1 # Comparative Analysis of Genetically-Modified Crops: Conditional

2 # Equivalence Criteria

3

4 ## Short Title: Conditional Equivalence Criteria for Genetically-Modified Crops

5

6 Changjian Jiang[1*], Chen Meng[1], Adam W. Schapaugh[1], Huizhe Jin[1]

7

8 [1] Bayer Crop Science
9  700 Chesterfield Parkway West
10  Chesterfield, Missouri, 63017-1732, USA

11

12

13 *Corresponding author

14

# Abstract

The comparative assessment of genetically-modified (GM) crops relies on the principle of

substantial equivalence, which states that such products should be compared to conventional

counterparts that have an established history of safe use. In an effort to operationalize this

principle, the GMO Panel of the European Food Safety Authority proposed an equivalence test

that directly compares a GM test variety with a set of unrelated, conventionally-bred reference

varieties with part of the difference as the known background of the test (the same as the given

control). The criterion of the EFSA test, however, is defined solely by genotypic differences

between the non-traited control and reference varieties (i.e. the background effect) while

assuming the so-called GM trait effect as zero. As the outcome of an EFSA equivalence test is

determined primarily by the similarity, or lack thereof, of the control and references, a

conditional equivalence criterion is proposed in this investigation that focuses on "unintended"

effects of a GM trait which is irrespective of the (random) genotypic value of a given control.

The new criterion also includes a mean-scaled standard similar to the 80-125% rule for

bioequivalence assessment practiced in the pharmaceutical industry as an alternative when the

reference variation is zero or close to zero. In addition, optional criteria are proposed with a step-

wise procedure to control the rate of false negatives (non-equivalence by chance) providing a

comprehensive assessment under multiple comparisons. An application to maize grain

composition data demonstrates that the conditional equivalence criterion provides effect-specific

and more robust assessment of equivalence than the EFSA criterion did, especially for GM traits

showing negligible or no unintended effects which are likely true for most traits in the current

market.

# Introduction

37

38    The comparative assessment of foods derived from genetically-modified (GM) crops relies

39    on the principle of substantial equivalence, which states that such products should be compared

40    to conventional counterparts that have an established history of safe use but are not required to

41    have zero difference from a near-isogenic control line absent of a GM trait in terms of "natural

42    variations" [1-4]. In an effort to operationalize this principle, the GMO Panel of the European

43    Food Safety Authority (EFSA 2010) proposed an equivalence criterion (thereafter called EFSA

44    equivalence criterion or limits) that compares a GM test variety with a set of conventionally-bred

45    references with part of the difference as the known genotypic background of the test (the same as

46    the near-isogenic control line) [5,6]. Similar criteria have also appeared in the literature [7,8].

47    Nevertheless, Codex states that "in achieving the objective of conferring a specific target trait

48    (intended effect) to a plant …, additional traits … could be lost or modified (unintended effects)"

49    and "the safety assessment of foods derived from recombinant-DNA plants involves methods to

50    identify and detect such unintended effects and procedures to evaluate their biological relevance

51    and potential impact on food safety" [2]. As described above EFSA's method requests an

52    assessment of differences of the test from a set of references regardless of differences from the

53    control. These differences would contain a trait effect if presents and, however, are certainly

54    driven by genotypic values of the control resulting from conventional plant breeding (as

55    described by Jiang et al. [9]). Thus, the result of EFSA equivalence testing in practice often is

56    unrelated to the trait effect, which should be the sole focus of the comparative assessment,

57    creating a series of discussions in the recent literature [7-15].

58    As the outcome of an EFSA equivalence test is determined primarily by the similarity, or

59    lack thereof, of the control and references (Fig 1), a conditional equivalence criterion is proposed

60    in this investigation that focuses on "unintended effects" of a GM trait irrespective of the

61    (random) genotypic values of a given control.

62

63    **Fig 1. Graphical illustration of EFSA and conditional equivalence criteria.** EFSA

64    equivalence criterion is defined primarily by a (random) control background effect and has no

65    specification of a GM trait effect, and a conditional equivalence criterion is defined solely for a

66    GM trait effect with a given control.

67

68       To our knowledge, no equivalence criterion of a GM trait effect (with a given genotypic

69    background) has been developed for the comparative assessment of GM crops and derived

70    food/feed. All criteria in the literature are defined on the basis of the background effect assuming

71    the absence of a GM trait effect [6,7,9] which have shown at least three limitations in practice:

72    first is the sensitivity of the equivalence conclusion to a given control (i.e., the same GM product

73    can generate different, completely contradictory, conclusions, based solely on the selected

74    control) (Fig 1); second is the incomplete coverage of a one-size-fits-all criterion for a wide

75    range of endpoints, with dramatic differences in their means and variations as explained in the

76    next section; and third is the lack of a strategy for a control of false negatives (non-equivalence,

77    in this case) in a comprehensive assessment under multiple comparisons as the criterion defined

78    by a fixed percentile without an adjustment for the number of comparisons.

79       Here, a conditional equivalence criterion is derived on the basis of an expected mean squared

80    difference of a GM crop from the references using a mixed model approach assuming the

81    random background variation, similar to that used by the Food and Drug Administration (FDA)

82    for individual bioequivalence in the presence of a random, individual-specific effect [16]. When

83    the reference variation for the background effect is too low to provide a valid (variation-scaled)

84    criterion, a mean-scaled criterion, similar to the 80-125% rule for the bioequivalence assessment

85    in the pharmaceutical industry, is recommended as an alternative. Due to the alleviated false

86    negative rate (much higher than the target level of 5% as the EFSA criterion defined by a 95%

87    confidence interval), a data-driven procedure is proposed for selecting criteria with optional

88    criteria to statistically control these errors. An application to a maize grain composition example

89    (used by EFSA) demonstrates that the proposed conditional equivalence criteria provides

90    substantial improvement for a true similarity measurement of GM crop over three limitations of

91    EFSA criterion and others.

92        The organization of the manuscript starts with basic assumptions of the principle of

93    substantial equivalence, followed by the derivation of a set of conditional equivalence criteria,

94    and then a data-driven procedure for selecting criteria across various endpoints in practice, and

95    finally an application to a maize grain composition example. The discussion includes additional

96    thoughts on each of these new criteria, and highlights areas of further research.

97

# The definition of substantial equivalence

## Assumptions and the current practice

100    Assume a Test-Control-Reference (TCR) trial where the test is the GM variety, the control is

101    an isogenic non-GM comparator with the genotypic background similar to the test (as monitored

102    by the molecular breeding technique [17,18]), and the references are comparators with different

103    genetic backgrounds. Let $(\mu_T, \mu_C, \mu_R)$ denote parameters of genotypic group means, $\sigma_g^2$ the

104      genetic variation among references, and ($\Delta_{TC}$, $\Delta_{CR}$, $\Delta_{TR}$) the mean differences among three

105      groups. For parameters of interest, $\Delta_{TC}$ represents a GM trait effect with a given control, while

106      $\Delta_{CR}$ is an effect of the genotypic background shared by the test and the control from the

107      traditional plant breeding. A simple difference of a GM test from the reference mean $\Delta_{TR}$ ( $= \Delta_{TC}$

108      $+ \Delta_{CR}$) as expected consists of both GM trait effect $\Delta_{TC}$ and background effect $\Delta_{CR}$.

109      The following are the underlying assumptions of the principle of substantial equivalence and

110      the current global regulatory approval procedure of a GM crop.

111      (a) GM crop safety assessment is solely interested in the effect of a GM trait regardless of the

112      genotypic background;

113      "It focuses on assessing the safety of any identified differences so that the safety of the

114      new product can be considered relative to its conventional counterpart" [2]. Though the

115      regulatory evaluation of a GM crop is on a given control background, upon approval, the GM

116      trait could be integrated into any conventional reference (in the current market or from a

117      breeding program) during the commercial application, and the background effect $\Delta_{CR}$ is

118      expected to vary from endpoint to endpoint for a given control or for the same endpoint

119      across different controls [17,18,19]. An equivalence of a GM crop should focus solely on a

120      GM trait effect $\Delta_{TC}$ ( $= \Delta_{TR} - \Delta_{CR}$) regardless of the genotypic background effect $\Delta_{CR}$ of a

121      given control (Fig 1).

122      (b) Substantial equivalence of a GM crop in statistics is a similarity measure to a distribution of

123      conventional references with a history-of-safe-use;

124      "Any observed differences should be assessed in the context of the range of natural

125      variations" [2] demonstrated by conventional references with no requirement of a trait effect

126    $\Delta_{TC} = 0$. In spite of a given control was applied in a TCR trial, the equivalence of a GM crop

127    in statistics is a similarity or distance measurement of the mean difference $\Delta_{TC}$ between two

128    probability distributions, one for GM crop with various genotypic backgrounds and one for

129    conventional references, in the scale of the reference variation $\sigma_g$.

130  (c) Equivalence conclusion of a GM crop (or a GM trait) relies on the totality of evidence across

131    key components when compared a given control background.

132    Codex guidelines state [2] that "A variety of data and information are necessary to assess

133    unintended effects because no individual test can detect all possible unintended effects or

134    identify, with certainty, those relevant to human health. These data and information, when

135    considered in total, provide assurance that the food is unlikely to have an adverse effect on

136    human health". In practice a comprehensive assessment has been performed over a wide

137    range of endpoints from various studies e.g. often > 50 analytes in a composition study alone,

138    and any experimental deviation from equivalence has been evaluated in terms of the "natural

139    variation" as well as the nominal level of the statistical significance. With such an

140    assessment, any limitation of the evidence due to a given control is minimized and an

141    equivalence of a GM trait if concluded could be assumed for different genotypic background

142    after the commercialization.

143    In summary, by OECD and WHO/FAO guidelines, GM crop safety assessment is

144  characterized by multiple comparisons of a GM crop across key endpoints with the conventional

145  references with a given (control) genotypic background. Three features of the assessment, focus

146  on the GM trait effect, a wide range of background variations, and multiplicity of the

147  comparisons, are considered in the following section for the derivation of a set of conditional

148  equivalence criteria.

149

## A conditional equivalence criterion for similarity

151      Let $(D_{TC}, D_{CR}, D_{TR})$ denote estimates of $(\Delta_{TC}, \Delta_{CR}, \Delta_{TR})$ with variances $\left(\sigma^2_{D_{TC}}, \sigma^2_{D_{CR}}, \sigma^2_{D_{TR}}\right)$. , At

152 first, equivalence criteria of EFSA [6] and Vahl and Kang [8] from the current literature were

153 formulated in the notation of this manuscript. Then a conditional equivalence criterion was

154 derived with a mixed model approach and relationships among three types of criteria were

155 discussed.

156      EFSA equivalence criterion (or limit) was defined by the following equation under EFSA

157 model 2, an ad hoc model assuming both test and control in a TCR trial as random and

158 independent varieties with the same variance as a reference but unspecified means.

159      $$\frac{|D_{TR}|}{\sigma_{D_{TR}}} < \theta_{EFSA} \tag{1}$$

160 where $D_{TR}$ as stated repeatedly is a mixture of the GM trait effect and the control background

161 effect, $\sigma_{D_{TR}}$ consists of both genetic and residual variations (due to sampling), and $\theta_{EFSA}$ was

162 specified by a 95% confidence limit of a t distribution with a sample estimate of $\sigma_{D_{TR}}$ in practice.

163 Clearly, EFSA equivalence criterion considered only the background variation of the test variety

164 (shared with the control) and no GM trait effect which, if presented, would have adopted a non-

165 central t distribution for $\theta_{EFSA}$.

166      Vahl and Kang's scaled average equivalence criterion is defined in a similar way but in the

167 scale of the genetic portion of the background variation [8].

168      $$\frac{|D_{TR}|}{\sigma_g} < \theta_{VK} \tag{2}$$

169 where $\sigma_g^2$ is the leading term of the genetic variation of $\sigma_{D_{TR}}^2$ in (1), and $\theta_{VK}$ was specified by a

170 95% confidence limit of a standard normal. Nevertheless, underlying assumptions, i.e. the GM

171 trait effect is zero and the background of a GM test variety being the same as a random

172 reference, are the same for both criteria (1) and (2).

173  Let $E_E$ and $E_C$ denote respective expectation with respect to the environmental effect, such as

174 site, replicate and residual, and the control background effect following the distribution of

175 references. A mixed model approach for a fixed effect $\Delta_{TC}$ and a random effect $\Delta_{CR}$ assumes $E_C$

176 $[E_E(D_{TR})] = E_C[\Delta_{TR}] = E_C[\Delta_{TC} + \Delta_{CR}] = \Delta_{TC}$, and $E_C\{[E_E(D_{TR})]^2\} = E_C\{\Delta_{TR}^2\} = \Delta_{TC}^2 + \sigma_g^2$.

177 When applying to Vahl and Kang's scaled average equivalence, the following equation can be

178 derived

179 $$\frac{E_C\{[E_E(D_{TR})]^2\}}{\sigma_g^2} = \frac{\Delta_{TC}^2}{\sigma_g^2} + 1.$$

180 While the first expectation $E_E(\cdot)$ indicates the test statistic and the parameter of interest the

181 same as for criteria (1) and (2), the second expectation $E_C(\cdot)$ reveals the change of hypothesis if

182 a control background were assumed as random. Clearly, $\Delta_{TR}$ is for an equivalence with the given

183 control background as one part of the assessment, and $\Delta_{TC}$ is for a marginal equivalence or a

184 conditional equivalence for a random background if absence of an interaction. Since $E_E$

185 $[E_C(D_{TR})] = E_C[E_E(D_{TR})] = \Delta_{TC}$, $E_C(D_{TR})$ would represent an estimate of $\Delta_{TC}$, and an

186 equivalence criterion for a GM trait effect could be defined implicitly by

187 $$\frac{|E_C(D_{TR})|}{\sigma_g} \leq \theta_c \qquad\qquad (3)$$

188 where $\theta_c$ is a conditional equivalence criterion to be discussed in the following. This mixed

189 model approach has been applied by FDA for individual equivalence [16] and by Vahl and Kang

190     for a distribution-wise equivalence in GM crop assessment [8]. Therefore, a conditional

191     equivalence is defined solely for a GM trait effect $\Delta_{TC}$, and a criterion $\theta_c$ could be derived as a

192     function of the background variation (Fig 1).

193     Let $E_C(D_{TR}) = (D_{TR} | D_{CR} = 0)$ denote a conditional difference with a mean $\mu_{D_{TC}|D_{CR}=0}$ which

194     has been shown to be the best estimator of the trait effect based on the correlation structure

195     among three genotypic group means in a TCR trial [15]. Note that $\Delta_{CR} = 0$ is a key assumption

196     in defining the natural variation for GM crop safety assessment. Thus, $\mu_{D_{TR}|D_{CR}=0}$ and $\Delta_{TC}$ are

197     interchangeable in terms of the parameter value (or the equivalence criterion), but a statistic for

198     $\mu_{D_{TR}|D_{CR}=0}$ would be applied for an equivalence testing.

199     Parallel comparisons of criteria (1), (2), and (3) could be made using mean-squares of

200     expected differences i.e. $[E_E(\cdot)]^2$ to demonstrate differences in the statistical hypotheses.

201
$$\begin{cases} \Delta_{TC}^2 < \theta_c^2 \sigma_g^2 & Conditional\ Criterion \\ \Delta_{TC}^2 < \theta_{VK}^2 \sigma_g^2 - 2\Delta_{TC}\Delta_{CR} - \Delta_{CR}^2 & Criterion\ of\ Vahl\ and\ Kang \\ \Delta_{TC}^2 < \theta_{EFSA}^2 \sigma_{D_{TR}}^2 - 2\Delta_{TC}\Delta_{CR} - \Delta_{CR}^2 & EFSA\ Criterion \end{cases}$$

202     First, by criteria of EFSA and Vahl and Kang, equivalence of a trait effect $\Delta_{TC}$ would be

203     largely determined by the background effect $\Delta_{CR}$, not only the magnitude but also its direction.

204     Opposite signs of $\Delta_{TC}$ and $\Delta_{CR}$ would be much more likely to be concluded as equivalent than

205     those with the same sign do. In addition, the probability thresholds for $\theta_{EFSA}$ and $\theta_{VK}$ assume $\Delta_{TC}$

206     $= 0$, not a requirement of the principle of substantial equivalence. In contrast, the conditional

207     equivalence criterion does not assume $\Delta_{TC} = 0$, but $\Delta_{TC} = 0$ is expected to provide a maximum

208     chance of concluding equivalence regardless of the sign and magnitude of $\Delta_{CR}$ for a given

209     control.

210    Second, a conditional equivalence standard could be derived as $\theta_c = \sqrt{\theta_{VK}^2 - 1}$ e.g. $\theta_c =$

211    $\sqrt{z_{0.975}^2 - 1} \approx 1.69$ of which an interpretation could be provided as follows. For a GM test with a

212    trait effect $|\Delta_{TC}| < 1.69\sigma_g$, when averaging over the background effect, the GM test crop would

213    be within the 95% confidence interval of a reference. Therefore, the conditional criterion $\theta_c\sigma_g$ is

214    for a trait effect and defined by the range of reference variation thus follows the OECD

215    guidelines [1,2,3].

216    Third, EFSA equivalence criterion is a much loosely defined criterion as a function of the

217    experimental design (i.e. numbers replicates and sites and total number of references) with $\theta_{EFSA}$

218    $> \theta_{VK}$ and $\sigma_{D_{TR}}^2 > \sigma_g^2$, which compromises the efficacy of detecting an unintended effect if

219    presents as pointed out by Vahl and Kang [8, p23], and tends to encourage a trial with lower

220    number of reference and consequently lead to an arbitrary conclusion of equivalence as

221    concerned by Ward et al. [10].

222    In summary, a conditional equivalence criterion independent of the background was derived

223    in this section in the scale of the reference variation, thus called a variation-scaled criterion.

224    However, while criteria of EFSA and Vahl and Kang are all variation-scaled criteria, if applied

225    as a one-fits-all criterion, certain problems are inevitable. Firstly, even though EFSA classifies

226    those cases with zero estimate of $\sigma_g$ as "Equivalence Not Concluded", an arbitrary conclusion is

227    expected as $\sigma_g$ becomes less than certain threshold (relative to the residual variation) due to a

228    large proportion of close to zero criterion. A second problem is that, with a criterion defined by a

229    95% confidence limit, false negative (i.e. non-equivalence by chance) is expected to be at least

230    5% for each endpoint and would be much higher due to the proof-of-equivalence. While a

231    comprehensive assessment requires a totality of evidence, optional criteria become necessary.

232

## Alternative criteria in a comprehensive assessment

### An empirical mean-scaled criterion when $\sigma_g$ is low

235    Alternative criteria are discussed in this section when the reference variation is too low for a

236    variation-scaled criterion. When $\sigma_g^2$ is low, references in a TCR trial become similar to each

237    other including the control. Consequently, the equivalence of a GM crop may become the same

238    as the bioequivalence of a generic drug to a brand-named reference in pharmaceutical industry in

239    terms of the comparison between the test and the control and the absence of references.

240    The 80-125% rule has long been adopted in pharmaceutical industry for two drugs being

241    "similar to such a degree that their effects, with respect to both efficacy and safety, will

242    essentially be the same" [20,21]. This standard is also recommended for GM crop equivalence

243    assessment in this research. That is, for endpoints with low $\sigma_g$, a conditional equivalence would

244    be defined in a mean-scale by

245 $$\delta_c = \frac{\Delta_{TC}}{\mu_R} \tag{4}$$

246    Under the 80-125% rule, standards were differentiated between $\Delta_{TC} < 0$ or $\Delta_{TC} > 0$. For

247    simplicity $\delta_c = 0.25$ is applied in this research, and under a log-transformation as recommended

248    by FDA, the criterion becomes $\Delta_{TC} = \pm \log(1 + \delta_c) = \pm \log(1.25)$.

249    While the variation-scaled standard $\theta_c = 1.69$ is based on the concept of equivalence under

250    natural variation with a history-of-safe-use, the mean-scaled standard $\delta_c = 0.25$ is empirical.

251    Two standards appear to be independent in theory, but in practice they are highly correlated as

252     will be shown in the following maize grain composition example. Let $CV_g = \sigma_g/\mu_R$ denote a

253     coefficient of genetic variation among references in the original scale of a TCR trial data. Two

254     equivalence standards would be equal, i.e. $1.69\sigma_g = 0.25\mu_R$, at $CV_g = 14.8\%$. When a log-

255     transformation is applied, a log-normal distribution of $y = log\,(x)$ is assumed and two standards

256     become the same in a log-scale, i.e. $1.69\sigma_{gy} = \log(1.25)$, at $CV_g = 13.3\%$ using $\sigma_{gy}^2 = log$

257     $(CV_g^2 + 1)$, $\sigma_{gy}^2$ as the reference variance in the log-scale. As will be shown in the example, $CV_g$

258     $= 13.3{\sim}14.8\%$ are near the center of the maize grain composition data, and two standards 1.69

259     $\sigma_g$ and $0.25\mu_R$ are in fact largely overlapped due to a narrow range of $CV_g$ and a wide separation

260     of means across analytes.

261

## 262     Optional criteria for multiple comparisons

263      Equivalence criteria of EFSA and Vahl and Kang in the previous section did not apply a

264     whole range of reference variation, and with a 95% confidence limit a minimum 5% false

265     negative is expected even with no trait effect and could be much higher due to the proof-of-

266     equivalence (as shown in the following example). The same is true for a conditional equivalence.

267     Therefore, an optional criterion $\theta_c = \sqrt{z_{0.995}^2 - 1} \approx 2.38$ corresponding to a 99% confidence

268     limit is recommended to control the number of false negative.

269      For a use of whole range of "natural variation" in a proof-of-equivalence, OECD provided

270     summary of some historic data including mean and range of maize and soybean compositions

271     [22,23]. In the meantime, EFSA adopted an intuitive evaluation using box plots of the test, the

272     control and references for analytes failed to conclude equivalence by the equivalence testing.
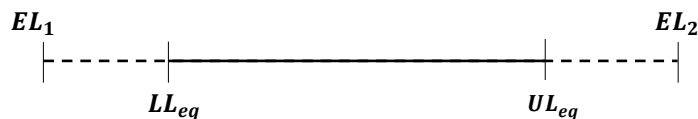
273    However, these approaches were unable to separate the true (genetic ) "unintended effect" of a

274    GM trait from those due to environments. In contrast, $\theta_c = 2.38$ considers only genetic variation

275    and provides a statistically interpretable standard of equivalence.

276        In addition, $\delta_c = 0.5$ is also recommended as an optional criterion for endpoints with low

277    reference variation following the practice of pharmaceutical industry for high variant endpoints

278    [24].

279

## A data-driven procedure for criterion selection

281        The following diagram describes the proof-of-equivalence approach as requested by EFSA,

282    similar to those in pharmaceutical industry [16,24,25]. The diagram can be applied on a mean of

283    a GM test or a difference between a GM test and references.



285    where $EL_1$ and $EL_2$ are lower and upper limits of an equivalence criterion, ($LL_{eq}$, $UL_{eq}$) is a

286    (confidence) interval containing critical values of the estimated mean or difference for

287    concluding equivalence. The dash lines consist of the estimated mean or difference failed to be

288    concluded as equivalent, called the burden of proof-of-equivalence, due to the margin-of-error in

289    comparisons with $EL_1$ and $EL_2$ at the designated significance level.

290        Assume a TCR trial with eight sites each of four replicates per site/treatment currently

291    requested by EFSA and generally accepted by most international regulatory agencies. The trait

292     effect is best estimated by a conditional difference with a mean $\mu_{D_{TC}|D_{CR}=0} = \Delta_{TC} + \frac{\sigma_{\hat{D}_{TC}}^2}{2\sigma_{\hat{D}_{CR}}^2}\Delta_{CR}$ and

293     a variance $\sigma_{\hat{D}_{TR}|D_{CR}}^2 = \sigma_{\hat{D}_{TC}}^2\left(1 - \frac{\sigma_{\hat{D}_{TC}}^2}{4\sigma_{\hat{D}_{CR}}^2}\right)$ as derived in Jiang et al. [15]. For simplicity, no genotype

294     by environmental interactions were assumed. With a given $\sigma_g{:}\sigma_e$ (or a given $\sigma_g$ at $\sigma_e = 1$), a

295     normal approximation was applied in the following for an asymptotic equivalence analysis using

296     an interval $(LL_{eq}, UL_{eq})$ as functions of $(EL_1, EL_2)$ defined by equations (3) and (4).

297        Let $EL_c = EL_2 = -EL_1$ for a GM trait effect. In this section, a threshold of $\sigma_g{:}\sigma_e$ for

298     alternating $EL_c = 1.69\sigma_g$ and $0.25\mu_R$, a key question in criterion selection, is investigated.

299     Obviously, no equivalence could be concluded even for $\Delta_{TC} = 0$ with $EL_c = 1.69\sigma_g$ if $\sigma_g{:}\sigma_e$ is

300     close to zero, and the threshold of $\sigma_g{:}\sigma_e$ should be large enough for $EL_c = 1.69\sigma_g$ to provide an

301     80% power of equivalence under a regulatory requested design. Fig 2 presents numerical results

302     of an asymptotic equivalence analysis (ignored the variation in estimating $\sigma_g$). Variation settings

303     $\sigma_g{:}\sigma_e = (0 \sim 3)$ in the left plot and the residual coefficient of variation $CV_e = (0 \sim 30\%)$ in the

304     right plot is based on the maize grain composition example in the following section. In the left

305     plot, a grid search was performed over $\Delta_{TC} = (0 \sim 2.38\sigma_g)$ for each value of $\sigma_g{:}\sigma_e$, the

306     maximum trait effect for 80% power of equivalence was obtained with $EL_c = 1.69\sigma_g$ and $2.38\sigma_g$

307     as functions of $\sigma_g{:}\sigma_e$. Similarly, in the right plot is for $EL_c = 0.25\mu_R$ and $0.5\mu_R$ as functions of

308     $CV_e$.

309

310     **Fig 2. Asymptotic maximum GM trait effects for 80% power of equivalence under EFSA**

311     **requested design.** Left plot: Asymptotic maximum GM trait effects with $EL_c = 1.69\sigma_g$ and 2.38

312      $\sigma_g$ as functions of $\sigma_g{:}\sigma_e$; Right plot: Asymptotic maximum GM trait effects with $EL_c = 0.25\mu_R$

313      and $0.5\mu_R$ as functions of residual coefficient of variation $CV_e$.

314

315      Results in Fig 2 indicate that, in general, an 80% power of equivalence requires a trait effect

316      $\Delta_{TC}$ substantially less than $EL_c$ due to the proof-of-equivalence. Asymptotically, a threshold $\sigma_g{:}$

317      $\sigma_e = 1.0$ could be applied for alternating $EL_c = 1.69\sigma_g$ and $0.25\mu_R$ (Fig 2). When $\sigma_g{:}\sigma_e \geq 1.0$,

318      the maximum trait effect for equivalence is about $1.4\sigma_g$ for $EL_c = 1.69\sigma_g$ and $2.1\sigma_g$ for $EL_c =$

319      $2.38\sigma_g$. When $\sigma_g{:}\sigma_e < 1.0$, even a negligible trait effect might not have enough power to

320      conclude equivalence by either criterion. By $EL_c = 0.25\mu_R$, the maximum trait effect for

321      equivalence is about $0.15\mu_R$ when $CV_e = 15\%$. Note that Fig 2 did not include the variation in

322      estimating $\sigma_g$ which should be a function of the number of references and, when considered, the

323      estimated maximum trait effect in Fig 2 could be substantially lower.

324      Results in Fig 2 demonstrate no existence of a one-fits-all criterion under practical ranges of

325      $\sigma_g{:}\sigma_e$ and $CV_e$. Therefore, the following set of conditional equivalence criteria was proposed

326      with a three-step procedure for criterion selection.

327     
$$EL_c = \begin{cases} EL_{c,\,0.95} = \begin{cases} 1.69\sigma_g & \text{when } \sigma_g{:}\sigma_e \geq 1 \\ 0.25\mu_R & \text{when } \sigma_g{:}\sigma_e < 1 \end{cases} \\ EL_{c,\,0.99} = \begin{cases} 2.38\sigma_g & \text{when } \sigma_g{:}\sigma_e \geq 1 \\ 0.5\mu_R & \text{when } \sigma_g{:}\sigma_e < 1 \end{cases} \end{cases} \qquad (5)$$

328      Firstly, $EL_{c,0.95} = 1.69\sigma_g$ is applied as the primary criterion whenever a reasonable variation $\sigma_g$

329      is available, i.e. $\sigma_g{:}\sigma_e \geq 1$. Secondly, an alternative $EL_{c,0.95} = 0.25\mu_R$ is used for $\sigma_g{:}\sigma_e < 1$.

330      Thirdly, $EL_{c,\,0.99}$ is optional for endpoints failed to show equivalence with $EL_{c,\,0.95}$ expectedly

331      only for a small proportion of endpoints say 5 to 10% or less.

332    In practice, the procedure (5) depends on estimates of $\sigma_g$, $\mu_R$, and $\sigma_g{:}\sigma_e$ and variations of

333    these estimates will be a function of the experimental design.

334

# Application to a maize grain composition example

336    Maize grain composition data in Jiang et al. [15], originally applied by EFSA for method

337    demonstration, were reanalyzed with and without the log transformation. At first, reference

338    means and variations were estimated from 13 references for each analyte, and paired estimates

339    across all 53 analytes for $EL_c = (1.69\sigma_g, 0.25\mu_R)$ without transformation and $EL_c = (1.69\sigma_g,$

340    $\log(1.25))$ with transformation were plotted (Fig 3). In the left plot with no transformation, a

341    strong linear correlation can be observed between estimates of $\sigma_g$ and $\mu_R$ (with an estimated $r^2$

342    $= 0.9644$ in the log-scale). The observed $CV_g$ $(= \sigma_g/\mu_R)$ is highly consistent within a range

343    $(0 \sim 25.2\%)$ and a mean $10.2\%$ across a wide range of $\mu_R$. The observed $CV_e$ has a mean $7.2\%$

344    and a range $(0.64 \sim 28.2\%)$. In the right plot with the log transformation, $EL_c = 1.69\sigma_g$ indeed

345    tends to be independent of the mean. The mean estimate of $EL_c = 1.69\sigma_g$ is 0.17, slightly lower

346    than the line $EL_c = log(1.25) \approx 0.22$ in the plot.

347

348    **Fig 3. EFSA example: Estimated equivalence criteria $EL_c = 1.69\sigma_g$ (markers) and $0.25\mu_R$**

349    **(line) across 53 analytes.** Left plot: Without transformation; Right plot: With log

350    transformation.

351

352      Results from Fig 3 have at least two implications. Firstly, high comparability between

353    equivalence criteria $EL_c = (1.69\sigma_g, 0.25\mu_R)$ strongly supports the equivalence evaluation of a

354    GM trait effect as a similarity measurement for a whole range of the reference variation. While

355    the concept of substantial equivalence is generally understood as comparing a GM crop with the

356    reference variation, a mean-scaled criterion in the original unit is a natural alternative when $\sigma_g$ is

357    small and lack of a good estimate. Secondly, $EL_c = 1.69\sigma_g$ on average appears to be more

358    conservative than $EL_c = 0.25\mu_R$, an empirical support for $EL_c = 1.69\sigma_g$ in terms of mean percent

359    difference.

360        Table 1 summarizes the estimate of $\mu_{D_{TC}|D_{CR}=0}$ in scales of $\sigma_g$ and $\mu_R$ (i.e. the same scale as

361    the criterion) and $\sigma_{D_{TC}|D_{CR}}$ (i.e. the t value of a conditional difference test by results of EFSA

362    model 1 and model 2). Values in bold follow the procedure (5). Though only two analytes with

363    zero estimate for $\sigma_g$, 9 and 10 analytes are estimated as $\sigma_g:\sigma_e < 1$ when no transformation or

364    transformation was applied.

365

366    **Table 1. Summary of estimates of $\mu_{D_{TC}|D_{CR}=0}$ in ratios to $\sigma_g$, $\mu_R$, and $\sigma_{D_{TC}|D_{CR}}$ (results**

367    **following the procedure (5) marked as bold) and comparisons of three criteria with**

368    **difference $D = D_{TR}$ for EFSA and Vahl and Kang (VK) methods, and $D = \hat{\mu}_{D_{TC}|D_{CR}=0}$ for**

369    **conditional equivalence (CD).**

| $\sigma_g : \sigma_e$ | Mean estimate and range of $\lvert\mu_{D_{TC}|D_{CR}=0}\rvert$ in ratios to | | | # Analytes (with $\lvert D \rvert > EL$) | | |
|---|---|---|---|---|---|---|
| | # Analyte | $\sigma_g$ | $\mu_R$ | $\sigma_{D_{TC}|D_{CR}}$ | EFSA | VK | CD |
| In original unit | | | | | | |

| ≥ 1 | 44 | **0.34(0.00~1.12)** | 0.04(0.00~0.13) | 1.48(0.01~5.25) | 2 | 3 | 0 |
|---|---|---|---|---|---|---|---|
| < 1 | 9 | 1.35(0.59~2.12)* | **0.09(0.02~0.19)** | 1.90(0.93~3.80) | 1* | 2* | 0 |
| With log transformation | | | | | | | |
| ≥ 1 | 43 | **0.37(0.01~1.24)** | 0.04(0.00~0.14) | 1.61(0.03~6.47) | 2 | 5 | 0 |
| < 1 | 10 | 1.02(0.17~1.81)* | **0.09(0.02~0.22)** | 1.69(0.35~3.79) | 1* | 1* | 0 |

370     *: Two analytes with zero estimate of $\sigma_g$ are not included.

371

372     While large t values for estimates of $\mu_{D_{TC}|D_{CR}=0}$ in the ratio to $\sigma_{D_{TC}|D_{CR}}$ are evidence of non-

373     zero trait differences, results in bold demonstrate the good performance of the procedure (5). For

374     example, for those analytes with $\sigma_g{:}\sigma_e < 1$ and $EL_c = 0.25\mu_R$ without transformation or $EL_c =$

375     log (1.25) with transformation, estimates of $\mu_{D_{Tc}|D_{CR}=0}$ are all well within the limit (i.e. 0.25) in

376     the scale of $\mu_R$, but several of them would have exceeded the limit (i.e. 1.69) if $EL_c = 1.69\sigma_g$

377     was applied. Results are highly consistent with and without transformation.

378     The last three columns of Table 1 summarize comparisons of the estimated difference with

379     three criteria: EFSA, VK for Vahl and Kang, and CD for the conditional equivalence. With a

380     proof-of-equivalence, an estimated difference exceeding *EL* would automatically lead to a non-

381     equivalence conclusion. $|D_{TR}| > EL$ were observed for both criteria of EFSA and VK. Three

382     cases of $|D_{TR}| > EL$ for EFSA criterion represents almost exactly a 5% of non-equivalence by

383     chance, 2.6 (i.e. 5% of 51) in expectation which simply suggests no evidence of GM trait effects.

384     These results are the same as EFSA original analysis with the transformation, where three

385     analytes were classified as "Non-Equivalence More Likely Than Not" or "Non-Equivalence".

386     Yet a total of nine analytes (i.e. 17% of 53) failed to conclude equivalence including two with

387     zero estimate of $\sigma_g$, much higher than the nominal 5% level due to the proof-of-equivalence.

388    However, for the conditional criteria no $\left|\mu_{D_{TC}|D_{CR}=0}\right| > EL_c$ was observed. Note that although no

389    trait effect exceeds $EL_{c,0.95}$ in Table 1, $EL_{c,0.99}$ might still be necessary in a formal testing due to

390    the proof-of-equivalence.

391    In summary, despite of a formal statistical testing yet to be developed, simple comparisons in

392    the example demonstrate obvious advantages of the conditional equivalence criteria proposed in

393    this investigation over those of EFSA and Vahl and Kang in a comprehensive equivalence

394    assessment.

395

# Discussion

397    Under OECD guidelines GM crop safety assessment is characterized by the comparative

398    approach on a GM crop of known (control) genotypic background to conventional references

399    with a history-of-safe-use. The equivalence testing method prescribed by EFSA assesses

400    differences between the GM variety and a group of commercial reference varieties. These

401    differences, as discussed by Jiang et al. [9], may be driven by a trait effect, a known control

402    background effect, or both. The EFSA equivalence criterion consists of only the background

403    variation, and three direct consequences are worth noting. First, the EFSA criterion is entirely for

404    the random background effect, contradicting with the principle focusing on "unintended effect"

405    of a GM trait. Second, the EFSA criterion becomes degenerate when the reference variation is

406    low and estimated as zero or close to zero, a common case in composition studies. Third, it is the

407    inability to control the false negative rate, i.e. substantially higher than the target 5% level as

408    defined by criteria of EFSA [6] or Vahl and Kang [8] due to the proof-of-equivalence even in the

409    absence of a true GM trait effect.

410       A set of conditional equivalence criteria under the same assumptions as those of EFSA and

411       Vahl and Kang are derived in this manuscript. However, the new criteria are for a GM trait effect

412       $\Delta_{TC}$, which is independent of the genotypic background of a given control and thus

413       fundamentally different from those of EFSA and others.

414       The approach using the reference variation $\sigma_g$ as a scale in this manuscript, if available,

415       follows the principle of substantial equivalence of OECD guideline that "any observed

416       differences should be assessed in the context of the range of natural variations" [2]. When the

417       reference variation $\sigma_g$ is small, an alternative criterion was proposed to apply the scale of the

418       reference mean $\mu_R$ with the procedure (5). The parallel nature of the mean-scaled and the

419       variation-scaled criteria lies in the definition of equivalence to a fixed reference (with a mean

420       percentage difference as an empirical standard) or to a group of references (with a fold of

421       standard deviation defining the range of a history-of-safe-use). Re-analysis of the maize grain

422       composition example originally applied by EFSA illustrate that even though only two endpoints

423       have $\sigma_g$ estimated as zero, low values of $\sigma_g$ relative to the residual are common (about 20% in

424       the example by the threshold $\sigma_g{:}\sigma_e < 1$). In these cases, a mean-scaled criterion as defined in (5)

425       should be considered as a natural alternative to a variation-scaled criterion, which is strongly

426       supported by the close correlation between the mean and variation in the example (Fig 3)  and

427       commonly observed in biological literature (e.g. the log-transformation suggested by EFSA and

428       FDA in pharmaceutical studies).

429       Another type of criterion in the literature for endpoints with low values of $\sigma_g{:}\sigma_e$ could be

430       labeled as "phenotypic equivalence" due to including residual variation $\sigma_e$ (and other

431       environmental variations) as part of the "natural variation" in defining equivalence. One example

432       is the distribution-wise equivalence proposed by Vahl and Kang [8] and applied by Van der Voet

433 et al. [26] in the analysis of five studies in which GM crops were fed to rats. Another example is

434 the criterion applied by Schmidt et al. [27], based on one unit of reference standard deviation, i.e.

435 $\sqrt{\sigma_g^2 + \sigma_e^2}$ in expectation, following an example in the EFSA guideline [28]. An intuitive

436 interpretation of the "phenotypic equivalence" would be a large proportion of overlapping

437 between observed responses of the test and the control. However, no discussion could be found

438 on the level of "unintended effect" of a GM trait in the unit of $\sigma_e$, either in terms of the

439 regulatory policy or a biological interpretation. In addition, practical implications of these criteria

440 would depend on the type of the study (e.g. the applicable sample size) and the characteristic of

441 the endpoint (e.g. magnitudes of $CV_g$ and $CV_e$).

442     In their guideline, EFSA also proposed a simulation approach for evaluating equivalence by

443 an empirical distribution of the number of significant outcomes in the difference testing between

444 two independent references. Regardless of the residual variation, the absolute mean difference

445 between two references is $\sqrt{2}\sigma_g$ and a 95% confidence limit would be approximately $2.8\sigma_g$

446 under normality. From this perspective, $EL_c = 1.69\sigma_g$ and $2.38\sigma_g$ in (5) would be considered as

447 conservative. Therefore, even though, under certain circumstances, some assumptions would be

448 more plausible than others, the comparative assessment of GM crops should be a comprehensive

449 approach with false negatives under control.

450     A false negative in an equivalence testing is in many ways similar to the false positive in a

451 difference test. In the maize grain composition example, the EFSA criterion demonstrated an

452 almost exact 5% of the analytes with observed differences greater than the equivalence limit (i.e.

453 3/51 or 5.9%) (Table 1). However, due to the proof-of-equivalence approach, a much higher

454 proportion of analytes (i.e. 9/53 or 17%) failed to conclude equivalence by the EFSA method [6].

455    However, all analytes in the example are within the conditional limits $EL_{c,0.95} = 1.69\sigma_g$ or

456    $0.25\mu_R$. In addition, a step-wise procedure with optional criteria $EL_{c,0.99} = 2.38\sigma_g$ or $0.5\mu_R$ was

457    proposed in this investigation for endpoints failed to conclude equivalence at $EL_{c,0.95}$. In contrast

458    to use the box plots of EFSA or historic data [22,23,29], criteria $EL_{c,0.99}$ consist of no

459    environmental effects, thus are true criteria of "unintended effect" due to multiple comparisons.

460    References in the current TCR trial though limited in number often may be a more reliable

461    source of information due to difficulties in estimating the genotypic variation among references

462    from historic data.

463        With the criteria developed here, an immediate further research subject would be statistical

464    methods of conditional equivalence testing applying these criteria. Because these criteria must be

465    estimated from the reference data and the variation of the estimation must be taken into account

466    especially when the number of references is limited. The variation of the estimated criterion

467    would be compounded with those of the trait effect, not accounted for in the EFSA equivalence

468    testing, which partially contributed to its poor performance.

469

# Acknowledgments

472

# References

474   1.   OECD, Safety evaluation of foods derived by modern biotechnology: Concepts and
475        principles. 1993, Organisation for economic co-operation and development: Paris,
476        France.
477   2.   Codex Alimentarius, Guideline for the conduct of food safety assessment of foods
478        derived from recombinant-DNA plants. 2003, Codex Alimentarius Commission, Joint
479        FAO/WHO Food Standards Programme, Food and Agriculture Organization of the
480        United Nations: Rome, Italy.
481   3.   Codex Alimentarius, Foods derived from modern biotechnology. 2009, Codex
482        Alimentarius Commission, Joint FAO/WHO Food Standards Programme, Food and
483        Agriculture Organization of the United Nations: Rome, Italy.
484   4.   FAO. GM Food Safety Assessment: Tools for traners. Food and Agriculture
485        Organization of the United Nations Rome, 2008. ISBN 978-92-5-105978-4.
486   5.   EFSA, Commission Implementing Regulation (EU) No 503/2013 of 3 April 2013 on .
487        for authorisation of genetically modified food and feed in accordance with Regulation
488        (EC) No 1829/2003 of the European Parliament and of the Council and amending
489        Commission Regulations (EC) No 641/2004 and (EC) No 1981/2006 Text with EEA
490        relevance. 2013, European Food Safety Authority: Parma, Italy.
491   6.   EFSA, Scientific opinion on statistical considerations for the safety evaluation of
492        GMOs. EFSA Journal, 2010. 8: p. 1-59.
493   7.   Kang Q and Vahl CI. Statistical analysis in the safety evaluation of genetically-modified
494        crops: Equivalence tests. Crop Science, 2014. 54: p. 2183-2200.
495   8.   Vahl CI and Kang Q. Equivalence criteria for the safety evaluation of a genetically
496        modified crop: A statistical perspective. The Journal of Agricultural Science, 2015.
497        154(03): p. 383-406.
498   9    Jiang C, Meng C, Schapaugh A. Comparative analysis of genetically-modified crops:
499        Part 1. Conditional difference testing with a given genetic background. PLoS ONE,
500        2019. 14(1): e0210747. https://doi.org/10.1371/journal.pone.0210747
501   10.  van der Voet H, Perry JN, Amzal B, and Paoletti C. A statistical assessment of
502        differences and equivalences between genetically modified and reference plant varieties.
503        BMC Biotechnology, 2011. 11(15): p. 1-20.

504    11. Ward KJ, Nemeth MA, Brownie C, Hong B, Herman RA, and Oberdoerfer R.
505         Comments on the paper "A statistical assessment of difference and equivalences
506         between genetically modified and reference plant varieties" by van der Voet et al. 2011.
507         BMC Biotechnology, 2012. 12(13): p. 1-7.

508    12. Fernandez A and Paoletti C. Unintended Effects in Genetically Modified Food/Feed
509         Safety: A Way Forward. Trends in Biotechnology, 2018. 36(9): p. 872-875.

510    13. Simó C, Ibáñez C, Valdés A, Cifuentes A. and García-Cañas V. Metabolomics of
511         Genetically Modified Crops. Int. J. Mol. Sci. 2014. 15: p. 18941-18966;
512         Doi:10.3390/ijms151018941

513    14. Herman RA, Fast BJ, Scherer PN, Brune AM, de Cerqueira DT, Schafer BW, et al.
514         Stacking transgenic event DAS-O15O7-1 alters Herman RA less than traditional
515         breeding. Plant Biotechnology Journal, 2017. 15(10): p. 1264-1272.

516    15. Herman RA, Huang E, Fast B, and Walker C. EFSA Genetically Engineered Crop
517         Composition Equivalence Approach: Performance and Consistency. Journal of
518         Agricultral and Food Chemistry, 2019. 67(14): p. 4080-4088.

519    16. FDA, Guidance for industry: Statistical approaches to establishing bioequivalence. 2001,
520         U.S. Department of Health and Human Services, Food and Drug Administration, Center
521         for Drug Evaluation and Research.

522    17. Prado JR, Segers G, Voelker T, Carson D, Dobert R, Phillips J, et al. Genetically
523         Engineered Crops: From Idea to Product. Annu. Rev. Plant Biol. 2014. 65: p. 769–90.

524    18. Venkatesh TV, Cook K, Liu B, Perez T, Willse A, Tichich R, et al. Compositional
525         differences between near-isogenic GM and conventional maize hybrids are associated
526         with backcrossing practices in conventional breeding. Plant Biotechnology Journal
527         2015. 13: p. 200-210.

528    19. Horak MJ, Rosenbaum EW, Woodrum CL, Martens AB, Mery RF, Cothren JT, et al.
529         Characterization of Roundup Ready Flex Cotton, 'MON 88913', for use in ecological
530         risk assessment: Evaluation of seed germination, vegetative and reproductive growth,
531         and ecological interactions. Crop Sci 2007. 47: p. 268-77;
532         http://dx.doi.org/10.2135/cropsci2006.02.0063

533    20. Schall R and Endrenyi L. Bioequivalence: tried and tested. Cardiovasc J Afr 2010.
534         21(2): p. 69-70.

535    21. Kumar SA, Sanjita D. A Review article on Bioavailability and Bioequivalence Studies.
536        International Journal of PharmTech Research 2013. 5: p. 1711-1721.

537    22. OECD Environmental Health and Safety Publication Series on the Safety of Novel
538        Foods and Feeds. No. 6 Consensus Document on Compositional Considerations for New
539        Varieties of Maize (*Zea Mays*): Key Food and Feed Nutrients, Anti-nutrients and
540        Secondary Plant Metabolites. 2002. OECD Environment Directorate, Paris.

541    23. OECD Environmental Health and Safety Publication Series on the Safety of Novel
542        Foods and Feeds. No. 25 Revised Consensus Document on Compositional
543        Considerations for New Varieties of SOYBEAN [*Glycine max* (L.) Merr]: Key Food
544        and Feed Nutrients, Anti-nutrients, Toxicants, and Allergens. 2012. OECD Environment
545        Directorate, Paris.

546    24. Davit BM, Chen ML, Correr DP, Haidar SH, Kim S, Lee CH, et al. Implementation of a
547        Reference-Scaled Average Bioequivalence Approach for Highly Variable Genetic Drug
548        Products by the US Food and Drug Administration. The AAPS Journal 2012. 14(4): p.
549        915-924.

550    25. Hsu J, Hwang JTG, Liu HK and Ruberg SJ. Confidence intervals associated with tests
551        for bioequivalence. Biometrika 1994. 81: p. 103–114.

552    26. van der Voet H, Goedhart PW, Schmidt K. Equivalence testing using existing reference
553        data: An example with genetically modified and conventional crops in animal feeding
554        studies. Food and Chemical Toxicology 2017. 109: p. 472-485.

555    27. Schmidt K, Schmidtke J, Schmidt P, Kohl C, Wilhelm R, Schiemann J, et al. Variability
556        of control data and relevance of observed group differences in five oral toxicity studies
557        with genetically modified maize MON810 in rats. Arch Toxicol 2017. 91: p. 1977–2006.

558    28. EFSA Guidance on conducting repeated−dose 90-day oral toxicity study in rodents on
559        whole food/feed. EFSA J 2011. 9: p. 2438.

560    29. Herman RA, Scherer PN, Phillips AM, Storer NP, Krieger M. Safety composition levels
561        of transgenic crops assessed via a clinical medicine model. Biotechnol. J. 2010. 5(2): p.
562        172-182.

563

564

# Supporting information

565

**S1 Fig. Graphical illustration of EFSA and conditional equivalence criteria.** EFSA

566

equivalence criterion is defined primarily by a (random) control background effect and has no

567

specification of a GM trait effect, and a conditional equivalence criterion is defined solely for a

568

GM trait effect with a given control.

569

570

**S2 Fig. Asymptotic maximum GM trait effects for 80% power of equivalence under EFSA**

571

**requested design.** Left plot: Maximum GM trait effects with $EL_c = 1.69\sigma_g$ and $2.38\sigma_g$ as

572

functions of $\sigma_g$:$\sigma_e$; Right plot: Maximum GM trait effects with $EL_c = 0.25\mu_R$ and $0.5\mu_R$ as

573

functions of residual coefficient of variation $CV_e$.

574

575

**S3 Fig. EFSA example: Estimated equivalence criteria of $EL_c = 1.69\sigma_g$ (markers) and**

576

**$0.25\mu_R$ (line) across 53 analytes.** Left plot: Without transformation; Right plot: With log
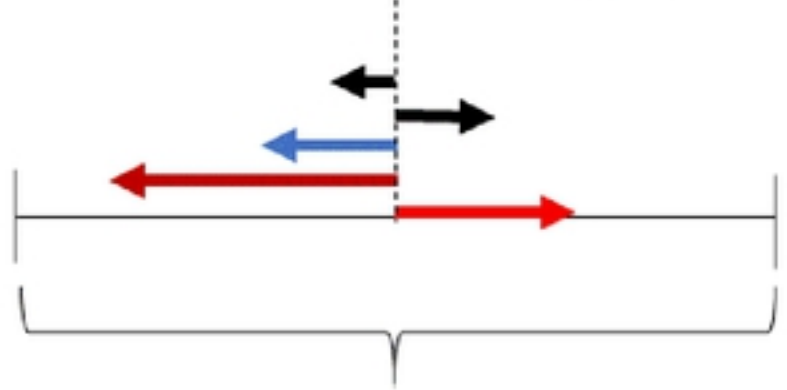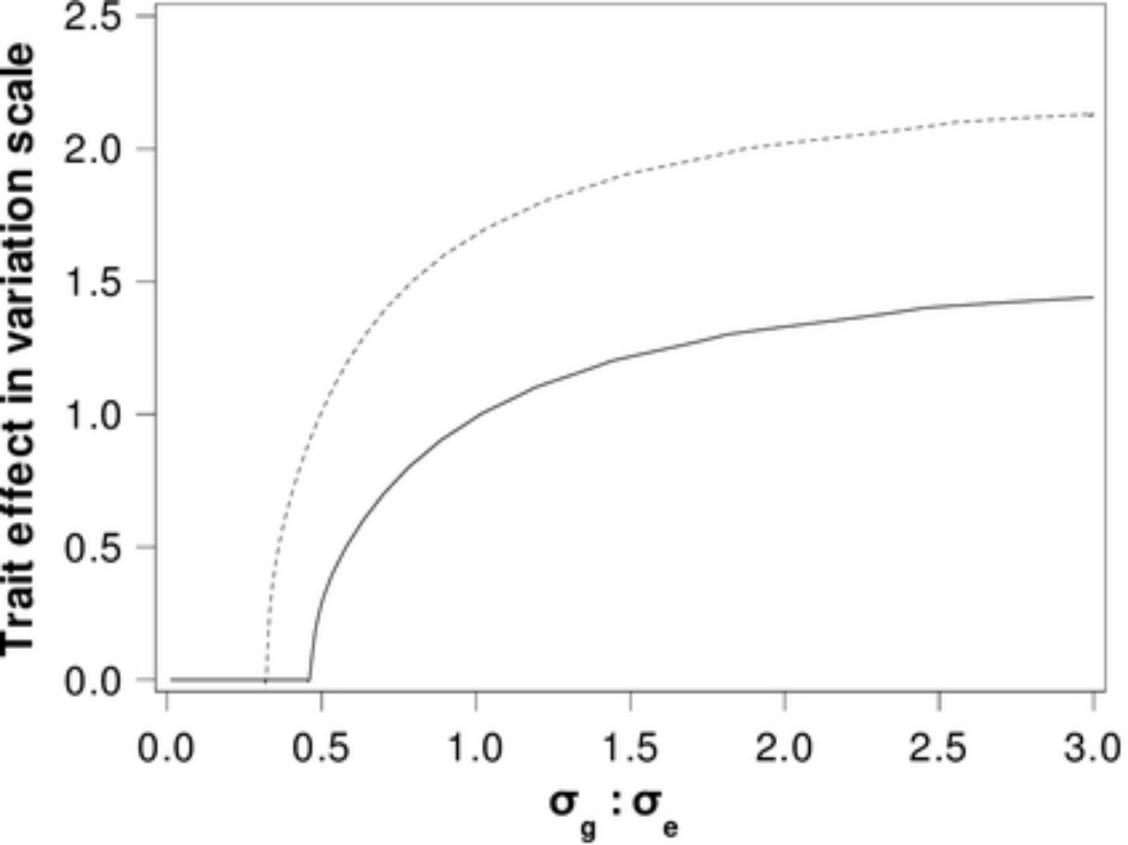
577

transformation.

578

579

580

Cover Letter

Cover Letter

Cover Letter