which low-level traits (e.g., shoot counts, canopy cover,
40 canopy height, or senescence) are extracted. After feature extraction, these low-level traits may be tracked over
41 time in a subsequent temporal modelling step (van Eeuwijk et al., 2019; Moreira et al., 2020). If monitored
42 across the lifetime of a plant, low-level traits often follow some sort of monotonically increasing function (e.g.,
43 canopy height or senescence) or concave functions (e.g., number of growing shoots or canopy cover), which
44 allows estimating a dynamics' characteristics, referred to as intermediate traits.

45 Estimating such intermediate traits from spatio-temporal measurements implies *a priori* knowledge of growth
46 processes, best summarized in crop growth models. The performance of these crop growth models can only
47 advance if they become validated with field-based data (Ramirez-Villegas et al., 2015). Crop models have rapidly
48 gained in complexity over time, culminating in the description of plants by 3-D functional-structural models
49 (Vos et al., 2010). Indoor platforms have proven useful to characterize the dynamics of such models (Tardieu
50 et al., 2017), but discrepancies between field and indoor experiments raised doubts if results are always directly
51 transferable (Poorter et al., 2016). Field-based phenotyping may help to bridge this gap (Araus et al., 2018).

52 While under controlled conditions environmental factors may be adequately controlled, the lack of control
53 over meteorological conditions poses a major challenge for field phenotyping. Several additional types of errors
54 need to be considered, which can be classified into those directly affecting the sensor reading, and those affecting
55 the plant development.

56 In HTFP there are attempts to quantify genotype-specific timing of phenology stages (Hurtado et al., 2012)
57 and response patterns to distinct environmental variables like temperature (Grieder et al., 2015; Kronenberg
58 et al., 2020a). A comparable approach in genomics uses functional mapping of quantitative trait loci (QTLs),
59 e.g., based on logistic growth curves (Ma et al., 2002; Malosetti et al., 2006). Ma et al. proposed to distinguish
60 three biological processes in such models: allometric laws, growth models, and reaction norms. Characterizing
61 crop model dynamics using field data becomes increasingly difficult as models become more complex. A solution

62 is to model the dynamic process of growth based on traits or scores that lack a clear physiological interpretation.
63 In phenomics, this was demonstrated using serial measurements as predictors for statistical learning (e.g. [Ubbens
64 et al., 2020](#); [Maimaitijiang et al., 2020](#); [Herrero-Huerta et al., 2020](#)). In genomics, comparable approaches are
65 based on functional principal component analysis, where curves are specified as linear combinations of basis
66 functions, and the corresponding scores then used as intermediate traits ([Kwak et al., 2016](#); [Moreira et al.,
67 2020](#)).

68 From a plant physiology point of view, such approaches represent a 'black box': Drawing conclusions on the
69 biological importance of the underlying traits is rather difficult. Therefore, we believe that a classical approach to
70 extract traits related to distinct crop ideotypes based on *a priori* knowledge is more suitable (see also [van Eeuwijk
71 et al., 2019](#); [Bustos-Korts et al., 2019](#)). This approach may then represent a standard to compare modern learning
72 approaches with.

73 Based on HTFP literature and the biological processes described in [Ma et al. \(2002\)](#), we identified three main
74 intermediate trait categories which can be related to ideotype concepts:

- 75 1. **Timing of key stages:** Turning points in the dynamics of numeric measurements which may be related to
76 phenology; e.g., beginning of stem elongation ([Kronenberg et al., 2017](#)), time of canopy closure ([Soltani
77 and Galeshi, 2002](#)), time of maximum canopy growth rate ([Borra-Serrano et al., 2020](#)), heading and flow-
78 ering ([Sadeghi-Tehran et al., 2017](#)), or onset and end of senescence ([Anderegg et al., 2020](#); [Aasen et al.,
79 2020](#)). Genotype-specific responses to environmental covariates and/or indices may help to predict key
80 stages; e.g., flowering time ([Millet et al., 2019](#)).
- 81 2. **Quantities at defined time points or periods:** Traits based on numeric measurements; either at a steady
82 state; e.g., canopy temperature between flowering and beginning of senescence ([Perich et al., 2020](#)), or
83 at well-defined time points; e.g., number of tillers at beginning of stem elongation ([Roth et al., 2020](#)) and
84 at harvest ([Jin et al., 2019](#)), number of ears at harvest ([Fernandez-Gallego et al., 2018](#)), or canopy cover
85 at maximum ([Borra-Serrano et al., 2020](#)). Area-under-the-curve traits may represent a special case of this
86 category where one summarizes quantities over a defined range of time points ([Blancon et al., 2019](#)).
- 87 3. **Dose-response curves:** Traits that describe developmental responses in dependence on covariates between
88 clearly defined boundary key stages. Dose-response experiments are classically conducted under controlled
89 conditions, e.g., by examining the response of leaves to temperature and water deficit ([Reymond et al.,
90 2003](#)) and to soil water deficiency and evaporative demand ([Welcker et al., 2011](#)) during their linear
91 growth phase. More recently, such experiments were conducted in the field; e.g., in the early, exponential
92 development phase of canopy cover between emergence and tillering ([Grieder et al., 2015](#)) or at the linear
93 development phase of canopy height between start and end of stem elongation ([Kronenberg et al., 2020a](#)).

94 Despite the differences in subsequent processing, the extraction of each of the three different trait categories
95 is a highly repetitive task which requires analysis routines with sufficient robustness and generality. While timing

96 of key stages and quantities belong to growth model processes, dose-response curves relate to reaction norm
97 processes (Via et al., 1995). Arguably, dose-response curves represent the most challenging modelling aspect in
98 field phenotyping, as they require quantifying growth and relate it to environmental covariates. We will cover this
99 aspect in a follow-up paper. However, a robust evaluation of such dose-response curves requires to determine the
100 boundaries between which a steady development takes place. Here, we aimed to develop a method to extract
101 such timing of key stages and quantity traits.

102 We start by developing a plot-level model for repeated measurements, with a focus on the outdoor field
103 phenotyping platform FIP (Kirchgessner et al., 2017). The FIP allows to densely monitor a large set of replicated
104 genotypes ($\geq 2 \times 300$) over a whole growing season with genotypes being the only treatment. The aim of such
105 experiments is to i) allow developing new traits and phenotyping methodologies; ii) characterize a specific target
106 environment including the targeted ideotypes; and iii) to serve as part of a multi-environment experiment that
107 covers a mega-environment. We propose a possible solution to analyze such experiments based on existing
108 statistical tools such as P-spline fitting and stage-wise linear mixed model analysis. We further evaluate and
109 demonstrate the suitability of the approach to extract the timing of key stages and quantities at defined time
110 points from low-level traits using simulated wheat canopy height data.

111 2. Materials and Methods

112 2.1. A plot-level model for repeated measurements

113 A planned experiment generally includes an experimental design (Figure 1b, green boxes) in which the treat-
114 ment factors to be tested are randomly assigned to experimental units (usually plots). For the FIP, the only
115 treatment factor are genotypes. The design comprises only one site but multiple years. The data for each year
116 holds a subset of treatment levels (genotypes) together with checks and design factors (blocks) to allow correcting
117 for spatial variability at the site. In the specific case of the FIP, a panel of on average 345 genotypes is replicated
118 twice per year and each replication is planted on a different lot in the FIP area. Each replication is augmented
119 with spatial checks in a 3×3 block arrangement.

120 Performing measurements includes the application of a sensing device collecting measurements from a plot
121 (Figure 2a-c). This process results in data which either directly represent a trait value (e.g., a point measurement
122 of temperature) or can be translated to one or several low-level traits by means of feature-extraction (Figure 1b,
123 blue boxes). Measurements of the same type on the same set of plots, but at subsequent time points, can be
124 summarized as campaigns (Figure 1b, red boxes). The sensing concept and the degree of automation determine
125 the time intervals at which measurements can be taken.

126 Nowadays, the measuring intervals of campaigns often last days rather than minutes or hours. Here, we
127 define a campaign time point as the time at which the whole set of plots (usually belonging to an experimental
128 design or year-site) is completely measured. Such a measurement may take seconds to hours depending on the

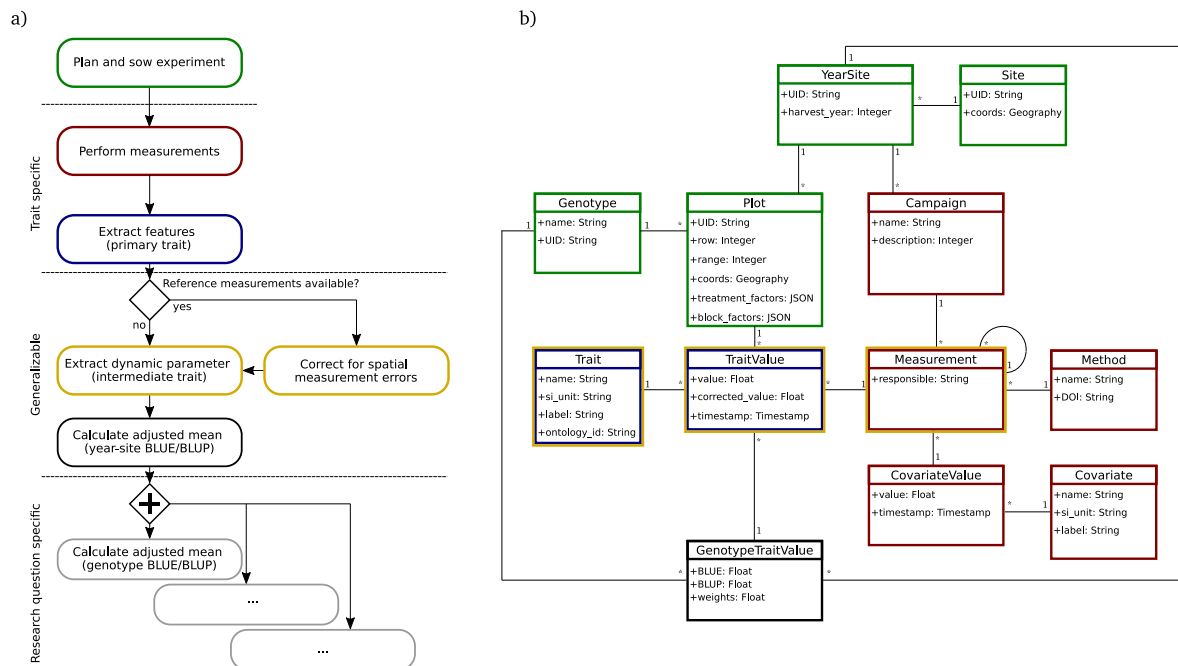


Figure 1: Minimal process-driven model for the FIP: a) Process model, b) Data model.

129 phenotyping methodology and size of the design. By contrast, a measuring time point (or timestamp) is the exact
 130 timing at which an experimental unit (usually a plot) is measured.

131 The same approach of campaigns and measurements also applies to covariate measurements. Covariates are
 132 usually measured at very short time intervals of minutes to hours. In contrast to traits, environmental covariates
 133 can be measured at various levels (e.g., year-site, plot, plant, or plant organ). The measurement level determines
 134 what is regarded as phenotypic heterogeneity caused by covariate variation. The FIP site comprises a site-specific
 135 local weather station, which corresponds to covariate measurements at year-site level (Figure 3b). Therefore,
 136 one must consider heterogeneity caused by covariate variations at plot, plant and organ levels and their effects on
 137 plant growth (Figure 3a)—namely variations of the timing of key stages (Figure 3c1) resulting in growth period
 138 condition variations (Figure 3c2) and consequently variations of quantities at defined time points (Figure 3c3).

139 In a phenotyping experiment one has to distinguish between nuisance factors affecting the growth and devel-
 140 opment of the plant (Figure 3a1–3), and measurement errors affecting the precision at which a certain phenotype
 141 is measured at a given time (Figure 3a4). The latter factors may affect whole campaign time points (i.e., at the
 142 day-to-day level, Figure 3c4, red one-sided arrow) but also individual measuring time points within a day (Figure
 143 3c4, red two-sided arrow). Nuisance factors affecting growth and development are, e.g., soil fertility inhom-
 144 geneities, spatial temperature gradients, mices, herbivore damages, and other abiotic and biotic factors varying
 145 spatially and temporally in the field. Such factors will in the best case lead to spatio-temporally correlated phe-
 146 notypic observations.

Phenomics data processing cheatsheet

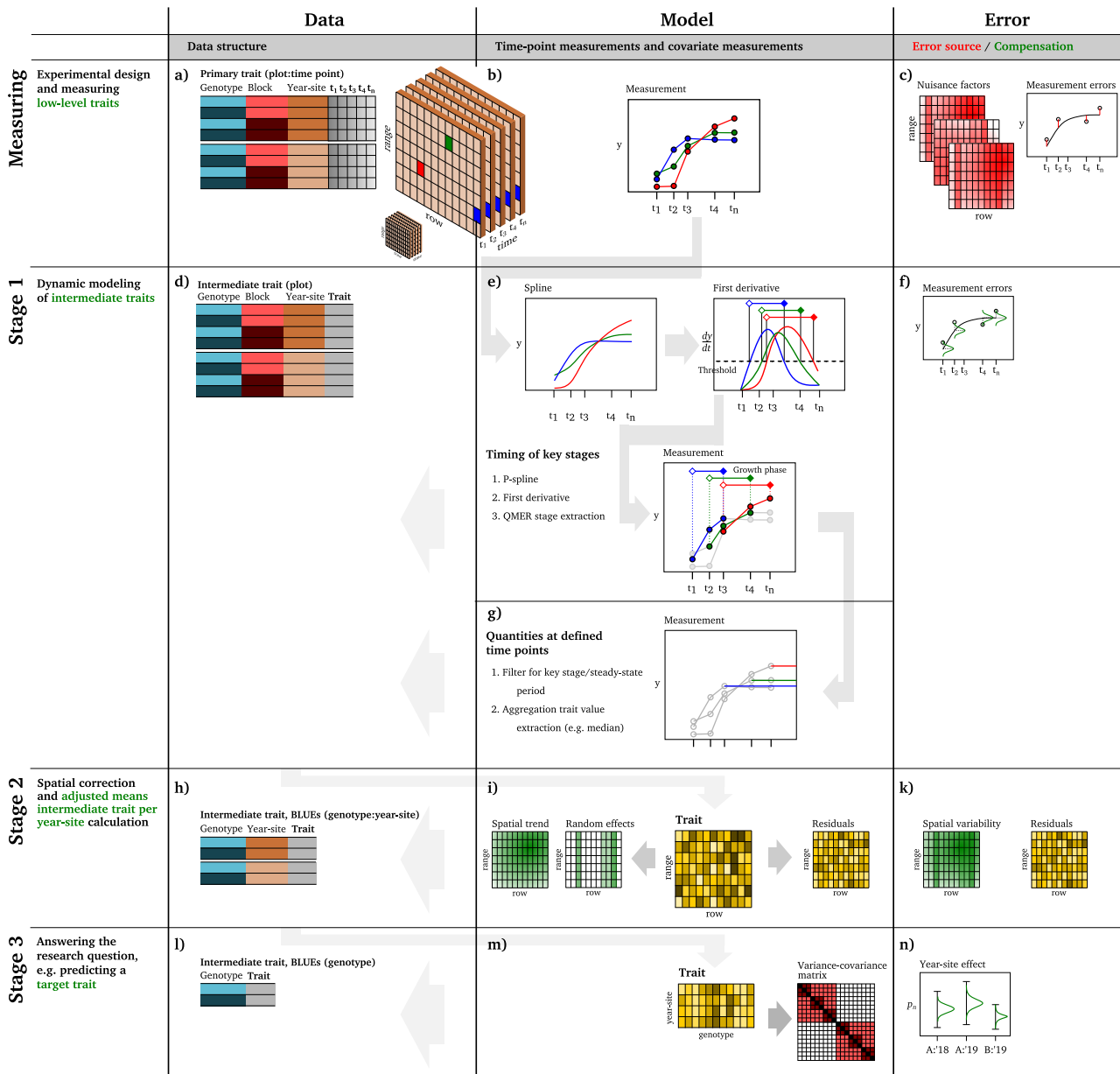


Figure 2: Phenomics data processing cheatsheet: Extraction of timing of key stages and quantities at defined time points from high-throughput field phenotyping data.

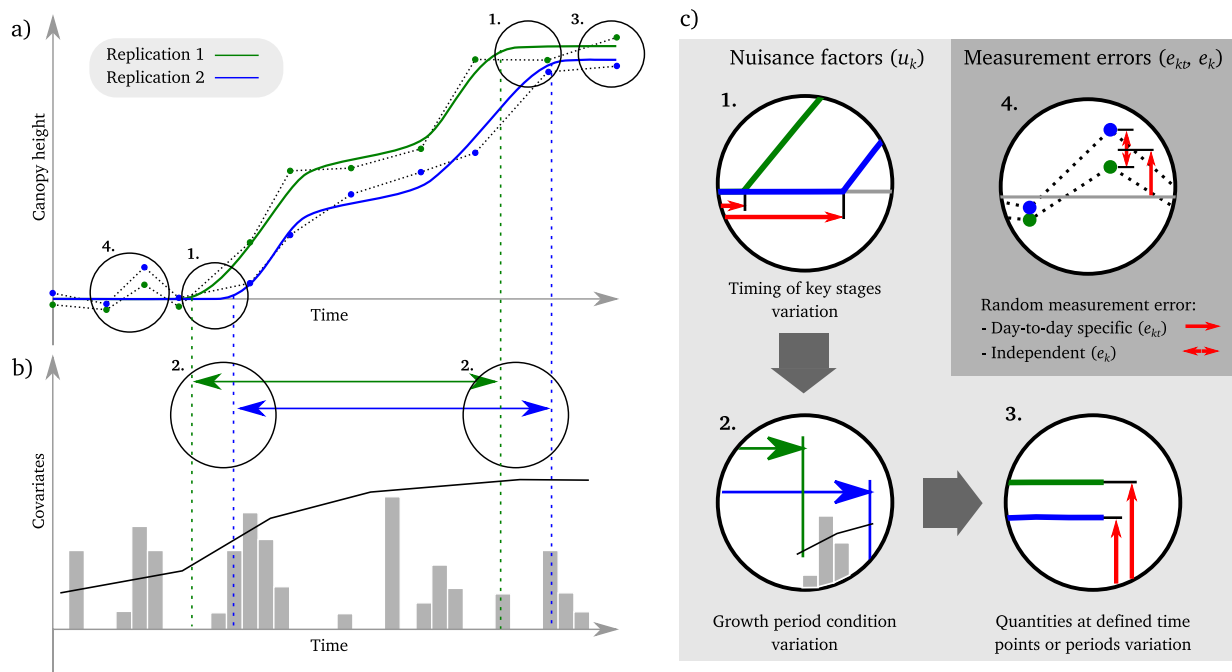


Figure 3: Sources of variation in HTPF on the example of canopy height measurements. (a) Canopy height development of two replications of the same genotype (green and blue lines) and realized measurement time points (green and blue points). (b) Covariate measurements during the growth phase of canopies (e.g., temperature and precipitation). (c) Sources of variation: (1) spatial and crop-husbandry effect leading to different timings of key stages, e.g., start and end of stem elongation; (2) timing of key stage variations leading to variations in the different gradients of environmental covariates, e.g., temperature gradients in the stem elongation phase; (3) spatial and crop-husbandry effects leading to quantitative variations in trait values; e.g., final height at the end of the stem elongation phase; (4) Day-to-day random measuring errors, e.g., related to differing conditions between measurement days; and independent random measuring errors, e.g., related to the measurement precision of the device.

147 Sources for measurement errors are, e.g., factors differing between campaign time points. These factors may
148 lead to day-specific under- or overestimation of measurements, e.g., due to positioning shifts of the sensor head,
149 re-adjustment of sensor settings between measurement campaigns, changes in canopy characteristics after rain
150 or during hot days, and differing illumination conditions (Figure 3c4, red one-sided arrow). Taking reference
151 measurements (e.g., by the use of calibration targets) allows correcting for some of these errors, but such mea-
152 sures may not always be feasible in crop phenotyping experiments. Apart from the effects related to the whole
153 campaign time point, changing conditions during the measuring sequence may lead to additional, temporally cor-
154 related measurement errors among measuring time points. Sources for such errors are, e.g., changing weather
155 conditions during a measurement that takes a considerable amount of time. Such temporal effects may translate
156 into apparent spatial effects within a campaign time point and therefore be confounded with nuisance factors. Fi-
157 nally, random measurement device errors (Figure 3c4, red two-sided arrow) represent another source of variation
158 in HTFP. These errors are usually assumed to be identically and independently normally distributed.

159 Consequently, we define a HTFP observation y_{kt} for the t -th time point on the k -th plot ($k = 1, \dots, K$) as the
160 result of a dynamic model g that is a function of time t and of a vector ($\vec{\cdot}$) of plot-specific crop growth parameters
161 $\vec{\theta}_{k(i)}$ associated with genotype i modulated ($;$) by a vector of time-varying covariates \vec{x}_t , and of a plot residual e_{kt}
162 that is i.i.d. ($\sim \mathcal{N}(0, \sigma_k^2)$),

$$y_{kt} = g(t, \vec{\theta}_{k(i)}; \vec{x}_t) + e_{kt}. \quad (1)$$

163 While e_{kt} will account for random measurement device errors, we assume here that g will absorb any spatio-
164 temporal correlation among measurements. Dynamic modeling (Equation 1) is done separately for each individ-
165 ual plot-based time series (Stage 1), i.e., (y_{k1}, \dots, y_{kT}) (Figure 2d-g).

166 Stage 1 therefore estimates plot-specific crop growth model parameters $\vec{\theta}_{k(i)}$. Those crop growth model pa-
167 rameter will become a phenotypic trait when measured / estimated at a set of genotypes. Correcting for spatial
168 correlations is done in a subsequent stage (Stage 2) of a stage-wise approach to get estimates of genotype specific
169 crop growth model parameters $\vec{\theta}_i$ (Figure 2h-k). This estimation step is done separately for each crop growth
170 model parameter in $\vec{\theta}_i$ based on fitting the linear model

$$\hat{\theta}_{k(i)} = \theta_i + u_k + e_k, \quad (2)$$

171 where $\hat{\theta}_{k(i)}$ is the crop growth parameter estimate from Stage 1, u_k a spatially correlated random component, and
172 e_k are plot residuals assumed to be normally distributed with zero mean and $\text{var}(e_k) = \sigma^2 w_k^{-1}$, where $w_k = (\text{s.e.})_k^{-2}$
173 are weights based on the standard error estimates (s.e.) from Stage 1. For a stage-wise approach with weights
174 based on variance estimations, one usually fixes σ^2 to unity. Nevertheless, if expecting proportionality of $\text{var}(e_k)$
175 to w_k^{-1} only—e.g., when the s.e.'s are derived from a correlated trait—it is required to estimate σ^2 . The spatially
176 correlated error term u_k will absorb any spatial correlation caused by random measurement errors and by physical
177 phenotypic differences, and e_k any plot-specific residual.

178 This approach is not limited to parametric or dimensionality reduction techniques but allows including ar-
179 bitrary dynamic models g with high complexity based on biologically meaningful traits. Nevertheless, it also
180 obviates modeling a spatio-temporally correlated residual term in its full extent by assuming that all serial cor-
181 relation is accounted for by the time-dependence of g . In the following, we hypothesize and exemplify with a
182 simulation that our approximation of the spatio-temporal correlation structure is well suited to extract interme-
183 diate traits with adequate precision from HTFP data.

184 2.2. Dynamic modeling of three trait categories

185 In dynamic modeling, one has to specify a method, based on g of Equation 1, to estimate a vector of meaning-
186 ful plot-level traits $\vec{\theta}_{k(i)}$ (for brevity we henceforth drop the index i for genotypes, referring to $\vec{\theta}_k$, it being under-
187 stood that a plot-level parameter is always genotype-specific) based on measured phenotypes y_{kt} and measured
188 covariates x at (potentially differing) time points t . In the following, we will provide theoretical considerations
189 and specific examples for each of the three trait categories defined in the introduction, (1) timing of key stages'
190 traits, (2) quantities at defined time points or periods, and (3) dose-response curve traits.

191 The first intermediate trait category—timing of key stages—describes growth as a sequence of key stages.
192 Such phenology traits are well-known in agronomy, e.g., the timing of jointing, heading, and flowering in wheat.

193 The second intermediate trait category—quantities at defined time points or periods' traits—describes pheno-
194 typic characteristics at key stages or steady state phases. Hence, such traits include a time point definition, e.g.,
195 with traits of the first category. The number of tillers per plant at jointing, the number of ears per square unit at
196 harvest, or the average canopy cover between tillering and jointing are examples of such traits for wheat.

197 The third intermediate trait category—dose-response curves—describes phenotypes as the result of a dose-
198 response model dependent on a covariate course between key stages. Hence, also these traits require time point
199 definitions, e.g., with traits of the first category. The response of the stem elongation to temperature is an example
200 of such a trait for wheat.

201 To obtain traits of the first two categories, we favor semi-parametric approaches (e.g., spline fitting) over
202 parametric approaches (e.g., logistic regression) for the dynamic modeling based on the following considerations:
203 Taking the example of early canopy development of winter wheat, where one wants to extract a timing (1) or
204 quantity (2) trait at a specific stage, growth may fluctuate strongly due to the environment, leading to a “stepped”
205 growth curve (Figure 3a). While non-parametric approaches are able to follow such growth curves, parametric
206 approaches would require to modify the timescale to, e.g., thermal time. Despite the fact that thermal time is a
207 widely accepted concept in agriculture (Parent et al., 2019), it is nevertheless based on model assumptions such
208 as the existence of a base temperature, and the linearity of the response. Using such a scale is therefore at odds
209 with the research aim to identify the model behind timing and quantity traits.

210 When using a semi-parametric approach (e.g., P-splines), one approximates g with a plot-specific model, i.e.
211 a smooth function of time $s(t)$. To extract traits of the first category—timing of key stages—from such a smooth

212 function, a set of methods q_n ($n = 1, \dots, N$) to estimate timing traits $\theta^{T(n)}$ (e.g., to approximate the end of the
213 stem elongation phase) from s has to be defined,

$$g(t, \vec{\theta}_k; x_t) \hat{=} s_k(t), \quad (3)$$

$$\theta_k^{T(n)} = q_n(s_k), \quad (4)$$

214 where $\hat{=}$ indicates that s_k estimates g for the k -th plot.

215 Extracting traits of the second category—quantities at defined time points or periods—builds on the spline
216 function s (Equation 3) and extracted timing of key stages (Equation 4) but inverts the approach of extracting
217 key stages: If $\theta^{T(n)}$ represent timing of key stages (e.g., the end of stem elongation), then quantities at defined
218 time points $\theta^{Q(n)}$ (e.g., canopy cover at the approximated end of stem elongation) may be extracted from the
219 spline s as

$$\theta_k^{Q(n)} = s_k(\theta_k^{T(n)}). \quad (5)$$

220 It is important to note that the underlying low-level traits for the timing trait $\theta^{T(n)}$ and the spline s in Equation 5
221 may differ, giving rise to a vast amount of possible trait combinations, e.g., when combining canopy height
222 timing traits with canopy cover quantity traits. While Equation 5 extracts quantities at points in time, extracting
223 aggregated quantities (e.g., normalized area-under-the-curve traits) for a period of time may be of interest as
224 well. If $\theta^{T(a)}$ and $\theta^{T(b)}$ represent two cautiously chosen timings of key stages' traits where $\theta^{T(a)} < \theta^{T(b)}$ (e.g.,
225 approximated start and end of flowering), then a quantity at defined time period trait $\theta^{Q(a...b)}$ (e.g., average
226 temperature at approximated flowering) may be extracted from s as

$$\theta_k^{Q(a...b)} = \frac{1}{\theta_k^{T(b)} - \theta_k^{T(a)}} \int_{\theta_k^{T(a)}}^{\theta_k^{T(b)}} s_k(t) dt. \quad (6)$$

227 If either $\theta_k^{T(a)}$ or $\theta_k^{T(b)}$ corresponds to a time series boundary (e.g., end of stem elongation to end of time series),
228 the trait may represent an initial or final trait value (e.g., final height).

229 For the third trait category—dose-response curves—one describes a phenotype as the result of a dose-response
230 model \dot{g} that relates growth rates to a covariate course x_t and a corresponding set of crop growth model param-
231 eters $\theta^C = (\theta^{C(1)}, \theta^{C(2)}, \dots, \theta^{C(L)})$ where L is the total number of parameters of the dose-response curve,

$$g(t, \vec{\theta}_k; x_t) = \int_{\theta_k^{T(a)}}^{\theta_k^{T(b)}} \dot{g}(\theta_k^C, x_t) dt. \quad (7)$$

232 Similar to quantities at defined time periods' traits (Equation 6), dose-response curve traits require the definition
233 of a corresponding growth phase, characterized by a start ($\theta^{T(a)}$) and a stop ($\theta^{T(b)}$). Therefore, a preliminary
234 extraction of traits of the category one (Equation 4) is required. Subsequently, θ^C may be estimated.

235 The striking similarity of Equation 6 and 7 is no coincidence: The area-under-the-curve of a defined growth
236 period can be seen as a direct cause of a response to covariates in this growth phase. The two approaches differ

237 in how they include covariates: While dose-response curves model an explicit dependency to covariates, an area-
238 under-the-curve quantifies implicitly the result of such a dependency.

239 An example for a dose-response curve \dot{g} at a defined growth phase is the stem elongation rate of wheat in
240 relation to temperature. Extracting such a dose-response curve implies that one is interested in fitting a specific
241 non-linear function.

242 2.3. Combining multi-year measurements

243 HTFP platforms such as the FIP are usually run on a continuous basis, thus increasing the number of year
244 measurements with each year of operation since inauguration. Experimental designs and genotype sets may
245 change to some extent along the years. The question is how to combine such multi-year measurements in a way
246 that one can process years in stages, which is of high benefit for both documentation purpose and processing
247 requirements.

248 The problem of stage-wise analysis we are addressing here has a long history (Cochran, 1954) and is well
249 known in plant breeding (Smith et al., 2005; Piepho et al., 2012) and also in other contexts, most notably in
250 meta-analysis (Whitehead, 2002; Borenstein et al., 2009). Most commonly, the problem arises in settings where
251 information needs to be combined across several experiments, whereas in the present work we consider the case
252 where different pieces of information need to be combined across units in a single experiment. Despite these
253 differences in scale, the statistical challenges are the same. To illustrate, consider a simple setting in which a set
254 of replicated genotypes is tested for yield at a number of years in a platform. The response of the i -th genotype
255 on the k -th plot at the j -th year can be written as

$$y_{ijk} = \mu + g_i + v_j + (gv)_{ij} + e_{ijk}, \quad (8)$$

256 where μ is an intercept, g_i is the main effect of the i -th genotype, v_j the main effect of the j -th year, assumed to
257 be normal with zero mean and variance σ_v^2 , $(gv)_{ij}$ is the interaction of the i -th genotype and j -th year assumed to
258 be normal with zero mean and variance σ_{gv}^2 , and e_{ijk} a residual error assumed to be normal with mean zero and
259 year-specific variance $\sigma_{e(j)}^2$. An objective among others in field phenotyping platforms is to estimate genotype
260 means across years, $\eta_i = \mu + g_i$ and their differences.

261 This can be done in a single stage by fitting the model (Equation 8) directly to plot data y_{ijk} . Alternatively,
262 we may proceed in two stages and first estimate genotype means per year using sample means $\bar{y}_{ij\cdot}$. These means
263 have variance $\text{var}(\bar{y}_{ij\cdot}) = r_{ij}^{-1} \sigma_{e(j)}^2$, where r_{ij} is the number of replications of the i -th genotype in the j -th location.
264 In the second stage, we can fit the model

$$\bar{y}_{ij\cdot} = \mu + g_i + v_j + (gv)_{ij} + \bar{e}_{ij\cdot}, \quad (9)$$

265 where $\text{var}(\bar{e}_{ij\cdot}) = r_{ij}^{-1} \sigma_{e(j)}^2$, which is the conditional variance of the genotype means computed in the first stage.
266 The estimates of genotype means, $\eta_i = \mu + g_i$, are identical for single-stage and two-stage analysis, provided the

267 variance components are known (Piepho et al., 2012). Differences arise in practice because variances need to
268 be estimated. Stage-wise analysis entails an approximation of the gold standard of single-stage analysis because
269 variances $\text{var}(\bar{y}_{ij.}) = r_{ij}^{-1} \sigma_{e(j)}^2$ as estimated in the first stage are treated as known quantities in the second stage,
270 disregarding the degrees of freedom associated with these estimates and their uncertainty. A key feature of
271 stage-wise analysis is that the inverses of these estimated variances act as weights in the second-stage analysis.
272 A major challenge in any stage-wise analysis is how to best determine the weights and how to account for the
273 uncertainties associated with them.

274 The situation faced in the analysis of HTFP is comparable in that it proceeds in stages with necessity because
275 a single-stage analysis is in conflict with performance and generalization demands (i.e., multi-year HTFP data
276 may comprise a number of differing experimental designs that require individual processing in a first stage) and
277 that the primary interest is the genotype main effect g_i , which equals θ_i in HTFP (Figure 21-n). The statistical
278 challenges are rather more daunting, however, for several reasons: (i) HTFP involves high-frequency time series in
279 which observations are serially correlated; (ii) summarizing time-series data usually requires nonlinear regression
280 models; (iii) analyses of field trials are often done exploiting spatial correlations among neighboring plots; (iv)
281 remote or proximate sensed data are affected by environmental conditions (wind, illumination) and may change
282 during the course of a measurement; (v) the number of stages required for the full analysis process is much
283 greater than two. These additional features make the determination of appropriate weights to be carried forward
284 from one stage to the next even more challenging than in the simple example given above.

285 Here, we propose a weighing approach for the intermediate trait category (1) (timing of key stages) and (2)
286 (quantities at defined time points or periods) only for brevity, and illustrate its application using a simulation
287 study described in the following section. Traits of the third category (dose-response curves) will be considered
288 in a follow-up paper.

289 2.4. Simulation of canopy height data

290 To demonstrate the extraction of traits of the first two categories (timing of key stages and quantities at
291 defined time point or periods), winter wheat canopy height data were simulated implementing a temperature
292 dose-response curve (trait category three, Equation 7). The temperature response of the stem elongation phase
293 was assumed to follow a dose-response curve with break points (Figure 4),

$$r_{BP}(T, \theta^C) = \begin{cases} 0, & T < T_{\min} \\ r_{\max}, & T > T_{\text{opt}} \\ r_{\max} \cdot \frac{T - T_{\min}}{T_{\text{opt}} - T_{\min}}, & \text{otherwise,} \end{cases} \quad (10)$$

294 where T_{\min} is the base temperature below which the elongation rate r is zero and T_{opt} the optimum temperature
295 above which the elongation rate reaches the maximum hourly elongation rate r_{\max} , while $\theta^C = (r_{\max}, T_{\min}, T_{\text{opt}})$
296 (Figure 4).

297 As starting point for the simulation, existing experimental designs of three consecutive years at the ETH
 298 research station of agricultural sciences in Lindau Eschikon, Switzerland (47.449 N, 8.682 E, 556 m a.s.l.) were
 299 used. The experiment consisted of 352 wheat genotypes, replicated twice per year on two spatially separated
 300 fields, both augmented with spatial checks in a 3×3 block arrangement.

301 To simulate canopy height time series, existing weather data were used to introduce a close-to-realistic
 302 stochastic behavior. Canopy growth was simulated for a measurement interval of one per day and for a pe-
 303 riod between first of March and 20th of July ($d = 1, \dots, 142$) for each of the three simulated years j ($j =$
 304 2016, 2017, 2018). Growth between daily campaign time points t was modeled as cumulative response to hourly
 305 temperature measurements T_{jdh} ($h = 1, \dots, 24$). The canopy height y_{ijkt} of genotype i ($i = 1, \dots, 352$) at plot k
 306 ($k = 1, \dots, 704$) in the year j at a specific time point ($t = 1, \dots, 142$) was then simulated as

$$y_{ijkt} = g_T(t, \theta_{ijk}^C, \theta_{ijk}^T; T_{jdh}) + e_{jkt}, \quad (11)$$

307 where g_T depends on r_{BP} in Equation 10 (see below) and simulates growth as a function of temperature T_{jdh} ,
 308 time point t , a vector of plot-specific crop growth model parameters $\theta_{ijk}^C = (r_{max}, T_{min}, T_{opt})$, and a vector of
 309 plot-specific timing traits $\theta_{ijk}^T = (tPH_{start}, tPH_{stop})$. The error term e_{jkt} simulates plot and time point residuals.
 310 The growth function g_T was specified as

$$g_T(t, \theta^C, \theta^T; T_{dh}) = \sum_{d=1}^t \begin{cases} \sum_{h=1}^{24} r_{BP}(T_{dh}, \theta^C) & tPH_{start} < d < tPH_{stop} \\ 0, & \text{otherwise} \end{cases}, \quad (12)$$

311 where r_{BP} represents a dose-response as function of hourly temperatures T_{dh} and a vector of crop growth model
 312 parameters θ^C (Equation 10), tPH_{start} the time point where canopy growth started, and tPH_{stop} the time point
 313 where canopy growth stopped.

314 This approach produced realistic-looking canopy growth curves (compare Figure 5 with, e.g., real data in
 315 [Kronenberg et al. 2017, 2020a](#)) with a characteristic start of growth (tPH_{start}) and a stop of growth (tPH_{stop}),
 316 corresponding to the first intermediate trait category (timing of key stages). Additionally, growth curves indi-

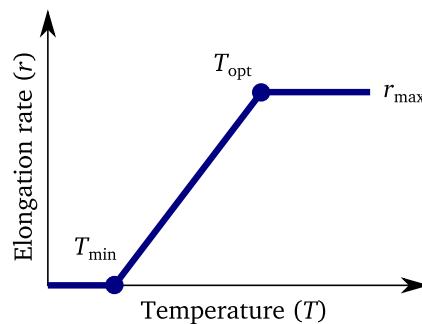


Figure 4: Schematic drawing of the dose-response model (\dot{g} of Equation 7) implemented as break-point model (r_{BP} , Equation 10) used for the simulation of canopy height time series based on temperature courses.

317 cated a characteristic final height (PH_{\max}), corresponding to the second intermediate trait category (quantities at
318 defined time points or periods).

319 Noise as specified in Section 2.1 was introduced on a genotype, plot and time point level. Genotype-year
320 interactions were not explicitly introduced, as it was assumed that they will emerge as the result of random θ_i^C
321 and θ_i^T combinations applied to year-specific temperature courses.

322 To add noise to genotype traits, the crop growth model parameters θ_{ijk}^C and the timing traits θ_{ijk}^T were further
323 decomposed in genotypic and spatially correlated parts,

$$\theta_{ijk}^C = \theta_i^C + \theta_{jk}^C \quad (13)$$

$$\theta_{ijk}^T = \theta_i^T + \theta_{jk}^T, \quad (14)$$

324 where θ_i^C and θ_i^T were simulated using normal distributions ($\sim \mathcal{N}(\mu, \sigma^2)$). Trait-specific μ and σ^2 were chosen
325 based on literature if available, and otherwise based on own unpublished field data. θ_{jk}^C and θ_{jk}^T were spatial
326 correlated heterogeneity components for those traits ($AR(1)_x \otimes AR(1)_y$), where $AR(1) \otimes AR(1)$ is a two-dimensional
327 first-order autoregressive model in row (x) and range (y) direction, mimicking the influence of other covariates
328 and therefore spatial heterogeneity. Note that a high autocorrelation in row and range direction with $\rho = 0.95$
329 and half the variance of the corresponding input parameter (Appendix Table 2) was assumed, which appeared
330 reasonable for cereal experiments (e.g. Velazco et al., 2017).

331 The plot residual e_{jkt} was simulated as sum of three error terms,

$$e_{jkt} = e_{jkt,1} + e_{jk} + e_{jkt,2}. \quad (15)$$

332 The first error term $e_{jkt,1}$ corresponds to the serial correlation of measurement errors ($AR(1)_t$) that g in Equation
333 1 presumably absorbs. The second error term e_{jk} mimics a systematic spatially correlated measurement error,
334 e.g., after an incomplete correction with reference measurements ($AR(1)_x \otimes AR(1)_y$). We note that adding this
335 error introduces an intentional discrepancy between the analysis model and the simulation: the proposed plot-
336 level model for repeated measurements does not include such a systematic error in the first stage (dynamic
337 modeling). Consequently, estimating the spatial correlation in the second stage will confound measurement errors
338 and nuisance factors, which corresponds to a situation we frequently encounter in HTFP. The third error term $e_{jkt,2}$
339 corresponds to e_{kt} in Equation 1 and represents a plot-based i.i.d. residual ($\sim \mathcal{N}(0, \sigma^2)$). The first error term was
340 assumed to cause most of the known measurement error, wherefore σ was set accordingly to 0.01 m (Roth et al.,
341 2020), while for the second and third error term σ was significantly reduced. The autocorrelation parameter ρ
342 was arbitrary set to 0.7. All simulation input parameters and sources for the aforementioned assumptions are
343 summarized in the Appendix (Table 2).

344 A total of 500 simulation runs were performed. These simulated time series with a measurement interval
345 of one day were then further thinned to intervals of three, five, seven and 14 days to study the effect of lower
346 frequencies.

347 We note that the simulation (Equation 11) comprised θ^T , i.e. traits of the first category, and θ^C , i.e. traits of
348 the third category. The second trait category θ^Q was dependent on the first and third category and year specific
349 temperature courses, and therefore only an indirect input parameter of the simulation. Therefore, the simulation
350 allowed extracting traits of all three categories, and validating traits of category one (θ^T) and three (θ^C) with
351 genotypic input data, and traits of category two (θ^Q) with plot-level (indirect) input data. Here, we illustrate the
352 extraction of θ^T and θ^Q only for brevity. The extraction of θ^C and therefore dose-response curve parameters of
353 a crop growth model will be considered in a follow-up paper.

354 We further note that all simulation input parameters for a given genotype i (θ_i^T and θ_i^C) were uncorrelated.
355 In reality, genetic effects and artificial selection have certainly resulted in weak to strong correlations for those
356 parameters. Dynamic modeling may introduce new, artificial correlations of parameters. When examining a real-
357 world genotype set, e.g., a breeding population, these effects will be confounded, but using a simulation with
358 uncorrelated input parameters allows quantifying the extraction artifacts.

359 2.5. Stage 1: Extracting the timing of key stages and quantities at define time points

360 To extract timing of key stages, a monotonically increasing P-spline was fitted to plot time series using the R
361 package *scam* (Pya, 2019), thus implementing $s_k(t)$ of Equation 3. The package fits shape constrained generalized
362 additive models (GAM) (Pya and Wood, 2015). A Bayesian approach to uncertainty quantification is used to
363 obtain standard errors of predictions.

364 The number of knots was set proportional to 3/4 of the observations. In a next step, the start and end of
365 stem elongation (tPH_{start} and tPH_{stop}) were extracted based on the quarter of maximum elongation rate (QMER)
366 method, which in brief extracts key time points with elongation rates greater than a threshold of 1/4 of the
367 maximum elongation rate. Thus, the QMER method represents an implementation of $q_n(s_k)$ of Equation 4.

368 In detail, in a first step spline predictions for canopy heights \hat{y}_t and standard error estimates $s.e.(\hat{y}_t)$ were
369 calculated separately for each plot at hourly time steps using the prediction function of the *scam* package. There-
370 after, hourly growth rates \hat{r}_t were derived from the difference between subsequent predictions, $\hat{r}_t = \hat{y}_t - \hat{y}_{t-1}$
371 (Figure 2e). Then, the following algorithm was applied to extract intermediate traits and corresponding weights
372 w based on standard errors of spline predictions:

373 1. Determine maximum elongation rate:

$$374 \hat{r}_{\text{max}} = \max(\hat{r}_t)$$

375 2. Filter \hat{r}_t for data points with an elongation rate greater than 1/4 of the maximum elongation rate:

$$376 \hat{r}_{t,\text{set1}} = \hat{r}_t \text{ where } \hat{r}_t \geq 1/4 \cdot \hat{r}_{\text{max}}$$

377 3. Define the earliest time points that is left after filtering as the start of growth:

$$378 tPH_{\text{start}} = t \text{ of first}(\hat{r}_{t,\text{set1}})$$

$$379 w_{tPH_{\text{start}}}^{-1/2} = s.e.(\hat{y}_t) \text{ where } t = tPH_{\text{start}}$$

380 4. Filter \hat{r}_t for data points with an elongation rate lower than $1/4$ of the maximum elongation rate and a
381 minimum distance of 40 days to the approximated start of growth:

$$382 \hat{r}_{t,\text{set2}} = \hat{r}_t \text{ where } \hat{r}_t \leq 1/4 \cdot \hat{r}_{\text{max}} \wedge t - t\text{PH}_{\text{start}} \geq 40$$

383 5. The earliest value that is left after filtering indicates the approximated end of growth:

$$384 t\text{PH}_{\text{stop}} = t \text{ of first}(\hat{r}_{t,\text{set2}})$$

$$385 w_{t\text{PH}_{\text{stop}}}^{-1/2} = \text{s.e.}(\hat{y}_t) \text{ where } t = t\text{PH}_{\text{stop}}$$

386 Note that the weights for timing of key stages' traits in this work were based on the standard errors of spline
387 predictions \hat{y} . We will address the conditions that should be met to justify our approach in the following section.

388 We extracted the growth stages start and end of stem elongation ($t\text{PH}_{\text{start}}$ and $t\text{PH}_{\text{stop}}$) and corresponding
389 standard error estimates based on the quarter of maximum elongation rate (QMER) method. To compare the
390 QMER method with the approach taken by [Kronenberg et al. \(2017\)](#), we additionally determined the time points
391 where 15% ($t\text{PH}_{15}$) and 95% ($t\text{PH}_{95}$) of final height was reached (for details, see [Kronenberg et al., 2017](#)). In
392 Figure 2e, we depict only the QMER method.

393 The quantity at a defined time point final height (PH_{max}) was calculated as the median of 24 hourly spline
394 predictions after the estimated stop of growth:

395 1. Filter \hat{y}_t for data points after reaching final height:

$$396 \hat{y}_{t,\text{final}} = \hat{y}_t \text{ where } t\text{PH}_{\text{stop}} \leq t \leq t\text{PH}_{\text{stop}} + 24 \text{ h}$$

397 2. Aggregate data points:

$$398 \text{PH}_{\text{max}} = \text{median}(\hat{y}_{t,\text{final}})$$

$$399 w_{\text{PH}_{\text{max}}}^{-1/2} = \text{s.e.}(\hat{y}_t) \text{ where } t = t\text{PH}_{\text{stop}}$$

400 2.6. Weighting based on estimated standard errors

401 The chosen implementation of the QMER method did not provide standard errors for the derived growth rate
402 (\hat{r}) and time points (t). Therefore, weighting for further processing after the dynamic modeling was based on
403 standard errors of spline-based predictions of the response ($\text{s.e.}(\hat{y}_t)$).

404 Using weights based on the standard errors of spline predictions is intuitive for quantities at defined time
405 points or periods' traits (e.g., PH_{max}), as both $\text{s.e.}(\hat{y}_t)$ and \hat{y}_t share the same unit. Nevertheless, for timing
406 of key stages (e.g., $t\text{PH}_{\text{start}}$ and $t\text{PH}_{\text{stop}}$), such a weighting approach requires a positive and high association
407 between the true weights for t and y for a given (to be determined) time point. Alternatively, one could use
408 an inverse regression approach (e.g., the Fieller's theorem ([Seber, 2003](#)) or the delta method ([Johnson et al.,](#)
409 [1993](#))) to determine weights for two means with different units. Such an inverse regression approach becomes
410 non-trivial when involving a combination of statistical tools—e.g., P-splines and the QMER method. Therefore,
411 using an inverse regression approach may contradict the requirement to provide a seamless workflow to integrate
412 arbitrary complex dynamic models g (Equation 2).

413 Consequently, we decided to assume proportionality of weights for standard errors of spline predictions and
414 timing of key stage estimations. The factor of proportionality was estimated via the residual variance (σ^2), which
415 was estimated in each analysis, rather than fixed at unity, as is customary in standard weighted analysis, where the
416 inverse weights are taken to be the known residual variances (Piepho et al., 2012). Our assumption is based on
417 considerations on a concrete example (see Appendix). In addition, standard errors of spline predictions suppose
418 that observations of plot-based time series are independent. As this is—at least for the simulation—not true
419 (see Section 2.1), the calculated standard errors of the estimates will be biased. To test whether weighting was
420 advantageous, despite possible bias in the weights and imperfect proportionality for timing of key stage traits,
421 we optionally provided the weights in the next processing step.

422 2.7. Stage 2: Calculating adjusted genotype means per year

423 The extraction of dynamics characteristics resulted in measurement time point independent trait values at a
424 plot level (Stage 1). These plot values were subsequently processed in a two-stage linear mixed model analysis
425 (Stage 2 and 3), where the second stage averaged over within-year effects (e.g., spatial heterogeneity) and the
426 third stage over between-year effects.

427 We used SpATS (Rodríguez-Álvarez et al., 2018) to fit a model with a smooth bivariate surface defined over
428 spatial coordinates of plot centers ($f(x(jk), y(jk))$) and added fixed genotype effects (θ_{ij}) and random effects of
429 plot rows and ranges ($p_{r(jk)}$ and $p_{c(jk)}$),

$$\hat{\theta}_{jk} = \theta_{ij} + f(x(jk), y(jk)) + p_{r(jk)} + p_{c(jk)} + e_{jk}. \quad (16)$$

430 Model parameters are listed and explained in Table 1 (Stage 2). Stage 2a and 2b are two nested models;
431 Stage 2b corresponds to Stage 2a but additionally includes weights. Equation 16 was applied to all intermediate
432 traits to calculate BLUEs of genotype means per year.

433 2.8. Stage 3: Genotypic marginal means calculation

434 The second stage already covered aspects such as spatial heterogeneity and design-specific characteristics
435 such as row and range arrangements, and allowed obtaining adjusted year genotype means $\hat{\theta}_{ij}$ (BLUEs). In the
436 third stage, those means were further processed with a model based on Equation 9,

$$\hat{\theta}_{ij} = \mu + \nu_j + \theta_i + (\theta\nu)_{ij} + e_{ij}. \quad (17)$$

437 The model assumes that genotype-environment effects can be partitioned into genotype response effects (θ_i) and
438 genotype-year interaction effects ($(\theta\nu)_{ij}$) (Piepho et al., 2012) while the residual errors (e_{ij}) are assumed to be
439 identically and independently normally distributed. Model parameters are listed and explained in Table 1 (Stage
440 3). Stage 3a and 3b are two nested models; Stage 3b corresponds to Stage 3a but additionally includes weights.
441 Models were fitted using the R package R-asreml (Butler, 2018).

Table 1: Model parameters for the second and third stage of the stage-wise linear mixed model analysis. k denotes the k -th plot, j the j -th year, and i the i -th genotype.

Stage	Term	Description	Part
2)	$\hat{\theta}_{jk}$	Plot response based on dynamic modeling	Response
	θ_{ij}	Year genotype response	Fixed
	$P_{c(jk)}$	Range effect on field (main working direction, e.g., for sowing)	Random
	$P_{r(jk)}$	Row effect on field (orthogonal to main working direction)	Random
	$f(x(jk), y(jk))$	Smooth bivariate surface in spatial x and y coordinates (mapping real distances on the field) consisting of a bivariate polynomial and a smooth part (for details see Rodríguez-Álvarez et al., 2018)	Spatial
a)	e_{jk}	Residuals with $\text{var}(e) = \sigma^2$	Residual
b)	e_{jk}	Residuals with $\text{var}(e) = \sigma^2 w^{-1}$, where w are weights based on the standard error estimations from the previous dynamic modeling step (Stage 1), and σ^2 the residual variance parameter	Weights
3)	$\hat{\theta}_{ij}$	Adjusted year genotype mean (BLUE) from Stage 2	Response
	μ	Global intercept	Fixed
	ν_j	Year effect	Random
	θ_i	Genotype response	Fixed
	$(\theta\nu)_{ij}$	Genotype year interaction	Residual
a)	e_{ij}	Residuals with $\text{var}(e) = \sigma^2$	Residual
b)	e_{ij}	Residuals with $\text{var}(e) = \sigma^2 w^{-1}$, where w are weights based on the square rooted diagonal of the variance-covariance matrix from Stage 2, and σ^2 the residual variance parameter	Weights

442 Separating the dynamic modeling step from further processing steps prevents implementing the gold standard
 443 of a one-stage analysis. Nevertheless, subsequent processing stages can be summarized in one stage, hence result-
 444 ing in a two-stage approach. To allow comparing such an approach with a three-stage approach, the estimated
 445 intermediate traits from Stage 1 were additionally processed using a two-stage model,

$$\hat{\theta}_{jk} = \mu + v_j + \theta_i + (\theta v)_{ij} + p_{r(jk)} + p_{c(jk)} + f(r(jk), c(jk)) + e_{jk}, \quad (18)$$

446 where μ is a global intercept, v_j a year intercept, θ_i the genotype response, $(\theta v)_{ij}$ genotype year interactions,
 447 $p_{r(jk)}$ and $p_{c(jk)}$ range and row effects, $f(r(jk), c(jk))$ year specific AR(1) \otimes AR(1) interactions based on ranges
 448 ($c()$) and rows ($r()$) of plots, and e_{jk} plot residuals with year-specific variances.

449 2.9. Simulation validation

450 Bias, variance, root-mean squared error (RMSE) and Pearson's correlation were calculated both after dynamic
 451 modeling (Stage 1) and after the stage-wise linear mixed model analysis (Stage 2 and 3) separately for each
 452 simulation run.

453 3. Results

454 A total of 176,000 genotypes replicated on 1,056,000 plots (500 runs \times 3 years \times 2 replications \times 352
 455 genotypes) containing 149,952,000 data points (number of plots \times 142 measurement days) were simulated. In
 456 the following, we give insights on the precision of extracted traits influenced by the choice of method, weighting,
 457 and measurement interval.

458 3.1. Dynamic modeling

459 P-splines model fits converged for all simulated plot time series and produces smooth-looking growth curves
 460 (Figure 5). Start and end of stem elongation estimations were successfully extracted using the QMER method as
 461 well as the final height percentile method.

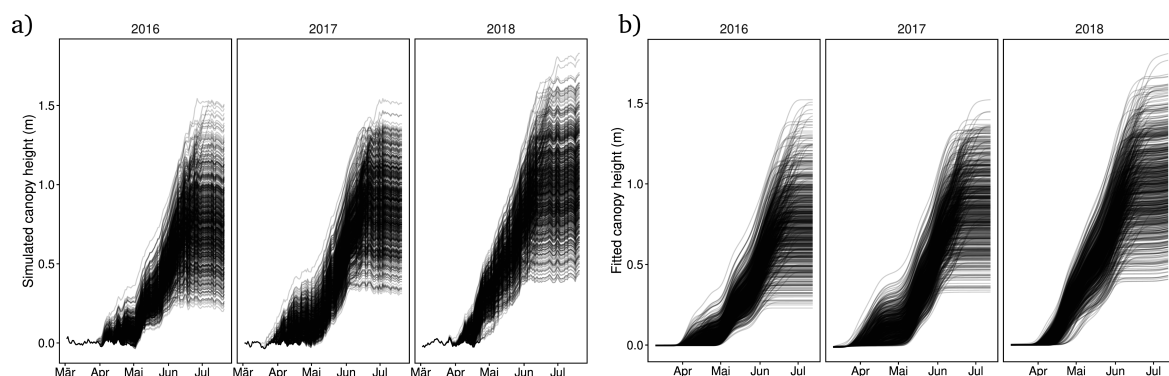


Figure 5: Simulated canopy heights (a) and fitted canopy height splines (b) for one simulation run with 352 genotypes, two replications per year, and three years.

462 The timing of the key stage trait tPH_{start} was better estimated by the P-spline/QMER method with a lower
463 median RMSE and lower median bias (Figure 6). Nevertheless, in comparison to the final height percentile
464 method, the median variance was higher, and larger outliers for RMSE and variance were found. The trait
465 tPH_{stop} was better estimated by the final height percentiles method with lower median bias, median RMSE and
466 median variance than by the P-spline/QMER method, but the percentiles method also produced larger outliers
467 for variance and RMSE than the P-spline/QMER method.

468 Both the P-spline/QMER and final height percentile methods performed comparably and were able to predict
469 tPH_{start} with a strong and tPH_{stop} with a very strong correlation to input values (Figure 7), but also for both
470 methods, the estimated start of stem elongation (tPH_{start}) was weakly biased by the input trait base temperature.
471 Nevertheless, the correlation between the extracted start and end of stem elongation—an artifact of the method,
472 as the simulation input was uncorrelated—was much higher for the Percentile method than for the P-spline/QMER
473 method. Based on these findings, the P-spline/QMER model was selected for further processing in the stage-wise
474 analysis.

475 3.2. Required measurement intervals

476 Estimating tPH_{stop} and PH_{max} using the P-spline/QMER or Percentile method was not affected by increased
477 or reduced measurement intervals unless reduced from 7 to 14 days, where the correlation for both tPH_{start} and
478 tPH_{stop} dropped (Figure 8). The estimation of tPH_{start} was, in contrast to the two other traits, sensitive to reduced
479 measurement intervals above five days for the P-spline/QMER method. The prediction of final height was not
480 affected by increased measurement intervals.

481 3.3. Stage-wise linear mixed model analysis

482 For both traits tPH_{start} and tPH_{stop} , calculating overall adjusted genotype means reduced the median variance
483 and median bias if compared to plot-based values for the P-spline/QMER method (Figure 6) and improved the
484 median RMSE for tPH_{start} but not for tPH_{stop} (Figure 9, Appendix Table 3). Based on bias and variance, for the
485 three-stage model (dynamic modeling followed by a two-stage linear mixed model analysis), weighting Stage 2
486 with errors of the prediction from dynamic modeling (Stage 1) was of advantage for tPH_{start} . The lowest median
487 bias was found for the combination of weighting Stage 2 as well as Stage 3 and the lowest median variance
488 for the combination of weighting Stage 3 but not Stage 2 (Figure 9, Appendix Table 3). For tPH_{stop} , median
489 differences between weighting combinations for Stage 2 and 3 were overall very small, but weighting stage 2
490 reduced outliers for bias and RMSE.

491 When compared to a two-stage model (dynamic modeling followed by a one-stage linear mixed model anal-
492 ysis), using a three-stage model was of disadvantage for tPH_{start} , indicated by a lower median RMSE and larger
493 outliers (Figure 9). For tPH_{stop} , the median RMSE was slightly higher for the two-stage model than for the
494 three-stage model, but outliers were more frequent for the three-stage model.

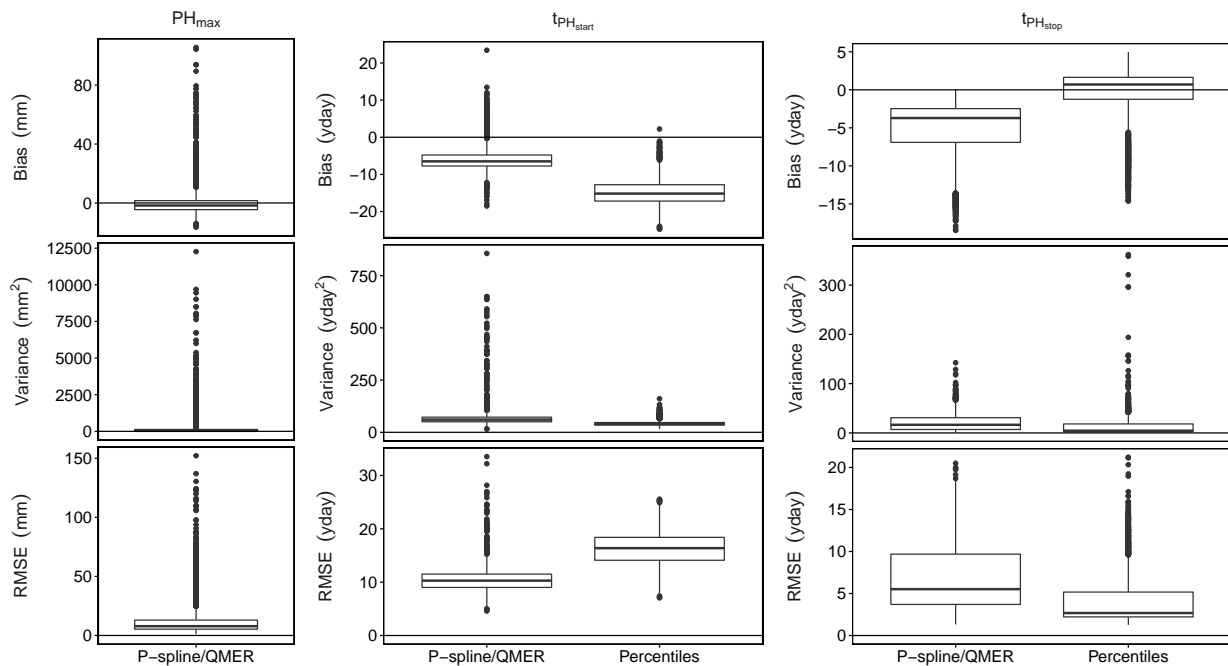


Figure 6: Box plots for the 500 simulated datasets of plot-based bias, variance and root-mean squared error (RMSE) of two timing of key stages models (P-spline/QMER model and final height percentiles).

	Simulation input						Percentiles		P-spline/QMER		
	r_{max}	T_{min}	T_{Opt}	PH_{max}	tPH_{start}	tPH_{stop}	tPH_{start}	tPH_{stop}	PH_{max}	tPH_{start}	tPH_{stop}
r_{max}	1	0	0	0.33	0	0	0	-0.02	0.33	0	0
T_{min}	0	1	-0.01	-0.33	0	0	0.15	0.04	-0.33	0.21	0.06
T_{Opt}	0	-0.01	1	-0.32	0	0	0.05	0.02	-0.32	0.07	0.02
PH_{max}	0.33	-0.33	-0.32	1	-0.33	0.53	-0.17	0.44	1	-0.32	0.36
tPH_{start}	0	0	0	-0.33	1	0.12	0.74	0.19	-0.33	0.73	0.24
tPH_{stop}	0	0	0	0.53	0.12	1	0.47	0.98	0.53	0.24	0.92

Figure 7: Pearson's correlations of plot time series traits. Provided are simulated input parameters and extracted timing of key stages' traits for the P-spline/QMER and final height percentile model. Black bold boxes indicate correlations between predicted and true values for identical traits.

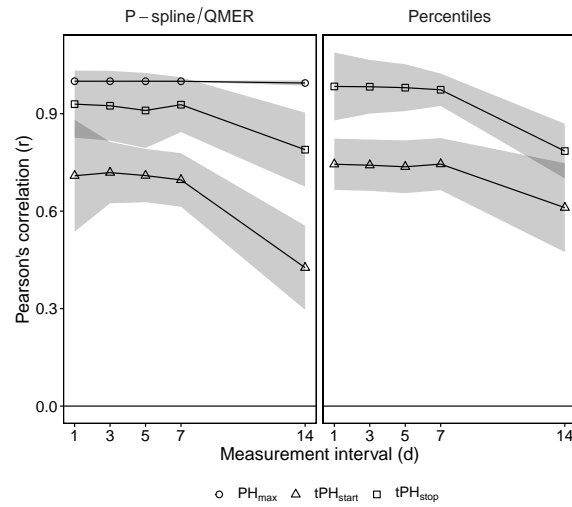


Figure 8: Pearson's correlations for differing measurement intervals for the timing of key stages based on splines (P-spline/QMER method) and final height percentiles (Percentiles method).

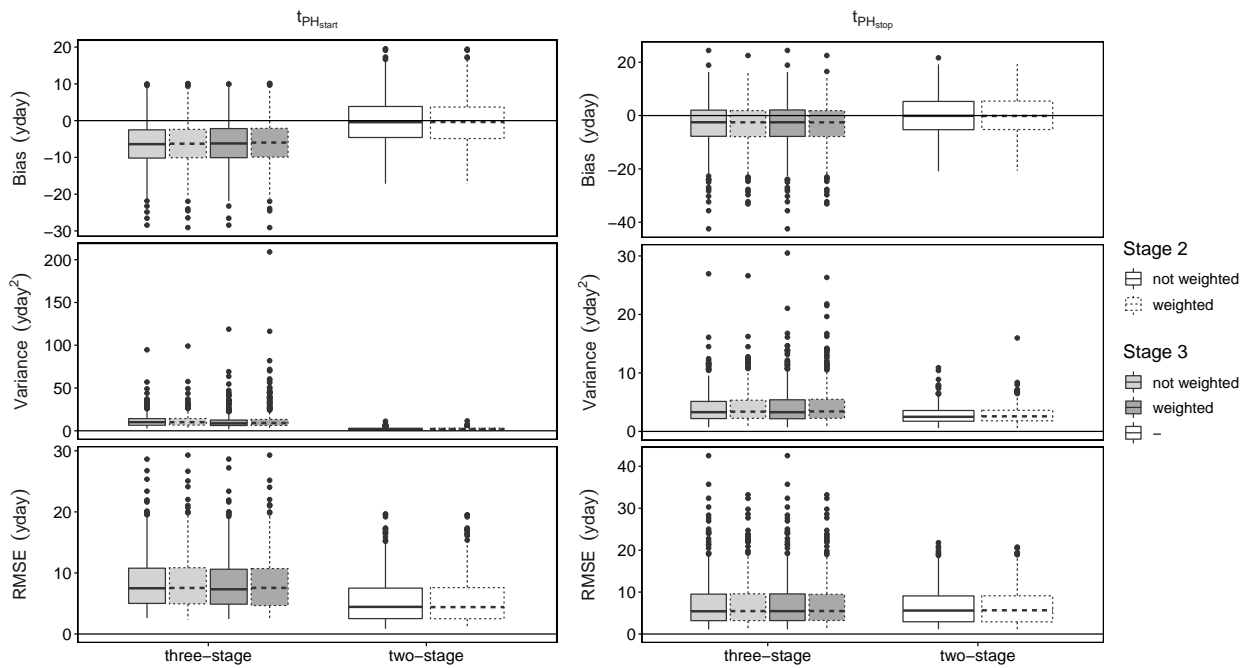


Figure 9: Box plots for the 500 simulated datasets of genotype based bias, variance, and root-mean squared error (RMSE) for the key stages P-spline/QMER model for the three-stage model and the two-stage model.

495 4. Discussion

496 4.1. Data processing in stages

497 The overall workflow of HTFP requires a joint effort of disciplines (Cobb et al., 2013; Araus and Cairns, 2014)
498 which may be separated into three main domains: (1) automation and sensing including feature extraction
499 from sensor readouts, (2) applied phenotyping including dynamic modeling and trait extraction from sensor-
500 derived features, and (3) analysis of designed agricultural experiments or breeding experiments. Plant phenomics
501 must bridge these three disciplines with the overall aim to characterize phenotypes as the result of genotype,
502 environment and management. A plot-level model for repeated measurements may help to link the highly specific
503 domains of sensing and the analysis of experiments. The link to genomic information in breeding and quantitative
504 genetics further increases the complexity of the topic, but is only marginally addressed in this study.

505 Here, we presented a strategy to process HTFP data. Based on the evaluated sources of variation, we decided
506 to process in stages, starting with dynamic modeling, followed by two stages of a linear mixed model analysis.
507 This approach is to some extent the reverse of van Eeuwijk et al. (2019) who suggested correcting time point
508 measurements in a first stage of a stage-wise linear mixed model analysis, followed by dynamic modeling and
509 modeling of environmental dependencies, and a second stage of a stage-wise linear mixed model analysis to cal-
510 culate adjusted means across years. Both options—correcting for spatial or temporal correlations first—represent
511 valid alternatives. In the present case, we decided not to correct for spatial gradients before dynamic modeling
512 for two reasons:

513 (1) Calculating adjusted genotype means in a first stage will correct for spatially correlated measurement
514 errors, but also for effects caused by start and lag phase variations, quantitative trait variations, and environment
515 variations due to start and lag phase variations (Figure 3). While correcting for measurement errors is a desired
516 effect, correcting for other effects will bias the result of dynamic modeling by altering the time point variances
517 of low-level traits. Using measurement references and correcting for day-to-day random measurement errors
518 outside the framework of the stage-wise processing may therefore be of advantage. (2) In a linear mixed model
519 analysis interlaced by a dynamic modeling part, weighting becomes a daunting task. In opposite, weighting is an
520 integral part of the stage-wise analysis strategy presented in this study.

521 For the P-spline/QMER method, processing multiple years using a linear mixed model analysis reduced the
522 variance and bias of predictions while slightly increasing the RMSE. Weighting the first stage recovered the RMSE.
523 For the second stage, using weights based on estimated variances to approximate the gold standard of a single-
524 stage analysis proved to be of advantage for all traits if using meaningful weights for the first stage as well. These
525 findings indicate that our assumption about dynamic modeling was justified: the spatio-temporal correlation
526 caused by unconsidered covariates yields spatial correlated intermediate traits $\vec{\theta}_{ijk}$. Nevertheless, using a one-
527 stage linear mixed model with an $AR(1) \otimes AR(1)$ autocorrelation structure outperformed the stage-wise approach
528 for tPH_{start} and to some extent for tPH_{stop} .

529 Providing robust and reusable analysis routines represented an essential objective of the proposed approach.
530 The resulting generalization requirements may be in conflict with well-established analysis principles. This
531 conflict became well visible when formulating a linear mixed model for Stage 2: The philosophy “analyse-as-
532 randomise” would require to include all randomization factors—e.g., incomplete blocks—in the analysis. A gen-
533 eralized model as used in this work that includes besides a smooth bivariate surface just row and range effects
534 is certainly less efficient, but may nevertheless be suitable to draw correct conclusions on the outcome of the ex-
535 periment. Proposing a robust and reusable processing workflow therefore always represents a trade-off between
536 generalization and most efficient modeling.

537 4.2. *Intermediate trait categories*

538 In this study, we proposed three different trait categories: (1) Timing of key stages, (2) quantities at defined
539 time points or periods, and (3) dose-response curves. A fundamental difference between traits of the first two
540 categories and dose-response curve traits is how they include covariate dependencies. Dose-response curve traits
541 describe an explicit dependency on covariates. In contrast, timing of key stages’ traits include the effects of
542 covariates implicitly through the dependency on the timescale: Favorable conditions in spring may for example
543 accelerate the development of plants and therefore early key stages. Quantities at defined time points or periods’
544 traits may show a similar behavior, but here the directions are less clear: Early jointing in cereals due to favorable
545 conditions in spring may for example reduce the early canopy cover in the corresponding phase because of a
546 reduced growing time span. Nevertheless, one may also argue that favorable conditions in this reduced time
547 span may increase canopy cover. Both categories have in common that they describe an implicit reaction to a set
548 of covariate courses.

549 Consequently, to analyze traits of the first two categories, one reduces growing seasons with their charac-
550 teristic covariate courses to environments (E) and quantifies the influence of genotypes (G) and environments
551 on measured traits in a subsequent $G \times E$ analysis (for an overview see [van Eeuwijk et al., 2016](#)). In contrast,
552 dose-response curve traits are less affected by—but rather drivers of— $G \times E$. This difference may require differing
553 processing steps. We will cover dose-response curves in a follow-up paper.

554 4.3. *Limitations of dynamic modeling*

555 Clear limitations of the proposed approach became visible: Although all input parameters of the simulation
556 were completely uncorrelated, the extracted traits were to varying extents correlated. The simulation consisted of
557 500 independent simulation runs, and correlations were aggregated over all runs. Therefore, the observed effects
558 are a systematic result of the extraction methods and should be seen as corresponding limitations. When using
559 P-splines to extract key points of the stem elongation, the estimated start of the stem elongation may be biased
560 by the base temperature of growth. Nevertheless, this effect presumably applies to any other method including
561 the Percentile method, as both early start and low base temperature may result in a comparable phenotype in
562 early stages.

563 An increased length of the measurement interval may save considerable time and labor costs which may be
564 invested in larger number of tested genotypes. If aiming to extract timing of key stages, high frequencies are
565 to some extent superfluous if using P-splines, as the spline approach is presumably able to interpolate critical
566 measurement time points. Therefore, one to two measurements a week are sufficient, providing that the total
567 number of measurements does not drop below eight data points (as fitting a shape constrained P-spline using the
568 *scam* package to a time series with less than eight data points becomes to our experience challenging).

569 4.4. Limitations of processing in stages

570 A salient feature of our suggested approach is to proceed in several stages, starting with an analysis of time
571 series per plot. Because of this feature, our approach does not explicitly account for gross day-dependent errors
572 operating across all plots, although such errors represent an issue in real field data (Kronenberg et al., 2020b).
573 Explicitly accounting for such errors while also modelling the temporal trajectory would require joint spatio-
574 temporal modelling of the time series across all plots simultaneously. There are several approaches for spatio-
575 temporal modelling of environmental data that could be used here. As we are using splines for modelling both
576 the temporal and the spatial dimension, the most immediate option would be to use three-dimensional tensor
577 spline smoothing (Wood, 2017; Verbyla et al., 2018; Pérez et al., 2020). However, most of these are rather
578 more complex and computationally demanding and as such less suited for a seamless implementation for routine
579 analysis.

580 5. Conclusion

581 Processing repeated plot-level measurements using a well-defined process and data model revealed insights
582 on best practice in phenomics data handling. The results confirmed that HTFP measurements allow extracting
583 genotype specific timing of key stages and quantities at defined time points. P-splines combined with the QMER
584 method allowed extracting the timing of key stages and quantities at define time points with a precision that is
585 suitable for, e.g., plant breeding purposes.

586 Weighting turned out to be essential if processing HTFP data in stages, and linear mixed model analysis was
587 suitable to account for heterogeneity introduced by not considered covariates. Clear restrictions of the proposed
588 data processing strategy became obvious: Correlations between extracted traits cannot only arise from data, but
589 also from the extraction method itself. Therefore, care has to be taken when interpreting such correlations.

590 Yet, overall, the scientific community dealing with crop phenotyping has not come up with generally accepted
591 procedures how to organize the workflow from raw data generation to extraction of physiologically meaningful
592 results. Hopefully, the herein introduced modeling framework can contribute to achieving this aim; not only for
593 the merit of increased scholarly knowledge generation, but in the interest of a more efficient workflow for crop
594 breeding to improve global nutrition aspects in times of climate change.

595 **Appendix**

596 *5.1. Table: Simulation input parameters*

Table 2: Model input parameters for the simulation

	Distribution	Values	Sources
θ_t^C	$\mathcal{N}(\mu, \sigma^2)$	$T_{\min}: \mu = 8, \sigma = 2$ $T_{\text{opt}}: \mu = 18, \sigma = 2$ $r_{\max}: \mu = 0.9, \sigma = 0.2$	Kemp and Blacklow (1982) Kemp and Blacklow (1982) Own data
θ_{jk}^C	$\text{AR}(1)_{\text{row}} \otimes \text{AR}(1)_{\text{range}}$	$\rho = 0.95, \sigma_T = \frac{\sigma}{2\sqrt{2}}$	Velazco et al. (2017)
θ_t^T	$\mathcal{N}(\mu, \sigma^2)$	2016: $\mu_{t\text{PH}_{\text{start}}} = 108, \sigma_{t\text{PH}_{\text{start}}} = 2.8$ 2017: $\mu_{t\text{PH}_{\text{start}}} = 103, \sigma_{t\text{PH}_{\text{start}}} = 3.0$ 2018: $\mu_{t\text{PH}_{\text{start}}} = 101, \sigma_{t\text{PH}_{\text{start}}} = 3.1$ 2016: $\mu_{t\text{PH}_{\text{stop}}} = 165, \sigma_{t\text{PH}_{\text{stop}}} = 2.5$ 2017: $\mu_{t\text{PH}_{\text{stop}}} = 162, \sigma_{t\text{PH}_{\text{stop}}} = 3.5$ 2018: $\mu_{t\text{PH}_{\text{stop}}} = 158, \sigma_{t\text{PH}_{\text{stop}}} = 4.0$	Kronenberg et al. (2020a) Kronenberg et al. (2020a) Own data Kronenberg et al. (2020a) Kronenberg et al. (2020a) Own data
θ_{jk}^T	$\text{AR}(1)_{\text{row}} \otimes \text{AR}(1)_{\text{range}}$	$\rho = 0.95, \sigma = \frac{\sigma_{\text{PH}}}{2\sqrt{2}}$	Velazco et al. (2017)
$e_{jkt,1}$	$\text{AR}(1)_t$	$\rho = 0.7, \sigma_m = 0.01$	Roth et al. (2020)
e_{jk}	$\text{AR}(1)_{\text{row}} \otimes \text{AR}(1)_{\text{range}}$	$\rho = 0.7, \sigma = \frac{\sigma_m}{50}$	Assumption
$e_{jkt,2}$	$\mathcal{N}(\mu, \sigma^2)$	$\mu = 0, \sigma = \frac{\sigma_m}{100}$	Assumption

597 *5.2. Table: Median bias, variance and root-mean squared errors for the P-spline/QMER method*

Table 3: Genotype based bias, variance, and root-mean squared error (RMSE) for the key stages obtained using the P-spline/QMER method, with weighting as option for the second and third stage of the three-stage model, and weighting as option for the second stage of the two-stage model. Results report the median values over the 500 simulated datasets. For sake of completeness, plot-based median values for the P-spline/QMER method are reported as well.

Trait	Model	Weighted?	Bias	Variance	RMSE
		Stage 2 Stage 3	(yday)	(yday ²)	(yday)
$t\text{PH}_{\text{start}}$	Plot-based		-6.8	57.6	10.3
	Three-stage	no no	-6.4	10.1	7.49
		no yes	-6.21	8.68	7.32
		yes no	-6.27	10	7.54
		yes yes	-5.97	9.09	7.55
	Two-stage	no -	-0.417	2.01	4.44
yes -		-0.37	2.07	4.41	
$t\text{PH}_{\text{stop}}$	Plot-based		-2.85	11.5	4.43
	Three-stage	no no	-2.52	3.29	5.42
		no yes	-2.52	3.28	5.46
		yes no	-2.53	3.38	5.47
		yes yes	-2.54	3.42	5.49
	Two-stage	no -	-0.11	2.5	5.6
yes -		-0.139	2.58	5.66	

598 *5.3. A thought on weighting for traits of the second category (timing of key stages)*

599 Splines can be thought of as polynomials, or other functions that are linear in the regression parameters,
 600 pieced together at the knots. Thus, to gain some insight, we here consider a quadratic polynomial as a simple
 601 concrete example: $f(t) = \mu + \beta_1 t + \beta_2 t^2$. We observe data $y_i(t) = f(t) + e_i$ ($i = 1, \dots, n$), where $e_i \sim \mathcal{N}(0, \sigma^2)$.
 602 The model is linear and can be written in general for as $y = X\beta + e$, where $e \sim \mathcal{MVN}(0, I_n \sigma^2)$. Parameters can
 603 be estimated by ordinary least squares using $\hat{\beta} = (X^T X)^{-1} X^T y$ with

$$\text{var}(\hat{\beta}) = (X^T X)^{-1} \sigma^2. \quad (19)$$

604 A prediction at a particular value of t is obtained from $\hat{y}(t) = \hat{f}(t) = k^T \hat{\beta}$ with $k^T = (1 \quad t \quad t^2)$, and this has
 605 variance

$$\text{var}(k^T \hat{\beta}) = k^T (X^T X)^{-1} k \sigma^2. \quad (20)$$

606 Now assume that the aim is to find the value of t at which the response $f(t)$ is maximized. For simplicity,
 607 we take for granted that a maximum indeed occurs in the relevant range for t . At the maximum, the slope of
 608 the curve, i.e. the first derivative equals zero. This can be used to determine the optimal input level: $\frac{\partial f(t)}{\partial t} =$
 609 $\beta_1 + 2\beta_2 t = 0 \Leftrightarrow t_{\text{opt}} = -\frac{\beta_1}{2\beta_2}$. This can be estimated by $\hat{t}_{\text{opt}} = -\frac{\hat{\beta}_1}{2\hat{\beta}_2}$.

610 Now what can be said about the variance of this estimator, which would be needed for weighting? Here, we
 611 may use the delta method (Johnson et al., 1993) to find

$$\text{var}(\hat{t}_{\text{opt}}) \approx \left(\frac{\partial t_{\text{opt}}}{\partial \beta_1} \right)_{\beta_1=\hat{\beta}_1}^2 \text{var}(\hat{\beta}_1) + \left(\frac{\partial t_{\text{opt}}}{\partial \beta_2} \right)_{\beta_2=\hat{\beta}_2}^2 \text{var}(\hat{\beta}_2) + 2 \left(\frac{\partial t_{\text{opt}}}{\partial \beta_1} \right)_{\beta_1=\hat{\beta}_1} \left(\frac{\partial t_{\text{opt}}}{\partial \beta_2} \right)_{\beta_2=\hat{\beta}_2} \text{cov}(\hat{\beta}_1, \hat{\beta}_2). \quad (21)$$

612 From Equation 19, this is a linear function of σ^2 . Now Equation 20 is also linear in σ^2 . This suggests that
 613 the weights for \hat{t}_{opt} will be positively associated with those for $\hat{y}(t_{\text{opt}})$. Exact proportionality cannot be expected,
 614 however, because whereas $k^T (X^T X)^{-1} k$ in Equation 20 is constant across plots, the variance in Equation 21
 615 depends on regression parameters that are plot-specific. However, so long as these parameters are not very
 616 variable between plots, the association between weights for \hat{t}_{opt} and $\hat{y}(t_{\text{opt}})$ may be expected to be positive and
 617 high.

618 **Acknowledgement**

619 We like to acknowledge Helge Aasen, Lukas Kronenberg and Norbert Kirchgessner (ETH Zurich) for feedback
 620 on an early version of the manuscript.

621 **Funding**

622 LR received funding from Innosuisse (<http://www.innosuisse.ch>) in the framework for the project
 623 “Trait spotting” (grant number: KTI P-Nr 27059.2 PFLS-LS). MXRA was supported by project MTM2017-82379-
 624 R (AEI/FEDER, UE), by the Basque Government through the BERC 2018-2021 program, and by the Spanish

625 Ministry of Science, Innovation, and Universities (BCAM Severo Ochoa accreditation SEV-2017-0718). HPP was
626 supported by DFG grant PI 377/24-1.

627 **Declaration of Competing Interest**

628 The authors declare no conflict of interest.

629 **CRediT authorship contribution statement**

630 **Lukas Roth:** Conceptualization, Methodology, Software, Formal analysis, Visualization, Writing - Original
631 Draft. **María Xosé Rodríguez-Álvarez:** Software, Writing - Review & Editing **Fred van Eeuwijk:** Writing -
632 Review & Editing. **Hans-Peter Piepho:** Conceptualization, Methodology, Writing - Original Draft, Review &
633 Editing. **Andreas Hund:** Conceptualization, Supervision, Project administration, Funding acquisition, Writing -
634 Review & Editing.

635 References

- 636 H. Aasen, N. Kirchgessner, A. Walter, and F. Liebisch. PhenoCams for Field Phenotyping: Using Very High Temporal Resolution Digital
637 Repeated Photography to Investigate Interactions of Growth, Phenology, and Harvest Traits. *Frontiers in Plant Science*, 11(593), 2020.
638 ISSN 1664462X. doi: 10.3389/fpls.2020.00593.
- 639 J. Anderegg, K. Yu, H. Aasen, A. Walter, F. Liebisch, and A. Hund. Spectral Vegetation Indices to Track Senescence Dynamics in Diverse Wheat
640 Germplasm. *Frontiers in Plant Science*, 10(1749), 2020. ISSN 1664462X. doi: 10.3389/fpls.2019.01749.
- 641 J. L. Araus and J. E. Cairns. Field high-throughput phenotyping: The new crop breeding frontier. *Trends in Plant Science*, 19(1):52–61, 2014.
642 ISSN 13601385. doi: 10.1016/j.tplants.2013.09.008.
- 643 J. L. Araus, S. C. Kefauver, M. Zaman-Allah, M. S. Olsen, and J. E. Cairns. Translating High-Throughput Phenotyping into Genetic Gain.
644 *Trends in Plant Science*, 23(5):451–466, 2018. ISSN 13601385. doi: 10.1016/j.tplants.2018.02.001.
- 645 J. Blancon, D. Dutartre, M. H. Tixier, M. Weiss, A. Comar, S. Praud, and F. Baret. A high-throughput model-assisted method for phenotyping
646 maize green leaf area index dynamics using unmanned aerial vehicle imagery. *Frontiers in Plant Science*, 10(685), 2019. ISSN 1664462X.
647 doi: 10.3389/fpls.2019.00685.
- 648 M. Borenstein, L. V. Hedges, J. P. Higgins, and H. R. Rothstein. *Introduction to Meta-Analysis*. 2009. ISBN 978-0-470-05724-7.
- 649 I. Borra-Serrano, T. De Swaef, P. Quataert, J. Aper, A. Saleem, W. Saeys, B. Somers, I. Roldán-Ruiz, and P. Lootens. Closing the Phenotyping
650 Gap: High Resolution UAV Time Series for Soybean Growth Analysis Provides Objective Data from Field Trials. *Remote Sensing*, 12(1644),
651 2020. ISSN 20724292. doi: 10.3390/rs12101644.
- 652 D. Bustos-Korts, M. P. Boer, M. Malosetti, S. Chapman, K. Chenu, B. Zheng, and F. A. van Eeuwijk. Combining Crop Growth Modeling
653 and Statistical Genetic Modeling to Evaluate Phenotyping Strategies. *Frontiers in Plant Science*, 10(1491), 2019. ISSN 1664462X. doi:
654 10.3389/fpls.2019.01491.
- 655 D. Butler. *asreml: Fits the Linear Mixed Model*, 2018. URL www.vsnl.co.uk. R package version 4.1.0.93.
- 656 J. N. Cobb, G. DeClerck, A. Greenberg, R. Clark, and S. McCouch. Next-generation phenotyping: requirements and strategies for enhancing
657 our understanding of genotype-phenotype relationships and its relevance to crop improvement. *Theoretical and Applied Genetics*, 126:
658 867–887, 2013. ISSN 00405752. doi: 10.1007/s00122-013-2066-0.
- 659 W. G. Cochran. The Combination of Estimates from Different Experiments. *Biometrics*, 10(1):101–129, 1954.
- 660 J. A. Fernandez-Gallego, S. C. Kefauver, N. A. Gutiérrez, M. T. Nieto-Taladriz, and J. L. Araus. Wheat ear counting in-field conditions: High
661 throughput and low-cost approach using RGB images. *Plant Methods*, 14(22), 2018. ISSN 17464811. doi: 10.1186/s13007-018-0289-4.
- 662 C. Grieder, A. Hund, and A. Walter. Image based phenotyping during winter: a powerful tool to assess wheat genetic variation in growth
663 response to temperature. *Functional Plant Biology*, 42:387–396, 2015. ISSN 1445-4408. doi: 10.1071/fp14226.
- 664 M. Herrero-Huerta, P. Rodriguez-Gonzalez, and K. M. Rainey. Yield prediction by machine learning from UAS-based multi-sensor data fusion
665 in soybean. *Plant Methods*, 16(78), 2020. ISSN 17464811. doi: 10.1186/s13007-020-00620-6.
- 666 P. X. Hurtado, S. K. Schnabel, A. Zaban, M. Veteläinen, E. Virtanen, P. H. Eilers, F. A. van Eeuwijk, R. G. Visser, and C. Maliepaard. Dynamics
667 of senescence-related QTLs in potato. *Euphytica*, 183(3):289–302, 2012. ISSN 00142336. doi: 10.1007/s10681-011-0464-4.
- 668 X. Jin, S. Madec, D. Dutartre, B. de Solan, A. Comar, and F. Baret. High-Throughput Measurements of Stem Characteristics to Estimate Ear
669 Density and Above-Ground Biomass. *Plant Phenomics*, 2019(4820305), 2019. doi: 10.34133/2019/4820305.
- 670 N. L. Johnson, S. Kotz, and A. Kemp. *Univariate discrete distributions*. Wiley series in probability and mathematical statistics. Probability and
671 mathematical statistics. J. Wiley, New York, 2nd edition, 1993. ISBN 0471548979.
- 672 D. R. Kemp and W. M. Blacklow. The responsiveness to temperature of the extension rates of leaves of wheat growing in the field under
673 different levels of nitrogen fertilizer. *Journal of Experimental Botany*, 33(132):29–36, 1982. ISSN 00220957. doi: 10.1093/jxb/33.1.29.
- 674 N. Kirchgessner, F. Liebisch, K. Yu, J. Pfeifer, M. Friedli, A. Hund, and A. Walter. The ETH field phenotyping platform FIP: A cable-suspended
675 multi-sensor system. *Functional Plant Biology*, 44:154–168, 2017. doi: 10.1071/FP16165.
- 676 L. Kronenberg, K. Yu, A. Walter, and A. Hund. Monitoring the dynamics of wheat stem elongation: genotypes differ at critical stages. *Euphytica*,
677 213(157), 2017. doi: 10.1007/s10681-017-1940-2.
- 678 L. Kronenberg, S. Yates, M. P. Boer, N. Kirchgessner, A. Walter, and A. Hund. Temperature response of wheat affects final height and the
679 timing of stem elongation under field conditions. *Journal of Experimental Botany*, 2020a. doi: 10.1093/jxb/eraa471.
- 680 L. Kronenberg, S. Yates, S. Ghiasi, L. Roth, M. Friedli, M. E. Ruckle, R. A. Werner, F. Tschurr, M. Binggeli, N. Buchmann, B. Studer, and A. Walter.
681 Rethinking temperature effects on leaf growth, gene expression and metabolism: Diel variation matters. *Plant, Cell and Environment*, pages
682 1–15, 2020b. ISSN 0140-7791. doi: 10.1111/pce.13958.
- 683 I.-Y. Kwak, C. R. Moore, E. P. Spalding, and K. W. Broman. Mapping quantitative trait loci underlying function-valued traits using functional
684 principal component analysis and multi-trait mapping. *G3: Genes, Genomes, Genetics*, 6:79–86, 2016. ISSN 21601836. doi: 10.1534/g3.
685 115.024133.
- 686 C.-X. Ma, C. George, and W. Rongling. Functional mapping of quantitative trait loci underlying the character process: A theoretical framework.
687 *Genetics*, 161:1751–1762, 2002. ISSN 00166731.
- 688 M. Maimaitijiang, V. Sagan, P. Sidike, S. Hartling, F. Esposito, and F. B. Fritschi. Soybean yield prediction from UAV using multimodal data
689 fusion and deep learning. *Remote Sensing of Environment*, 237(111599), 2020. ISSN 00344257. doi: 10.1016/j.rse.2019.111599.
- 690 M. Malosetti, R. G. Visser, C. Celis-Gamboa, and F. A. van Eeuwijk. QTL methodology for response curves on the basis of non-linear mixed
691 models, with an illustration to senescence in potato. *Theoretical and Applied Genetics*, 113(2):288–300, 2006. ISSN 00405752. doi:
692 10.1007/s00122-006-0294-2.
- 693 E. J. Millet, W. Kruijer, A. Coupel-Ledru, S. Alvarez Prado, L. Cabrera-Bosquet, S. Lacube, A. Charcosset, C. Welcker, F. van Eeuwijk, and
694 F. Tardieu. Genomic prediction of maize yield across European environmental conditions. *Nature Genetics*, 51(6):952–956, 2019. ISSN
695 15461718. doi: 10.1038/s41588-019-0414-y.
- 696 F. F. Moreira, H. R. Oliveira, J. J. Volenec, K. M. Rainey, and L. F. Brito. Integrating High-Throughput Phenotyping and Statistical Genomic
697 Methods to Genetically Improve Longitudinal Traits in Crops. *Frontiers in Plant Science*, 11(681), 2020. ISSN 1664462X. doi: 10.3389/
698 fpls.2020.00681.

- 699 B. Parent, E. J. Millet, and F. Tardieu. The use of thermal time in plant studies has a sound theoretical basis provided that confounding effects
700 are avoided. *Journal of Experimental Botany*, 70(9):2359–2370, 2019. ISSN 14602431. doi: 10.1093/jxb/ery402.
- 701 D. M. Pérez, M. X. Rodríguez-Álvarez, M. P. Boer, E. J. Millet, and F. A. van Eeuwijk. Spatio-temporal and hierarchical modelling of high-
702 throughput phenotypic data. In *Proceedings of the 35th International Workshop on Statistical Modelling : July 20- 24, 2020 Bilbao, Basque*
703 *Country, Spain*, pages 394–397, Bilbao, 2020. URL <http://hdl.handle.net/10810/45863>.
- 704 G. Perich, A. Hund, J. Anderegg, L. Roth, M. P. Boer, A. Walter, F. Liebisch, and H. Aasen. Assessment of Multi-Image Unmanned Aerial Vehicle
705 Based High-Throughput Field Phenotyping of Canopy Temperature. *Frontiers in Plant Science*, 11(150), 2020. ISSN 1664-462X. doi:
706 10.3389/fpls.2020.00150.
- 707 H. P. Piepho, J. Möhring, T. Schulz-Streeck, and J. O. Ogutu. A stage-wise approach for the analysis of multi-environment trials. *Biometrical*
708 *Journal*, 54(6):844–860, 2012. ISSN 15214036. doi: 10.1002/bimj.201100219.
- 709 H. Poorter, F. Fiorani, R. Pieruschka, T. Wojciechowski, W. H. van der Putten, M. Kleyer, U. Schurr, and J. Postma. Pampered inside, pestered
710 outside? Differences and similarities between plants growing in controlled conditions and in the field. *New Phytologist*, 212:838–855,
711 2016. ISSN 14698137. doi: 10.1111/nph.14243.
- 712 N. Pya. *scam: Shape Constrained Additive Models*, 2019. URL <https://CRAN.R-project.org/package=scam>. R package version 1.2-5.
- 713 N. Pya and S. N. Wood. Shape constrained additive models. *Statistics and Computing*, 25:543–559, 2015. ISSN 15731375. doi: 10.1007/
714 s11222-013-9448-7.
- 715 J. Ramirez-Villegas, J. Watson, and A. J. Challinor. Identifying traits for genotypic adaptation using crop models. *Journal of Experimental*
716 *Botany*, 66(12):3451–3462, 2015. ISSN 14602431. doi: 10.1093/jxb/erv014.
- 717 G. J. Rebetzke, J. Jimenez-Berni, R. A. Fischer, D. M. Deery, and D. J. Smith. Review: High-throughput phenotyping to enhance the use of
718 crop genetic resources. *Plant Science*, 282:40–48, 2019. doi: 10.1016/j.plantsci.2018.06.017.
- 719 M. Reymond, B. Muller, A. Leonardi, A. Charcosset, and F. Tardieu. Combining Quantitative Trait Loci Analysis and an Ecophysiological
720 Model to Analyze the Genetic Variability of the Responses of Maize Leaf Growth to Temperature and Water Deficit. *Plant Physiology*, 131:
721 664–675, 2003. ISSN 0032-0889. doi: 10.1104/pp.013839.
- 722 M. X. Rodríguez-Álvarez, M. P. Boer, F. A. van Eeuwijk, and P. H. Eilers. Correcting for spatial heterogeneity in plant breeding experiments
723 with P-splines. *Spatial Statistics*, 23:52–71, 2018. ISSN 22116753. doi: 10.1016/j.spasta.2017.10.003.
- 724 L. Roth, M. Camenzind, H. Aasen, L. Kronenberg, C. Barendregt, K.-H. Camp, A. Walter, N. Kirchgessner, and A. Hund. Repeated Multiview
725 Imaging for Estimating Seedling Tiller Counts of Wheat Genotypes Using Drones. *Plant Phenomics*, 2020(3729715), 2020. doi: 10.34133/
726 2020/3729715.
- 727 P. Sadeghi-Tehran, K. Sabermanesh, N. Virlet, and M. J. Hawkesford. Automated Method to Determine Two Critical Growth Stages of Wheat:
728 Heading and Flowering. *Frontiers in Plant Science*, 8(252), 2017. doi: 10.3389/fpls.2017.00252.
- 729 G. A. F. G. A. F. Seber. *Linear regression analysis*. Wiley series in probability and statistics. Wiley-Interscience, Hoboken, N.J., 2nd ed. edition,
730 2003. ISBN 1-280-58916-7.
- 731 A. B. Smith, B. R. Cullis, and R. Thompson. The analysis of crop cultivar breeding and evaluation trials: An overview of current mixed model
732 approaches. *Journal of Agricultural Science*, 143:449–462, 2005. ISSN 00218596. doi: 10.1017/S0021859605005587.
- 733 A. Soltani and S. Galeshi. Importance of rapid canopy closure for wheat production in a temperate sub-humid environment: Experimentation
734 and simulation. *Field Crops Research*, 77:17–30, 2002. ISSN 03784290. doi: 10.1016/S0378-4290(02)00045-X.
- 735 F. Tardieu, L. Cabrera-Bosquet, T. Pridmore, and M. Bennett. Plant Phenomics, From Sensors to Knowledge. *Current Biology*, 27:R770–R783,
736 2017. ISSN 09609822. doi: 10.1016/j.cub.2017.05.055.
- 737 J. Ubbens, M. Cieslak, P. Prusinkiewicz, and I. Stavnness. Latent Space Phenotyping: Automatic Image-Based Phenotyping for Treatment
738 Studies. *Plant Phenomics*, 2020(5801869), 2020. doi: 10.1101/557678.
- 739 F. A. van Eeuwijk, D. Bustos-Korts, and M. Malosetti. What should students in plant breeding know about the statistical aspects of genotype
740 x Environment interactions? *Crop Science*, 56(5):2119–2140, 2016. ISSN 14350653. doi: 10.2135/cropsci2015.06.0375.
- 741 F. A. van Eeuwijk, D. Bustos-Korts, E. J. Millet, M. P. Boer, W. Kruijjer, A. Thompson, M. Malosetti, H. Iwata, R. Quiroz, C. Kuppe, O. Muller,
742 K. N. Blazakis, K. Yu, F. Tardieu, and S. C. Chapman. Modelling strategies for assessing and increasing the effectiveness of new phenotyping
743 techniques in plant breeding. *Plant Science*, 282:23–39, 2019. ISSN 18732259. doi: 10.1016/j.plantsci.2018.06.018.
- 744 J. G. Velazco, M. X. Rodríguez-Álvarez, M. P. Boer, D. R. Jordan, P. H. Eilers, M. Malosetti, and F. A. van Eeuwijk. Modelling spatial trends
745 in sorghum breeding field trials using a two-dimensional P-spline mixed model. *Theoretical and Applied Genetics*, 130:1375–1392, 2017.
746 ISSN 00405752. doi: 10.1007/s00122-017-2894-4.
- 747 A. P. Verbyla, J. De Faveri, J. D. Wilkie, and T. Lewis. Tensor Cubic Smoothing Splines in Designed Experiments Requiring Residual Modelling.
748 *Journal of Agricultural, Biological, and Environmental Statistics*, 23(4):478–508, 2018. ISSN 15372693. doi: 10.1007/s13253-018-0334-9.
- 749 S. Via, R. Gomulkiewicz, G. De Jong, S. M. Scheiner, C. D. Schlichting, and P. H. Van Tienderen. Adaptive phenotypic plasticity: consensus
750 and controversy. *Trends in Ecology & Evolution*, 10(5):212–217, 1995. ISSN 01695347. doi: 10.1016/S0169-5347(00)89061-8.
- 751 J. Vos, J. B. Evers, G. H. Buck-Sorlin, B. Andrieu, M. Chelle, and P. H. De Visser. Functional-structural plant modelling: A new versatile tool
752 in crop science. *Journal of Experimental Botany*, 61(8):2101–2115, 2010. ISSN 00220957. doi: 10.1093/jxb/erp345.
- 753 C. Welcker, W. Sadok, G. Dignat, M. Renault, S. Salvi, A. Charcosset, and F. Tardieu. A common genetic determinism for sensitivities to soil
754 water deficit and evaporative demand: Meta-analysis of quantitative trait loci and introgression lines of maize. *Plant Physiology*, 157(2):
755 718–729, 2011. ISSN 00320889. doi: 10.1104/pp.111.176479.
- 756 A. Whitehead. *Meta-Analysis of Controlled Clinical Trials*. John Wiley & Sons, Ltd, Chichester, 2002. ISBN 0-471-98370-5.
- 757 S. N. Wood. *Generalized additive models an introduction with R*. Chapman & Hall/CRC texts in statistical science. CRC Press/Taylor & Francis
758 Group, Boca Raton, second edi edition, 2017. ISBN 1498728332.