

1 **Title**

2 **Estimating fruit tree growth curves in breeding field using fragmented longitudinal data: An**  
3 **application to citrus hybrid seedlings**

4

5 **Running title:** Estimating growth curves in citrus seedlings

6

7 **Authors**

8 Soh Kimura<sup>1</sup>, Mai F. Minamikawa<sup>2</sup>, Keisuke Nonaka<sup>3</sup>, Tokurou Shimizu<sup>3, 4</sup>, Hiroyoshi Iwata<sup>1\*</sup>

9

10 <sup>1</sup>Laboratory of Biometry and Bioinformatics, Department of Agricultural and Environmental  
11 Biology, Graduate School of Agricultural and Life Sciences, The University of Tokyo, 1-1-1  
12 Yayoi, Bunkyo, Tokyo 113-8657, Japan

13 <sup>2</sup>Institute for Advanced Academic Research (IAAR), Chiba University, 1-33 Yayoi, Inage,  
14 Chiba, Chiba 263-8522, Japan

15 <sup>3</sup>Institute of Fruit Tree and Tea Science, NARO, Okitsu Nakacho, Shimizu, Shizuoka 424-  
16 0292, Japan

17 <sup>4</sup>Laboratory of Plant DNA Analysis, Kazusa DNA Research Institute, 2-6-7 Kazusa-kamatari,  
18 Kisarazu Chiba, 292-0818 Japan

19

20 \*To whom correspondence should be addressed:

21 Email: [hiroiwata@g.ecc.u-tokyo.ac.jp](mailto:hiroiwata@g.ecc.u-tokyo.ac.jp)

22 Tel: +81-03-5841-5069

23

24 The official email addresses of other authors:

25 Soh Kimura<sup>1</sup>: [kimura@ut-biomet.org](mailto:kimura@ut-biomet.org)

26 Mai F. Minamikawa<sup>2</sup>: [minamikawa@chiba-u.jp](mailto:minamikawa@chiba-u.jp)

27 Keisuke Nonaka<sup>3</sup>: [nonakak6@affrc.go.jp](mailto:nonakak6@affrc.go.jp)

28 Tokurou Shimizu<sup>3</sup>: [tshimizu@affrc.go.jp](mailto:tshimizu@affrc.go.jp)

29 **Abstract**

30 Vegetative and reproductive growth in fruit trees is interconnected, and analyzing this  
31 relationship can provide valuable insights into fruit quality. However, characterizing  
32 vegetative growth through growth models is challenging because of the difficulty in obtaining  
33 longitudinal data, given the slow growth rate. In breeding fields, in contrast, seedlings of  
34 different ages are planted, allowing for simultaneous measurements that yield a dataset  
35 resembling longitudinal data with missing values --termed “fragmented longitudinal data.”  
36 Because longitudinal data are obtained from a single measurement, they can potentially  
37 shorten the period required for growth curve estimation. Bayesian nonlinear models offer  
38 advantages in estimating curves from incomplete data. In this study, we generated fragmented  
39 longitudinal data using genome data with 45,929 markers from 624 citrus hybrid seedlings  
40 and applied a Bayesian nonlinear model to explore its potential. We also incorporated genomic  
41 information into the model to assess the impact of the estimation accuracy. Our simulations  
42 indicated that the Bayesian nonlinear model’s ability to interpolate missing values  
43 significantly improved the estimation performance. At best, the mean square error of the  
44 parameter characterizing the later growth stage was reduced by 84.3 mm<sup>2</sup>. Although the  
45 improvement from incorporating genomic information was modest, it still surpassed models  
46 that lacked genomic data. We also predicted the curves of untested individuals using the  
47 estimated parameters. Although the prediction accuracy of each parameter measured by the  
48 correlation coefficient was lower than 0.5, one parameter consistently showed a better  
49 accuracy. Further research is required to reveal the advantages of integrating genomic data  
50 for better predictions.

51 **Introduction**

52 Vegetative growth in fruit trees is closely linked to reproductive growth, acting as both a  
53 source and sink. Understanding this relationship has been a key focus for researchers because it  
54 directly influences fruit quality<sup>1</sup>. While physiological research has traditionally dominated the  
55 study of this relationship<sup>2,3</sup>, quantitative analysis of vegetative traits, such as trunk diameter and  
56 crown width, along with subsequent correlation studies, has highlighted the significant role of  
57 vegetative growth in determining fruit quality<sup>4-7</sup>. Additionally, these traits have been shown to  
58 correlate with the juvenile period and inbreeding depression<sup>8-10</sup>. With the advent of next-  
59 generation sequencing, these analyses are expected to expand to the genetic level. However, the  
60 complex growth patterns of trees and the dynamic nature of their phenotypes over time make the  
61 genetic analysis of these traits highly time-dependent. Therefore, it is essential to evaluate the  
62 evolving growth patterns and align genetic analyses accordingly.

63 The coexistence of trees of different ages in the same field presents specific challenges  
64 for genetic analyses, particularly for traits that are strongly influenced by time. In fruit tree  
65 breeding, orchards often contain a mix of candidates selected at various growth stages. For instance,  
66 in Japanese Citrus breeding, approximately 1,000 seedlings obtained from various crosses are  
67 grafted onto trifoliolate orange rootstocks annually and assessed for over 20 traits, including fruit  
68 quality and disease resistance<sup>11,12</sup>. Given that selection periods are typically under 8 years and new  
69 seedlings are planted annually, breeding orchards contain trees of diverse ages from various crosses.  
70 To analyze genetic variations in trees of various ages simultaneously, it is necessary to link data  
71 from different growth stages using a growth model and evaluate their growth patterns based on  
72 the parameters of that model.

73 Growth models such as the logistic and Gompertz models are effective in describing tree  
74 growth trajectories using a minimal number of time-independent parameters<sup>13</sup>. Their nonlinearity

75 allows them to outperform linear models, such as polynomial regression, in terms of requiring  
76 fewer interpretable parameters and providing more stable extrapolation<sup>14</sup>, which has led to  
77 extensive studies evaluating their applicability<sup>15</sup>. Growth models have been used to better  
78 understand growth patterns of fruits<sup>16</sup>. The parameters estimated from these models, which serve  
79 as time-independent indices of growth, are often used as target traits for genomic selection<sup>17,18</sup> or  
80 genome-wide association study (GWAS)<sup>19,20</sup>. Prior research has highlighted both the challenge of  
81 analysis based on time-dependent indices and the potential of the growth model, proposing QTL  
82 analysis based on stage-dependent indices derived from these models<sup>21</sup>. Unlike cross-sectional data,  
83 which can be obtained from a single measurement, longitudinal data, which are necessary for  
84 estimating growth curves, require repeated sampling, placing a significant burden on researchers.  
85 This burden is particularly high for trees, which have slow growth rates and require measurements  
86 over several years. Recently, unmanned aerial vehicles (UAVs), such as drones, have been used to  
87 measure phenotypic data and show promise in reducing this burden<sup>22,23</sup>. However, despite the  
88 anticipated reduction in measurement load owing to advancements in UAV technology, acquiring  
89 longitudinal data will continue to be challenging owing to the time required.

90         The use of an age-mixed field along with nonlinear mixed-effect models (NLMEMs) and  
91 Bayesian nonlinear models offers a potential solution for this challenge. In an age-mixed field,  
92 trees of different ages are planted together, allowing for simultaneous measurement of all  
93 individuals. This results in a dataset similar to longitudinal data, which we refer to as “fragmented  
94 longitudinal data.” Although these data contained substantial missing values for each tree  
95 (genotype), they were collected over a relatively short time span. If accurate growth curves can be  
96 estimated from these fragmented longitudinal data, the overall measurement time can be  
97 significantly reduced.

98         The nonlinear mixed effect model (NLMEM), which combines the characteristics of a

99 nonlinear model and a mixed-effect model, is frequently employed in the literature to handle  
100 incomplete longitudinal data<sup>24,25</sup>. NLMEM includes both fixed and random effects, where fixed  
101 parameters are shared among all subjects and random parameters are unique to each individual<sup>26</sup>.  
102 Because individual-specific values of random parameters are estimated based on a covariance  
103 structure derived from all individuals in the dataset, each individual tree (genotype) can “borrow”  
104 information from the others, allowing for the imputation of missing values<sup>13</sup>. Bayesian nonlinear  
105 mixed models extend NLMEMs by assuming that all model parameters, not only random effects,  
106 follow specific distributions<sup>27</sup>, which further aid in imputing missing values<sup>28</sup>. Therefore, we expect  
107 that NLMEM and Bayesian nonlinear modeling will play crucial roles in estimating growth curves  
108 from fragmented longitudinal data. Recent studies have incorporated genomic information into  
109 these models<sup>29,30</sup>. As growth parameters are genetically controlled to some extent, incorporating  
110 genomic information into these models may improve parameter estimation accuracy. Moreover,  
111 when genomic information is considered, the growth curves of untested individuals can be  
112 predicted based on genomic data in the context of genomic prediction and selection<sup>31</sup>.

113 In this study, we applied Bayesian nonlinear models to estimate growth curve parameters  
114 using fragmented longitudinal data collected from age-mixed fields without the need for time-  
115 consuming, labor-intensive repeated measurements. In addition, we evaluated whether  
116 incorporating genomic information into the models could improve the accuracy of parameter  
117 estimation. To validate the potential of the model and the role of genomic information, we used  
118 real genomic data from citrus breeding materials and simulated their vegetative growth curves. In  
119 these simulations, we compare the estimation accuracy of several methods and scenarios using  
120 fragmented longitudinal data. Furthermore, we assessed genomic prediction accuracy for untested  
121 individuals. Thus, we identified the potential of these methods to estimate growth curves based on  
122 short-term data for materials with varying age structures.

123

## 124 **Results**

### 125 **Generating fragmented longitudinal data from genome data**

126           The main purpose of this study was to estimate and predict growth curves from  
127 fragmented longitudinal data obtained from an age-mixed field. To accurately evaluate the  
128 estimation and prediction performance, rather than using raw data, fragmented longitudinal data  
129 were generated from artificially created longitudinal data based on the following measurement  
130 design: In this experiment, we assumed that measurements were conducted once a year over a 2-  
131 year period in a citrus breeding field. During the first measurement, trees grafted 2, 4, and 6 years  
132 ago, which were categorized as the youngest, middle, and oldest cohorts, respectively, were  
133 measured. The same trees were measured again in the second year, resulting in 6 years of  
134 longitudinal data with 2 data points per individual (Fig. 1A).

135           To investigate the influence of population structure on the estimation and prediction  
136 performance, two different longitudinal datasets were generated based on population composition:  
137 one from a single-family population and the other from a population of 11 families (Table 1). We  
138 refer to the former as the “single-family group” and the latter as the “multiple-family group.”  
139 During the generation of fragmented longitudinal data, the multiple-family group was further  
140 divided into two different groups; “multiple-family group I” and “multiple-family group II”  
141 depending on the allocation pattern of individuals. In single-family and multiple-family Group I,  
142 all individuals were randomly assigned to three cohorts. However, in actual breeding situations,  
143 the same cross is not repeated over multiple years, resulting in all individuals derived from the  
144 same cross concentrated within a specific year. To evaluate the estimation and prediction  
145 performance under these conditions, individuals in multiple-family group II were allocated to each  
146 cohort by family. Family allocations were determined randomly for each simulation trial (iteration).

147           The generated and prepared fragmented longitudinal data are shown in Fig. 1B. Despite  
148 the dynamic fluctuation caused by the measurement noise in each curve, the overall patterns of  
149 both the generated and fragmented longitudinal data follow a typical logistic growth curve. The  
150 phenotype data for each year formed a distribution that resembled a normal distribution owing to  
151 both genomic and residual variance. Although the position of the distribution shifted upward over  
152 time, the distributions of the middle and oldest cohorts overlapped significantly, indicating that  
153 the generated longitudinal data began to converge at an early stage.

154           The genetic structures of all individuals used to generate the longitudinal data were  
155 visualized using principal component analysis (Fig. 1C), with 11 families represented by 11  
156 distinct colors. Family 11, which formed a single-family group, was concentrated in the top-right  
157 section of the graph, whereas the 11 families that comprised the multiple-family group were widely  
158 dispersed across the graph, reflecting their greater genetic diversity.

159

#### 160 **Estimation accuracy of future/past growth using a nonlinear model**

161           To assess the difference in the estimation performance depending on the section of the  
162 estimated growth curve, two different scenarios were devised (Fig. 2A). In “Scenario 1,” the  
163 growth curve of the youngest cohort was estimated, while “Scenario 2” focuses on the oldest cohort,  
164 allowing the evaluation estimation performance for future. (younger), and past (older) growth. In  
165 Scenario 1, although the fragmented longitudinal data for the youngest cohort lacked information  
166 on later growth, the middle and oldest cohorts contained longitudinal data covering the later stages  
167 of growth. If the growth curve estimation for the youngest cohort can leverage information from  
168 the later growth of other cohorts, estimation performance is expected to improve. Similarly,  
169 estimating younger age growth using data from the oldest cohort is challenging. However, the  
170 estimation performance improves if the growth curve of the oldest cohort is jointly estimated with

171 those of the other two cohorts. To validate this assumption, three estimation methods were  
172 compared for each scenario: Method 1, in which the growth curve of the youngest/oldest cohort  
173 was estimated from only the youngest/oldest cohort data without genomic information; Method 2,  
174 in which the growth curve of the youngest/oldest cohort was estimated using data from all cohorts  
175 without genomic information; and Method 3, in which the growth curve of the youngest/oldest  
176 cohort was estimated using data from all cohorts with genomic information, which was expected  
177 to yield the highest estimation accuracy.

178 For the estimation, the fragmented longitudinal data were fitted to a logistic model with  
179 three parameters:  $A$ ,  $B$ , and  $C$ . Although parameter  $A$  characterizes the later growth, parameters  
180  $B$  and  $C$  determine the dynamics of the initial growth stage. The similarity between the estimated  
181 and true curves was evaluated at ages 1–7 years for each scenario using both the mean square error  
182 (MSE) and Pearson’s correlation coefficient (correlation coefficient), as shown in Fig. 2B.

183 In scenario 1, where the later growth of the youngest cohort was estimated, the MSEs of  
184 methods 2 and 3 decreased substantially compared with method 1, particularly at ages 4–7 years  
185 across all groups. While the medians of the MSE of method 3 at ages 4 to 7 years were consistently  
186 smaller than those of method 2, the degree of superiority varied depending on the group (Fig. 2B  
187 (a), (b), and (c)). Similarly, the estimation performance evaluated by the correlation coefficient at  
188 ages 4–7 years showed the worst results for Method 1 and the best for Method 3 (Supplementary  
189 Fig. 1 (a), (b), and (c)). It should be noted that the MSE for Method 1 reached as high as 8000  
190  $\text{mm}^2$ , but values beyond 100  $\text{mm}^2$  were not shown on the y-axis owing to its scaling.

191 In scenario 2, where the younger growth of the oldest cohort was estimated, the function  
192 used to calculate the initial value for the Bayesian nonlinear model function did not work properly;  
193 therefore, no growth curve was obtained for method 1. The MSE of both methods were high,  
194 particularly at ages two and three (Fig. 2B (d), (e), and (f)), and the correlation coefficient was



195 relatively low at ages one and two (Supplementary Fig. 1 (d), (e), and (f)), indicating imputation  
196 limitations. Similar to Scenario 1, the median MSE in Method 3 was consistently smaller than that  
197 in Method 2, and the median and interquartile range of the MSE in the single-family group were  
198 better than those in the other two groups.

199         Although the advantage of incorporating genomic information to improve estimation  
200 performance at ages 4–7 years was confirmed to be significant in scenario 1, the magnitude of this  
201 advantage was influenced by the population structure. To further investigate the relationship  
202 between population and the magnitude of improvement, we defined an index called “accuracy  
203 improvement.” This index was calculated for each simulation iteration as follows: (estimation  
204 accuracy in Method 3) minus (estimation accuracy in Method 2) for the correlation coefficient,  
205 and (estimation error in Method 2) minus (estimation error in Method 3) for the MSE. Because  
206 the estimation performance for Methods 2 and 3 was calculated simultaneously using the same  
207 created longitudinal data in each simulation iteration, taking the difference in each iteration makes  
208 it possible to evaluate the improvement in estimation performance by considering the genomic  
209 information more clearly. As a result, the accuracy improvement was found to be positive in most  
210 iterations and showed clear differences among groups in both the MSE and correlation coefficient.  
211 Multiple-family group I showed the best results, whereas the single-family group performed the  
212 worst (Fig. 2C).

213         The estimation performance of the latent parameters  $A$ ,  $B$ , and  $C$ , which control the  
214 overall growth behavior, was also investigated, focusing on multiple-family group II in scenario 1  
215 to identify the cause of accuracy improvement. Because the results of Method 1 sometimes  
216 included outliers that skewed the mean value, the estimation performance was summarized using  
217 the median of 1000 iterations instead of the mean, as shown in Table 2. In scenario 1, the MSEs  
218 of all parameters were confirmed to improve in methods 2 and 3, with a notable improvement in

219 methods  $A$  and  $B$  when genomic information was considered. However, in Scenario 2, although  
220 the MSE of parameters  $A$  and  $B$  were higher in Method 3 than in Method 2, the magnitude of the  
221 difference was not as pronounced as that in Scenario 1.

222

### 223 **Predicting untested progeny using the model trained on incomplete longitudinal data**

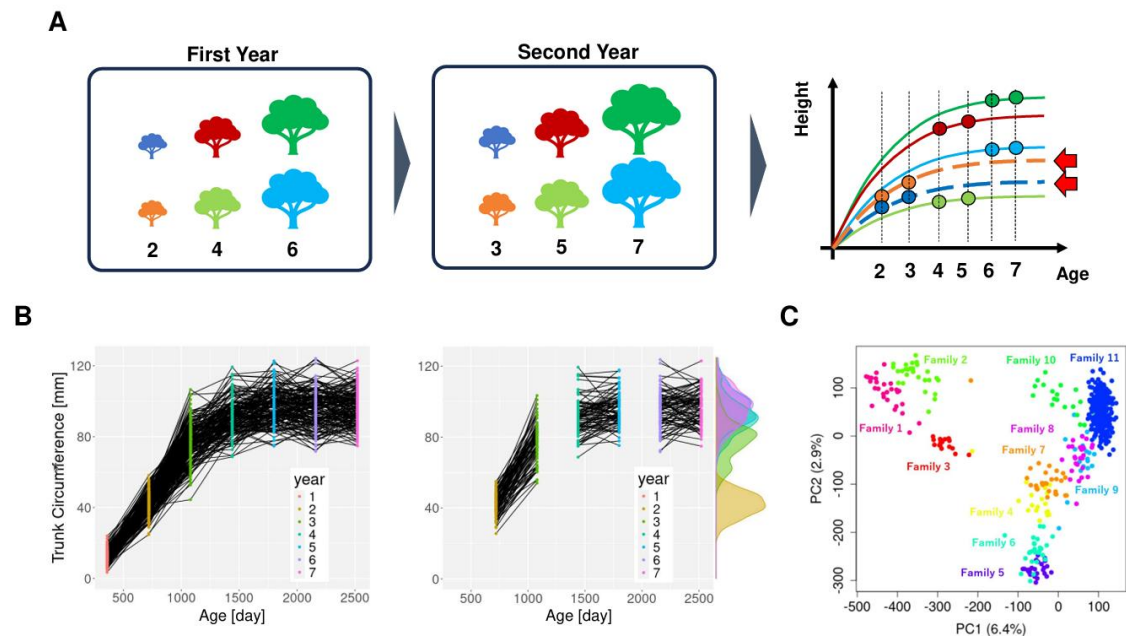
224 To date, estimation of the growth curve has been limited to individuals with phenotypic data.  
225 However, by incorporating the genomic information, it was possible to predict the growth curves  
226 of untested individuals (Fig. 3A). To assess the prediction performance of the model trained on  
227 fragmented longitudinal data, the prediction performances of latent parameters  $A$ ,  $B$ , and  $C$  were  
228 evaluated using cross-validation. The prediction accuracy (correlation coefficient) of parameter  $A$   
229 was higher than those of parameters  $B$  and  $C$  across all groups, with multiple-family group II  
230 showing the lowest accuracy for all parameters (Fig. 3B).

231

232 **Table 1. Number of individuals in each group**

Group	Family number										
	1	2	3	4	5	6	7	8	9	10	11
single- family group	0	0	0	0	0	0	0	0	0	0	266
multiple- family group	27	28	18	19	27	27	27	27	20	19	27

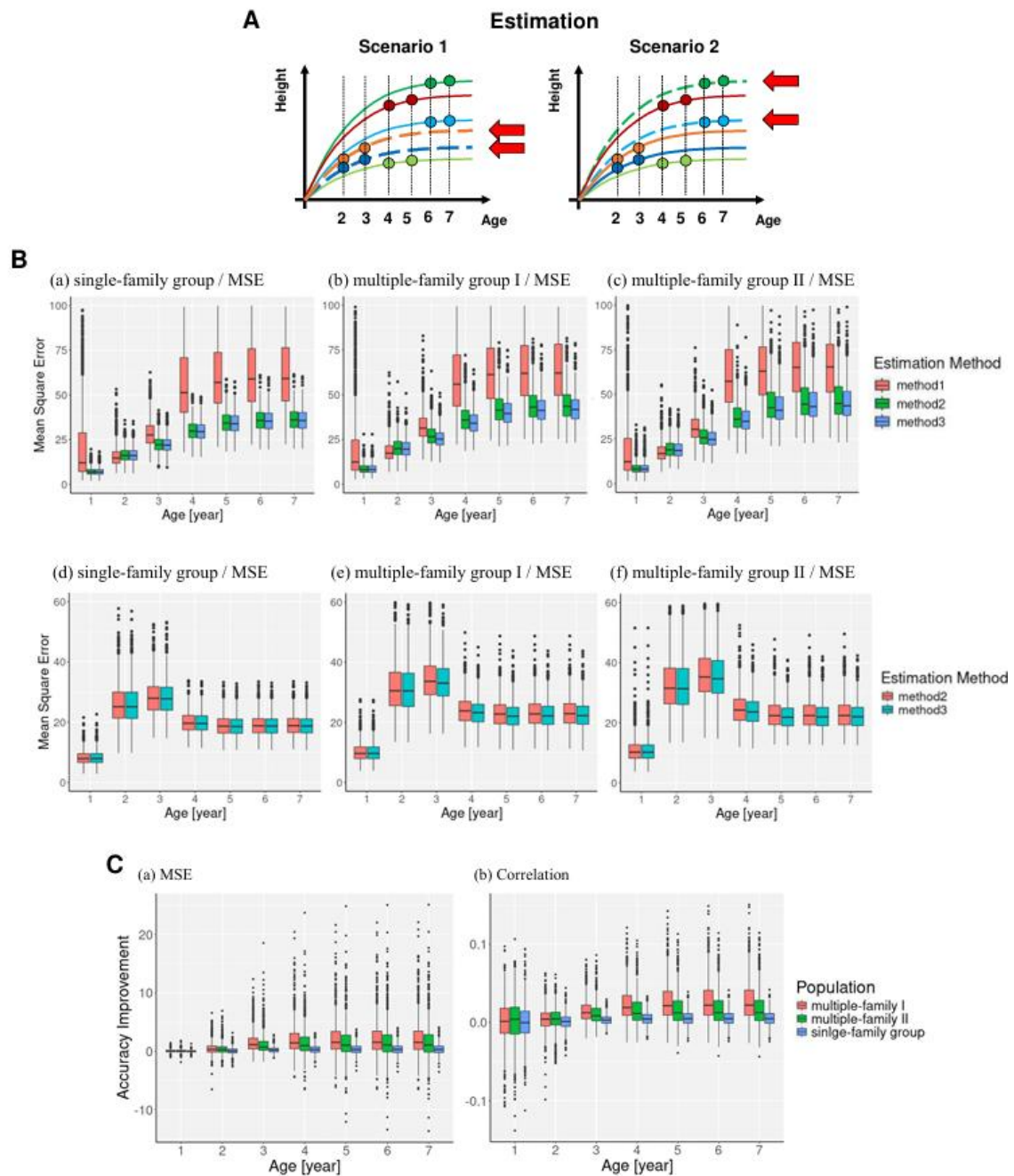
233



234

235 **Fig. 1. Measurement design and prepared fragmented longitudinal data**

236 (A) Detailed depiction of fragmented longitudinal data sampling and experimental scenario. The  
237 individuals with the same color are identical in the two left images, and the numbers below the  
238 trees represent the years after grafting. Measurements were conducted over 2 years for the same  
239 individuals, resulting in fragmented longitudinal data from years 2 to 7 within 2-year period. (B)  
240 Generated longitudinal data and the corresponding fragmented longitudinal data for multiple-  
241 family groups. The dots represent sampling data and are colored based on the year of sampling.  
242 (C) Principal component analysis of genomic data: all individuals were distinguished by 11 colors,  
243 corresponding to their respective families. The numbers on each axis shown in parentheses  
244 indicate the contribution rates.



245

246 **Fig. 2. Estimation scenario and estimation performance**

247 (A) In scenario 1, the later growth of the youngest cohort was estimated using three different

248 methods, and their estimation performance was compared. Similarly, in scenario 2, the old cohort

249 was targeted, and the estimation performance for younger growth was compared (B) The upper

250 three images (a), (b), and (c) show the results of scenario 1, whereas the lower three images (d),

251 (e), and (f) represent the results of scenario 2. Within each scenario, the left image corresponds  
 252 to the single-family group, middle to the multiple-family group I, and right to the multiple-family  
 253 group II. In each image, the estimation error of the growth curve was evaluated annually using the  
 254 three different methods, each represented by distinct colors. Boxplots are not shown when growth  
 255 curve could not be predicted. (C) Time variation of accuracy improvement: The x-axis represents  
 256 the age of the growth curve, and the y-axis represents the accuracy improvement at each age.  
 257 Estimation performance was evaluated using mean square error in the left image (a) and  
 258 correlation coefficient in the right image (b). Accuracy improvement was compared across groups,  
 259 with different colors representing each group.  
 260

261 **Table 2. Estimation performance of each parameter**

Scenario	Index	Method	Parameter		
			A	B	C
scenario 1					
Pearson Correlation Coefficient					
		method 1	0.61	0.55	0.33
		method 2	0.62	0.56	0.32
		method 3	0.65	0.57	0.33
Mean Square Error					
		method 1	127.9	2858.3	1733.0
		method 2	43.6	1406.1	610.6
		method 3	41.9	1399.7	611.3
scenario 2					

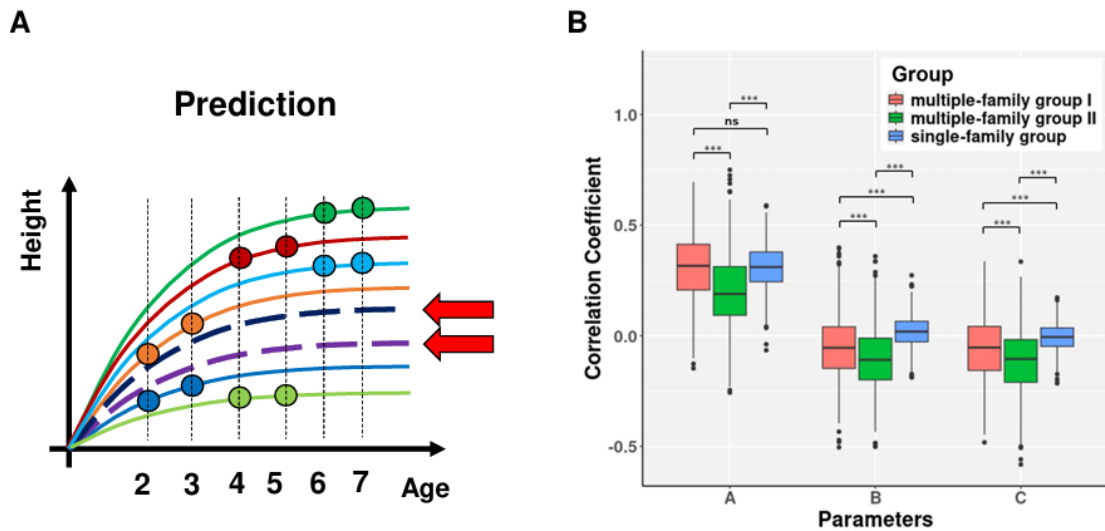
Pearson Correlation Coefficient

method 1	NA	NA	NA
method 2	0.65	0.00	0.00
method 3	0.65	0.02	0.01

Mean Square Error

method 1	NA	NA	NA
method 2	22.9	1722.7	657.8
method 3	22.3	1721.8	660.2

262

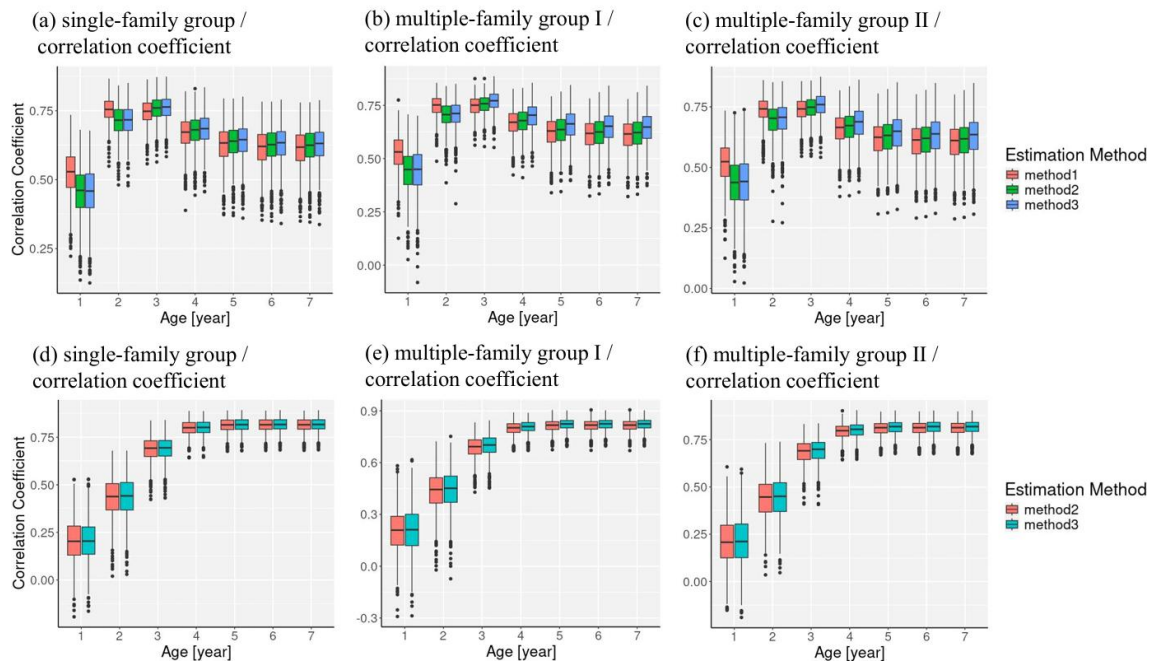


263

264 **Fig. 3. Prediction scenario and prediction performance**

265 (A) Individuals without phenotype were predicted using the model developed with fragmented  
 266 longitudinal data. (B) Prediction accuracy of all parameters in the logistic model. For each  
 267 parameter, the results for three different groups are represented by distinct colors. Statistical  
 268 significance is indicated by asterisk: \*\*\* for  $p \leq 0.001$  and ns for  $p > 0.05$ .

269



270

## 271 **Supplementary Fig. 1. Estimation accuracy in each group**

272 The upper three images (a), (b), and (c) show the results of Scenario 1, whereas the lower three  
273 images (d), (e), and (f) show the results of Scenario 2. Within each scenario, the left image  
274 corresponds to the single-family group, middle image to the multiple-family group I, and right  
275 image to the multiple-family group II. For each image, the estimation accuracy of the growth curve  
276 was evaluated annually using three methods, each represented by a distinct color. Box plots are  
277 not shown when the growth curve cannot be predicted.

278

## 279 **Discussion**

280 This study explored the feasibility of estimating the logistic growth curve from  
281 incompletely fragmented longitudinal data using a Bayesian nonlinear model to impute missing  
282 values. This approach aimed to shorten the measurement period and reduce the burden of  
283 constructing a growth curve. The fragmented longitudinal data comprised three cohorts (youngest,  
284 middle, and oldest), and the growth curve was estimated using three different methods and groups  
285 across the two scenarios. Our focus is on the following three points: 1. Extent to which the

286 estimation performance can be improved by considering data from other cohorts. 2. Whether the  
287 inclusion of genomic information enhances estimation performance. 3. How population structure  
288 impacts estimation performance. In addition, we evaluated the prediction performance of untested  
289 individuals based on the population structure.

290

### 291 **Improvements by leveraging data from individuals of other cohorts**

292 A comparison of Methods 1 and 2 addressed the first question. While Method 1 requires  
293 estimating the growth curve using only data from a single cohort, Method 2 enables simultaneous  
294 estimation using data from all cohorts, leveraging information from individuals in other cohorts.  
295 The comparison was conducted in two scenarios, and we confirmed an improvement in the  
296 estimation performance for later growth in Scenario 1.

297 Method 2 in Scenario 1 was designed to enhance the estimation performance for later  
298 growth, particularly for parameter A, which characterizes later growth. Interestingly, the MSEs  
299 improved not only for A, but also for B and C. Although parameter A primarily represents later  
300 growth, its influence extends to the initial growth phase, where B and C play critical roles.  
301 Therefore, improving the estimation performance of parameter A reduces uncertainty in the initial  
302 growth phase, leading to better estimates of parameters B and C. This improvement is attributed  
303 to the properties of the nonlinear Bayesian model. The Bayesian nonlinear model assumes that the  
304 parameters of all individuals follow the same distribution, thus preventing the parameters from  
305 deviating excessively from the overall behavior. In Method 1, parameter A for the youngest cohort  
306 became highly unstable due to the absence of later growth data, whereas the middle and oldest  
307 cohort parameters were estimated more reliably. This assumption in the Bayesian nonlinear model  
308 helps bind all the parameters together, thereby stabilizing the estimation performance for the  
309 youngest cohort.



310           Such a function to impute missing values can be expected not only in the Bayesian  
311 nonlinear model but also in the NLMEM, where random effects are assumed to follow the same  
312 distribution. The application of the NLMEM to incomplete data with missing values has been  
313 documented in various studies<sup>24,25</sup>. Additionally, some studies such as ours have attempted to  
314 estimate later growth from initial growth data<sup>32</sup>. In a shift from plant-related research, a functional  
315 linear mixed model (FLMM) was employed in another study to estimate adult growth from  
316 children's data, using previously obtained completed longitudinal data<sup>33</sup>. However, our study aims  
317 to shorten the measurement period, and is distinct in that it applies specifically structured missing  
318 data, known as fragmented longitudinal data.

319           Meanwhile, in Method 2 of Scenario 2, in which the growth curve of the older cohort was  
320 estimated using data from all cohorts, the estimation performance of parameters B and C, which  
321 characterize initial growth, dropped below 0.5. This decline can be attributed to the exclusion of  
322 the first-year data and uniform measurement intervals. We hypothesized that the targeted growth  
323 curve would represent the trunk circumference of the scion grafted onto the rootstock. Since  
324 measuring trunk circumference in the first year is difficult due to the thinness of the trunk, the  
325 first-year data were excluded from the fragmented longitudinal data, resulting in a lack of  
326 information on the initial growth phase. Additionally, assuming annual measurements over 6 years  
327 with constant intervals led to unbalanced sampling, as reflected in the distribution shift shown in  
328 Fig. 1B, where the initial growth was not sufficiently captured. The importance of the  
329 measurement interval has been well documented in the prior research<sup>34</sup>. Although our study  
330 indicates that estimating initial growth from fragmented longitudinal data is challenging, the  
331 performance can be improved by adjusting the measurement intervals.

332

333 **Improvement by genomic information**

334 To explore the potential of genomic information for imputing missing values, Method 3  
335 was designed by incorporating a genomic relationship matrix (GRM) into Method 2: Specifically,  
336 the improvement in the estimation performance by considering the GRM within the Bayesian  
337 nonlinear model was examined across three groups: the single-family group, multiple-family group  
338 I, and multiple-family group II. The superiority of Method 3 over Method 2 was confirmed,  
339 although the magnitude of superiority varied by group.

340 Generally, the longitudinal data at specific sampling points exhibit correlations with  
341 certain parameters. For example, late-stage longitudinal data are strongly correlated with  
342 parameter A, whereas initial-stage longitudinal data are more closely correlated with parameters  
343 B and C. Traditionally, parameter estimation has relied on these correlations. However, in this  
344 study, a large amount of data was missing, making it difficult to rely solely on correlation-based  
345 estimations. Nevertheless, when genomic information is considered, the parameters in the  
346 Bayesian nonlinear model are linked not only to phenotypes but also to genomic relationships;  
347 individuals who are genetically closer tend to have similar parameter values, and those who are  
348 genetically distant show more variation. In such cases, parameters can be estimated based on  
349 genomic relationships, which we believe will contribute to improvements in estimation  
350 performance. The algorithm for incorporating GRM into nonlinear models has been developed  
351 recently<sup>29,30,35,36</sup>, and to the best of our knowledge, this is the first study to demonstrate the potential  
352 of genomic information for imputing missing values within a nonlinear framework.

353

### 354 **Influence of population structure**

355 While the inclusion of genetic information has been confirmed to improve the estimation  
356 of growth curves from fragmented longitudinal data, the magnitude of this improvement, as  
357 evaluated by accuracy, depends heavily on the group structure. Multiple-family group I showed a

358 better improvement in estimation accuracy, whereas the single-family group yielded the poorest  
359 results. The key differences between these groups lie in the populations used to generate the  
360 longitudinal data and the strategy for allocating individuals to the three cohorts. Because both the  
361 single-family group and multiple-family group I used the same allocation strategies, their  
362 estimation performances can be compared based solely on the population structure. The genetic  
363 structure of the population used in this study, as visualized by PCA, clearly showed a wider range  
364 of genetic variation in the multiple-family group than in the single-family group. As we are  
365 evaluating the improvement guided by incorporating genomic information, a broader range of  
366 genetic variation suggests great potential for improvement, which likely contributes to the higher  
367 accuracy in multiple-family group I.

368         Meanwhile, multiple-family Groups I and II used the same population to generate  
369 longitudinal data but applied different allocation strategies, allowing for the comparison of  
370 estimation performance based solely on allocation patterns. Although the allocation strategy used  
371 in multiple-family group II was more realistic, its estimation performance was lower than that of  
372 multiple-family group I. When individuals are randomly assigned, as in the multiple-family group  
373 I, genetically related individuals, such as those from the same family, are distributed more evenly  
374 across the three cohorts, making it easier to estimate the growth curve. In contrast, when  
375 individuals from the same family are concentrated in a specific cohort, such as in the multiple-  
376 family group II, the phenotype data across cohorts seem less related, making it more difficult to  
377 estimate phenotypic values across cohorts. In genomic selection, the prediction performance is  
378 strongly influenced by the genomic relationship between the training and test populations, leading  
379 to numerous studies on the impact of the allocation strategy<sup>37</sup>. However, the experimental setup  
380 described in this section differs in that it focuses on longitudinal genetic relationships rather than  
381 on the relationship between the training and test populations. Another key distinction is that this

382 section emphasizes estimation rather than prediction. The prediction performance considering the  
383 allocation of the training and test populations is discussed in the subsequent section.

384

### 385 **Prediction of untested individuals**

386 Although our primary focus has been on the estimation performance of tested individuals,  
387 the incorporation of the GRM into the Bayesian nonlinear model also enables the direct prediction  
388 of untested individuals. In this study, we attempted to build a model using fragmented longitudinal  
389 data and evaluated the prediction performance for each group. However, the prediction  
390 performance was low, particularly for parameters B and C; even parameter A did not exceed 0.4,  
391 indicating the challenge of prediction from fragmented longitudinal data.

392 Among the three groups, the single-family group and multiple-family group I exhibited  
393 the best prediction performance for parameter A, whereas multiple-family group II exhibited the  
394 poorest performance across all parameters. The relationship between the prediction performance  
395 and population structure has been actively studied in the context of genomic selection. For  
396 example, a previous study categorized a wheat population based on training and test population  
397 structures and found that prediction performance declined when the family of the test population  
398 was excluded from the training population, attributing the decline to the localization of family  
399 specific genes<sup>37</sup>. Although the prediction performances of B and C were insufficient, they may be  
400 improved by optimizing the sampling time, as mentioned earlier, which warrants further  
401 investigation.

402

## 403 **Materials and Methods**

### 404 **Plant material**

405 In this study, we used 505 individuals with different genotypes from 11 families

406 generated by the National Agriculture and Food Research Organization, Institute of Fruit Tree  
407 and Tea Science (Okitsu, Shizuoka, Japan). The parents of these families had diverse backgrounds,  
408 including tangerines, lemons, and pomelos, and the families were numbered from 1 to 11. Family  
409 11 comprised 266 progeny, whereas the remaining 10 families had an average of 20 individuals.

410

411

#### 412 **SNP genotyping data**

413 We conducted GRAS-Di<sup>38</sup> analysis using Illumina HiSeq 4000, generating 150 bp paired-  
414 end reads according to the manufacturer's protocol, with an average of approximately 8 million  
415 reads per sample. We then mapped the trimmed clean reads to the clementine haploid reference  
416 genome v1.0  
417 (<https://data.jgi.doe.gov/refinedownload/phytozome?organism=Cclementina&expanded=Phytozome-182>)  
418 using GSNAP (version 2021-12-17)<sup>39</sup> with options `-k = 15` and `-max-mismatches=0.125`, resulting in BAM files for individual samples. After sorting the BAM files in  
419 coordinate order using Picard tools, the secondary mapped reads were marked and index files were  
420 generated. We called polymorphic sites using the GATK4 (4.5.5.0) HaplotypeCaller and created  
421 GVCF files for each sample. Next, we merged the GVCF files using GATK4 CombineGVCFs, and  
422 converted the merged GVCF file into a VCF file using GATK4 GenotypeGVCFs. The VCF file  
423 was filtered with thresholds of `meanDP >= 20` and `meanGQ >= 20` using `vcftools`<sup>40</sup>. Imputation  
424 was performed using Beagle 5.2<sup>41</sup>, selecting data based on a minor allele frequency (MAF) > 1%.  
425 Using Picard, we verified the Mendelian violations in 387 trios and excluded loci with  
426 inconsistencies in more than two trios. Finally, the linkage disequilibrium was calculated, and  
427 SNPs were selected based on a linkage disequilibrium threshold of less than 0.95 using "LD.thin"  
428 function of the R package RAINBOWR<sup>42</sup>, resulting in a total of 45,929 makers.

430

### 431 **Population structure analysis**

432 Principal component analysis was performed using the SNP genotyping data of 505  
433 individuals with the “prcomp” function of the R package stats.

434

### 435 **Longitudinal data generation**

436 The longitudinal data were generated using a growth model. Among various growth  
437 models, such as the Gompertz, von Bertalanffy, and Richards models, the logistic model was best  
438 suited to obtain longitudinal data of citrus trunk circumference of citrus breeding populations 2 to  
439 7 years after grafting onto 3-year-old trifoliolate orange rootstock (data not shown); therefore, it was  
440 used in our study. The logistic model was parameterized as follows:

$$441 \quad f(A, B, C, t) = \frac{A}{1 + \exp\{-(t - B)/C\}}$$

442 where,  $A$  is the asymptotic parameter, represents the maximum value of the growth curve,  $B$  is  
443 the inflection time, and  $1/C$  is the growth rate at the inflection point. We assumed that these  
444 parameters were unique to each individual, and the unique parameters  $A_i$ ,  $B_i$ , and  $C_i$  for the  $i$ -  
445 th individual were determined using the following equations:

$$446 \quad \varphi_{i,j} = \mu_j + u_{i,j}$$

$$447 \quad \mathbf{u}_j = (u_{1,j}, u_{2,j}, \dots, u_{n,j})^T \sim MVN(\mathbf{0}, K\sigma_j^2)$$

448 where  $\varphi_{i,j}$  is the  $j$ -th parameters of  $i$ -th individual,  $\mu_j$  is the mean of the  $j$ -th parameters, and  
449  $u_{i,j}$  is the deviation from the global mean, characterizing each individual.  $\sigma_j^2$  is the variance and  
450  $K$  is the genomic relationship matrix, calculated for both the single-family group and the multiple-  
451 family group using the “calcGRM” function of the RAINBOWR package. In the process of data  
452 generation, the mean  $\mu_j$  and the variance  $\sigma_j^2$  were first fixed, and then the individual-specific  
453 parameter  $\varphi_{i,j}$  was randomly generated. These parameters were applied to a logistic model to

454 determine individual-specific growth curves. The values of  $\mu_j$  and  $\sigma_j^2$  were determined to mimic  
455 the actually measured citrus trunk circumference growth data as follows:

$$456 \quad (\mu_1, \mu_2, \mu_3)^T = (100, 800, 200)^T$$

$$457 \quad (\sigma_1^2, \sigma_2^2, \sigma_3^2)^T = (100, 2500, 900)^T$$

458 Data for each year from 2 to 7 years were extracted from the generated individual-specific growth  
459 curves, and longitudinal data were completed by adding measurement noise. To ensure that  
460 heritability (i.e., the proportion of variance explained by genetic factors) remained constant over  
461 time, the variance of the longitudinal data was first calculated at each sampling point, and then the  
462 magnitude of variance was adjusted to maintain a heritability of 0.5.

463

#### 464 **Fragmented longitudinal data generation**

465 Fragmented longitudinal data were obtained from the artificially generated longitudinal  
466 data. Specifically, the 4th to 7th years' data were replaced with missing values for the youngest  
467 cohort, while the 2nd, 3rd and 6th, and 7th years' data were replaced with missing values for the  
468 middle cohort. Similarly, the 2nd to 7th years' data were replaced with missing values for the oldest  
469 cohort.

470

#### 471 **Estimation of latent variable**

472 To estimate the three parameters,  $A_i$ ,  $B_i$ , and  $C_i$ , of the logistic model, we applied a  
473 Bayesian nonlinear model. Individual growth curves were estimated by incorporating these  
474 parameters into the model functions. The Bayesian nonlinear model is defined as follows:

475 Let  $f$  represent a nonlinear growth model,  $\mathbf{D}$  the model input,  $\Phi$  the model  
476 parameters,  $\mathbf{Y}$  the observations, and  $i$  the genotype index. The model can be expressed as  
477 follows:

478 
$$\mathbf{Y}_i = f(\mathbf{D}_i, \Phi_i) + \boldsymbol{\varepsilon}_i$$

479 where  $\boldsymbol{\varepsilon}_i$  is the error term, assumed to follow a multiple normal distribution:

480 
$$\boldsymbol{\varepsilon}_i \sim MVN(\mathbf{0}, I\sigma_\varepsilon^2)$$

481 The prior distribution of the residual variances  $\sigma_\varepsilon^2$  is the Jeffreys' scale-invariant prior:

482 
$$p(\sigma_\varepsilon^2) \propto \frac{1}{\sigma_\varepsilon^2}$$

483 The  $j$ -th parameter of  $\Phi_i$  is regressed on genome-wide marker genotypes as:

484 
$$\varphi_{i,j} = g_j(G_i, \Theta_j, \Psi_j) + e_{i,j}$$

485 where  $G_i$  is the marker genotypes,  $g_j$  is the whole-genome regression function,  $\Theta_j$  is the

486 regression parameters of  $f_j$ ,  $\Psi_j$  is the hyperparameters of  $f_j$ , and:

487 
$$e_{i,j} \sim N(0, 1/\tau_{0,j}^2)$$

488 The prior for  $\tau_{0,j}^2$  follows

489 
$$p(\tau_{0,j}^2) \propto \frac{1}{\tau_{0,j}^2}$$

490 In this study, we used a GBLUP model for whole-genome regression, defined as

491 
$$\mathbf{g} = \mathbf{1}\mu + \mathbf{u}$$

492 with priors:

493 
$$p(\mu) \propto \text{constant}$$

494 
$$\mathbf{u} \sim MVN(\mathbf{0}, K\sigma_u^2)$$

495 where  $K$  is the additive genomic relationship matrix. The prior for  $\sigma_u^2$  is:

496 
$$\sigma_u^2 \sim \text{Scaled inverse schi square}(v_u, S_u^2)$$

497 The parameters to estimate  $\Theta$  are  $\mu, \mathbf{u}, \tau_0^2$ , and  $\sigma_u^2$ , and the hyperparameters  $\Psi$  are  $v_u$  and  $S_u^2$ .

498 For Methods 1 and 2, in which genomic information was not used, the additive genomic

499 relationship matrix  $K$  was replaced by the identity matrix  $I$ .

500 Estimation was performed using the "GenomeBasedModel" function of the



501 GenomeBasedModel package in R, with the initial value set by the “drm” function of the drc  
502 package<sup>43</sup>. In the initial value calculation, the longitudinal dataset of all individuals was treated as  
503 a single individually derived dataset. Iterations were run 1000 times for each generated dataset,  
504 and the estimation was evaluated using both Pearson’s correlation coefficient and the mean square  
505 error.

506

### 507 **Prediction of untested individuals**

508 When genomic information is incorporated into the Bayesian nonlinear model, the  
509 estimated parameters follow a multivariate normal distribution with a variance-covariance matrix  
510 proportional to the GRM. In this case, the logistic model parameters of untested (predicted)  
511 individuals,  $\boldsymbol{\varphi}_{j_{pred}} = (\varphi_{1.j_{pred}}, \dots, \varphi_{n.j_{pred}})^T$ , can be predicted based on the genomic relationship  
512 matrix and the estimated parameters of tested (observed) individuals,  $\boldsymbol{\varphi}_{j_{obs}} =$   
513  $(\varphi_{1.j_{obs}}, \dots, \varphi_{n.j_{obs}})^T$ , using the following equation<sup>44</sup>:

$$514 \quad \boldsymbol{\varphi}_{j_{pred}} = \mathbf{1}\mu_{j_{pred}} + K_{pred\ pred}^{-1}K_{pred\ obs}(\boldsymbol{\varphi}_{j_{obs}} - \mathbf{1}\mu_{j_{obs}})$$

515 where  $K_{pred\ pred}$  is the additive genetic relationship matrix corresponding to the predicted  
516 individuals, and  $K_{pred\ obs}$  is that of the predicted and observed individuals.  $\mu_{j_{pred}}$  and  $\mu_{j_{obs}}$  are  
517 the means of the predicted and estimated parameters (i.e., estimated parameters of an observed  
518 individual) respectively. In predicting  $\boldsymbol{\varphi}_{j_{pred}}$ ,  $\mu_{j_{pred}}$  is replaced by  $\mu_{j_{obs}}$ . A growth curve was  
519 then constructed using these parameters.

520 To evaluate the prediction performance of the model constructed using the fragmented  
521 longitudinal data, we conducted the following experiment: For the single-family and multiple-  
522 family group I, one eleventh of the individuals were used as the training population, and the rest

523 were used as the training population, which was further divided into three cohorts. Cross-  
524 validation was performed by rotating the test population 11 times and the prediction performance  
525 was calculated. This process was repeated 1000 times using different longitudinal data. In the case  
526 of multiple-family group II, individuals were allocated by family, and each family was selected as  
527 the test population for each rotation. It should be noted that because all individuals across the  
528 three cohorts were used for growth curve prediction, there was no distinction between Methods 1  
529 and 2.

530

### 531 **Significance test of the difference in prediction performance**

532 Because the variance in prediction performance across all simulation iterations differed  
533 by group, a statistical significance test for prediction performance was conducted using the Steel-  
534 Dwass asymptotic test. This was implemented with the “pSDCFlig” function in “NSM3” package  
535 in R.

536

### 537 **Acknowledgments**

538 We are grateful to all members of the National Agriculture and Food Research Organization  
539 Institute of Fruit Tree and Tea Science for maintaining the Citrus trees, as well as Kosuke  
540 Hamazaki and Kengo Sakurai for sharing their expertise in the simulation analysis. This research  
541 is supported by a grant from MAFF commissioned project study on “Smart breeding  
542 technologies to Accelerate the development of new varieties toward achieving “Strategy for  
543 Sustainable Food Systems, MIDORI”” and Japan Science and Technology Agency – OPERA  
544 (Program on Open Innovation Platform with Enterprises, Research Institute and Academia)  
545 Grant Number JPMJOP1851.

546

### 547 **Data Availability Statements**

548 Data available on request.

549

### 550 **Conflict of interests**

551 The authors declare that they have no conflict of interest.

552

553 **Competing financial interests**

554 The authors declare no competing financial interests.

555

556 **Author Contributions**

557 S. K. and H. I. conceived and designed the study. K.N. designed the study. T.S. extracted DNA  
558 and performed SNP genotyping. S.K., M.F.M., T.S., H.I., and K.N. conducted phenotyping. S. K.  
559 performed the simulations. H.I. provided technical help for statistical analysis. S. K. and T. S.  
560 drafted the manuscript. All the authors have read and approved the manuscript.

561

562 **References**

563

564 1 Yamaki S. Metabolism and accumulation of sugars translocated to fruit and their regulation.  
565 *J Japan Soc Hortic Sci* 2010; **79**: 1–15.

566 2 Koch KE. Translocation of photosynthetic products from source leaves to aligned juice  
567 segments in Citrus fruit. *HortScience* 1984; **19**: 260–261.

568 3 Koch KE. The path of photosynthate translocation into citrus fruit. *Plant Cell Environ* 1984;  
569 **7**: 647–653.

570 4 Rosati A, Paoletti A, Pannelli G, Famiani F. Growth Is Inversely Correlated with Yield  
571 Efficiency across Cultivars in Young Olive (*Olea europaea* L.) Trees. *HortScience* 2017; **52**:  
572 1525–1529.

573 5 Kumar D, Srivastava KK, Singh SR. Correlation of trunk cross sectional area with fruit yield,  
574 quality and leaf nutrient status in plum under North West Himalayan region of India. *J Hortic*  
575 *Sci* 2019; **14**: 26–32.

576 6 Fagherazzi AF *et al.* Initial crown diameter influences on the fruit yield and quality of  
577 strawberry Pircinque. *Agronomy (Basel)* 2021; **11**: 184.

578 7 Wang Zhijun and Lan PASF. Correlation Research on the Structure of the Apple Tree Vigor  
579 and Its Fruit Quality. In: Nakamatsu Kazumi and Kountchev RAAAAE-BNAHB (ed). *New*  
580 *Developments of IT, IoT and ICT Applied to Agriculture*. Springer Singapore: Singapore,  
581 2021, pp 55–63.

582 8 Watson A *et al.* Speed breeding is a powerful tool to accelerate crop research and breeding.  
583 *Nature Plants* 2018; **4**: 23–29.

584 9 Mitani N, Matsumoto R, Yoshioka T, Kuniga T. Citrus hybrid seedlings reduce initial time to

- 585 flower when grafted onto shiikuwasha rootstock. *Sci Hortic (Amsterdam)* 2008; **116**: 452–455.
- 586 10 DeBuse CJ, Shaw DV, DeJong TM. Response to inbreeding of early seedling growth and fruit  
587 traits in a *Prunus domestica* l. Breeding population. *Acta Hortic* 2013; : 87–95.
- 588 11 Omura M, Shimada T. Citrus breeding, genetics and genomics in Japan. *Breed Sci* 2016; **66**:  
589 3–17.
- 590 12 Shimizu T. Citrus breeding 2.0: A novel approach integrating deciphered parentage and  
591 genomics-assisted selection. *Jpn Agric Res Q* 2019; **53**: 81–85.
- 592 13 Zimmerman DL *et al.* Parametric modelling of growth curve data: An overview. *Test (Madr)*  
593 2001; **10**: 1–73.
- 594 14 Carey VJ, Wang Y-G. Mixed-effects models in S and S-plus. *J Am Stat Assoc* 2001; **96**: 1135–  
595 1136.
- 596 15 Salas-Eljatib C, Mehtätalo L, Gregoire TG, Soto DP, Vargas-Gaete R. Growth equations in  
597 forest research: Mathematical basis and model similarities. *Curr For Rep* 2021; **7**: 230–244.
- 598 16 Zadavec P, Veberic R, Stampar F, Schmitzer V, Eler K. Fruit Growth Patterns of Four Apple  
599 Cultivars Using Nonlinear Growth Models. .
- 600 17 Pong-Wong R, Hadjipavlou G. A two-step approach combining the Gompertz growth model  
601 with genomic selection for longitudinal data. *BMC Proc* 2010; **4 Suppl 1**: S4.
- 602 18 Toda Y *et al.* Genomic prediction of green fraction dynamics in soybean using unmanned  
603 aerial vehicles observations. *Front Plant Sci* 2022; **13**: 828864.
- 604 19 Duan X *et al.* Genome-wide association analysis of growth curve parameters in Chinese  
605 Simmental beef cattle. *Animals (Basel)* 2021; **11**: 192.
- 606 20 Xia H *et al.* Genome-wide association study of multiyear dynamic growth traits in hybrid  
607 *Liriodendron* identifies robust genetic loci associated with growth trajectories. *Plant J* 2023;  
608 **115**: 1544–1563.
- 609 21 Yin X, Kropff MJ, Stam P. The role of ecophysiological models in QTL analysis: the example  
610 of specific leaf area in barley. *Heredity (Edinb)* 1999; **82**: 415–421.
- 611 22 Chang A, Yeom J, Jung J, Landivar J. Comparison of canopy shape and vegetation indices of  
612 citrus trees derived from UAV multispectral images for characterization of citrus greening  
613 disease. *Remote Sens (Basel)* 2020; **12**: 4122.

- 614 23 Moe KT, Owari T, Furuya N, Hiroshima T, Morimoto J. Application of UAV photogrammetry  
615 with LiDAR data to facilitate the estimation of tree locations and DBH values for high-value  
616 timber species in northern Japanese mixed-wood forests. *Remote Sens (Basel)* 2020; **12**: 2865.
- 617 24 Huang S, Meng SX, Yang Y. Assessing the goodness of fit of forest models estimated by  
618 nonlinear mixed-model methods. *Can J For Res* 2009; **39**: 2418–2436.
- 619 25 Yang Y, Huang S. Comparison of different methods for fitting nonlinear mixed forest models  
620 and for making predictions. *Can J For Res* 2011; **41**: 1671–1686.
- 621 26 Lindstrom MJ, Bates DM. Nonlinear Mixed Effects Models for Repeated Measures Data.  
622 1990.
- 623 27 Lee SY. Bayesian nonlinear models for repeated measurement data: An overview,  
624 implementation, and applications. *Mathematics* 2022; **10**: 898.
- 625 28 Ibrahim JG, Molenberghs G. Missing data methods in longitudinal studies: a review. *Test*  
626 (*Madr*) 2009; **18**: 1–43.
- 627 29 Onogi A *et al.* Toward integration of genomic selection with crop modelling: the development  
628 of an integrated approach to predicting rice heading dates. *Züchter Genet Breed Res* 2016;  
629 **129**: 805–817.
- 630 30 Delattre M, Toda Y, Tressou J, Iwata H. Modeling soybean growth: A mixed model approach.  
631 bioRxiv. 2024; : 2023.06.13.544713.
- 632 31 Gianola D, Weigel KA, Krämer N, Stella A, Schön C-C. Enhancing genome-enabled  
633 prediction by bagging genomic BLUP. *PLoS One* 2014; **9**: e91693.
- 634 32 Hall DB, Bailey RL. Modeling and Prediction of Forest Growth Variables Based on Multilevel  
635 Nonlinear Mixed Models. *For Sci* 2001; **47**: 311–321.
- 636 33 Ran M, Yang Y. Optimal estimation of large functional and longitudinal data by using  
637 Functional Linear Mixed Model. *Mathematics* 2022; **10**: 4322.
- 638 34 Timmons AC, Preacher KJ. The Importance of Temporal Design: How Do Measurement  
639 Intervals Affect the Accuracy and Efficiency of Parameter Estimates in Longitudinal  
640 Research? *Multivariate Behav Res* 2015; **50**: 41–55.
- 641 35 Campbell MT, Grondin A, Walia H, Morota G. Leveraging genome-enabled growth models  
642 to study shoot growth responses to water deficit in rice. *J Exp Bot* 2020; **71**: 5669–5679.

- 643 36 Yu H, van Milgen J, Knol E, Fernando R, Dekkers J. 307. A Bayesian hierarchical model to  
644 integrate a mechanistic growth model in genomic prediction. In: *Proceedings of 12th World*  
645 *Congress on Genetics Applied to Livestock Production (WCGALP)*. Wageningen Academic  
646 Publishers, 2022, pp 1290–1293.
- 647 37 Norman A, Taylor J, Edwards J, Kuchel H. Optimising genomic selection in wheat: Effect of  
648 marker density, population size and population structure on prediction accuracy. *G3*  
649 (*Bethesda*) 2018; **8**: 2889–2899.
- 650 38 Enoki H, Takeuchi Y, Suzuki K. New genotyping technology, GRAS-Di, using next generation  
651 sequencer. In: *Proceedings of the Plant and Animal Genome Conference XXVI*. 2018.
- 652 39 Wu TD, Reeder J, Lawrence M, Becker G, Brauer MJ. GMAP and GSNAP for genomic  
653 sequence alignment: Enhancements to speed, accuracy, and functionality. *Methods Mol Biol*  
654 2016; **1418**: 283–334.
- 655 40 Danecek P *et al*. The variant call format and VCFtools. *Bioinformatics* 2011; **27**: 2156–2158.
- 656 41 Browning BL, Browning SR. A unified approach to genotype imputation and haplotype-phase  
657 inference for large data sets of trios and unrelated individuals. *Am J Hum Genet* 2009; **84**:  
658 210–223.
- 659 42 Hamazaki K, Iwata H. Rainbow: Haplotype-based genome-wide association study using a  
660 novel SNP-set method. *PLoS Comput Biol* 2020; **16**. doi:10.1371/journal.pcbi.1007663.
- 661 43 Ritz C, Baty F, Streibig JC, Gerhard D. Dose-response analysis using R. *PLoS One* 2015; **10**.  
662 doi:10.1371/journal.pone.0146021.
- 663 44 Hallander J, Waldmann P, Wang C, Sillanpää MJ. Bayesian inference of genetic parameters  
664 based on conditional decompositions of multivariate normal distributions. *Genetics* 2010;  
665 **185**: 645–654.
- 666