# Genomic Regions Associated with Important Seed Quality Traits in Food-grade Soybeans

Rachel M.Whiting, SepidehTorabi, Lewis Lukens, Milad Eskandari[*]

*Department of Plant Agriculture, University of Guelph, ON, Canada*

*Corresponding author Email:meskanda@uoguelph.ca

**Abstract**

**Background:** The production of soy-based food products requires specific physical and chemical characteristics of the soybean seed. Identification of quantitative trait loci (QTL) associated with these traits, such as seed weight,seed protein and sucrose concentrations could accelerate the development of competitive high quality soybean cultivars for the food-grade market through marker-assisted selection (MAS). The objectives of this study were to identify and validate QTL associated with these value-added traits in two high-protein recombinant inbred line (RIL) populations.

**Results:**Two RIL populations were derived from the high-protein cultivar 'AC X790P' (49% protein, dry weight basis), and two high-yielding commercial cultivars, 'S18-R6' (41% protein) and 'S23-T5' (42% protein). Fourteen large-effect QTL ($R^2>10\%$) associated with seed protein concentration were identified. Five of these protein-related QTLwere co-localized with QTL associated with seed sucrose concentration or seed weight. None of the protein-related QTL did not co-localize with seed yield QTL in either population. Sixteen candidate genes with putative roles in protein metabolism were identified within seven of these protein-related regions: qPro_Gm02-3, qPro_Gm04-4, qPro_Gm06-1, qPro_Gm06-3, qPro_Gm06-6, qPro_Gm13-4 and qPro-Gm15-3.

**Conclusion:**The use of RIL populations derived from high-protein parents created a unique opportunity to identify novel QTL that may have been masked by large-effect QTL segregating in populations developed from diverse parental cultivars. Nine QTL associated with seed protein concentration were identified and validated in both high-protein RIL populations. These QTL may be useful in the curated selection of new soybean cultivars for optimized soy-based food products.

**Key words:**Food-gradesoybean, protein,sucrose,seed weight, linkage analysis, candidate genes

## Background

Soybean [*Glycine max* (L.) Merrill] is a major source of plant-based dietary protein. An increased demand for whole-bean soy-based food products, such as tofu and soymilk, in western countries has attracted the attention of researchers, soybean growers and soy-based food processors. Soy-based products require specific physical and chemical characteristics of the soybean seed, including optimal seed protein concentration, seed sucrose concentration and seed weight [1-7], that are not of importance to commodity soybean breeding programs. As food processors require consistent seed composition to maintain production procedures, the development of environmentally stable, high yielding soybean cultivars with optimal value-added traits has become an important breeding objective.

Seed composition traits and yield are complex traits and affected by numerous genes and environmental factors[8-13]. Seed protein concentration shares a well-documented negative association with seed yield, which has hampered the development of competitive high-protein soybean cultivars [9, 14-23]. Additional value-added traits, such as high seed sucrose concentration and high seed weight, are also of interest to soy-food processors. Sucrose concentration is known to influence the palatability and texture of many soy-food products [24].However, seed protein and sucrose concentrations share a significant inverse relationship [25]. This relationship can be detrimental for soy-foods, such as tofu, that require high concentrations of both protein and sucrose for op0074imal production [5]. The identification and use of quantitative trait loci (QTL) associated with elevated seed protein concentration and additional value-added traits could accelerate the development of competitive high-protein soybean cultivars for the North American food-grade market by accumulating desirable alleles into a common genetic background.

Numerous studies have sought to determine the genetic basis of seed protein accumulation in soybean. SoyBase has indexed 248 bi-parental QTL associated with seed protein concentration, which encompass the results of more than 35 independent studies [26]. These QTL are located on every soybean chromosome, although chromosomes 6, 15, 18 and 20 are particularly favoured [27]. A QTL-meta analysis conducted by Qi et al. [28] also identified 51 consensus QTL across numerous genetic backgrounds and growing environments, which were located on all linkage groups except Chromosome 16. Many factors, such as large confidence intervals, small additive effects, negative associations with other desirable traits, poor environmental stability and QTL-by-genetic background interaction effects, have limited the usefulness of these QTL in marker-assisted selection programs [29-33]. Numerous QTL have also been identified for other traitsofinterest, including 318 seed weight-related QTL

2

29     identified in over 50 independent studies, and 188 seed yield-related QTL identified in 32 independent studies [26].

30     Sucrose concentration has received considerably less attention, with 37 sucrose-related QTL identified in 4

31     independent studies [26].

32          A global analysis of RNA-seq data revealed that Kunitz trypsin inhibitor 1, lectin family proteins, seed

33     storage 2S albumin superfamily proteins, bZIP homologues and MYB-like transcription factors were associated with

34     seed protein accumulation [28]. These transcripts were also associated with seed protein accumulation in previous

35     studies [34-36]. Specific genes, such as *ABI3, ABI4* and *LEC1* have also been associated with seed protein

36     accumulation [37, 38].

37          One method ofdetecting QTL that may be of use in improving polygenic traits is to utilize segregating

38     populations derived from elite parents [39]. Previous studies aimed at detecting protein-related QTL have mostly

39     used mapping populationsderived from exotic germplasm or parental cultivars with large phenotypic differences for

40     the desired traits [40]. Utilizing populations derived from elite lines may increase the chance of detecting novel QTL

41     that were masked by common large-effect QTL in diverse populations. These QTL have a higher change of being

42     beneficial for the development of new high-protein soybean cultivars.

43          In the present study, two recombinant inbred line (RIL) populations derived from crosses involving three

44     high-yieldingsoybean cultivars with high to moderately high-protein content were used to identify QTL associated

45     with traits important for food-grade soybean. Significant genomic regions associated with seed protein concentration

46     were examined for their relationship with seed sucrose concentrations, seed weight and yield. Identifying genomic

47     regions that underlie multiple value-added traits would be beneficial for the simultaneous improvement of desirable

48     traits in new food-grade soybean cultivars. To better understand the underlying mechanisms that regulate seed

49     storage protein accumulation in soybeans, these regions were also screened for putative candidate genes.

50

51     **Results**

52     **Phenotypic Analyses of Protein and Other Value-added Food-grade Traits**

53          The RIL populations were evaluated for seed weight, yield, protein and sucrose concentrations in multi-

54     environment trials during the 2015 and 2016 field seasons (Fig. 1; Supplementary Table 1-4).  Contrasts were noted

55     for seed protein concentration between the parental cultivars in both populations. In POPn_1, 'AC X790P' had an

56     average protein concentration of 48.08% (± 0.19%, standard error) across the five testing environments, while 'S18-

R6' had an average of 40.93% (± 0.19%). In POPn_2. 'AC X790P' had an average protein concentration of 48.24% (± 0.21%) across the five testing environments, while 'S23-T5' had an average of 42.60% (± 0.21%).

Differences in protein concentration between the RIL lines in each population were significant in the individual and combined environments (Fig. 1; Supplementary Table 1). In POPn_1, seed protein concentration varied from 41.53% to 45.27%, with an average protein concentration of 43.31% (± 0.03%). In POPn_2, seed protein concentration varied from 41.93% to 47.46%, with an average protein concentration of 44.61% (± 0.03%) (Fig. 1; Supplementary Table 1). Transgressive segregation was observed in some individual environments but was not observed when the combined environment data was considered (Supplementary Table 1). The normally distributed (Fig. 2) LSMEAN estimates for genotypes indicate that protein concentration is controlled by many genes.

The parental cultivars also differed for seed yield, seed weight and seed sucrose concentration, and considerable variation was also noted within the combined multi-environment data for both populations (Fig. 1). In POPn_1, entry seed weightestimates (grams per 100 seeds) varied from 18.08 grams to 23.88 grams, with an average seed weight of 21.18 grams (± 0.055 grams). Seed yield also varied from 2.55 tonnes ha$^{-1}$ to 4.49 tonnes ha$^{-1}$, with an average seed yield of 3.57 tonnes ha$^{-1}$ (± 0.025 tonnes ha$^{-1}$) and seed sucrose concentration varied from 5.44% to 6.82%, with an average sucrose concentration of 6.06% (± 0.016%; Supplementary Table 2-4).Similar variability was noted in POPn_2 (Fig. 1). Seed weight varied from 17.67 grams to 22.95 grams, with an averageseed weight of 20.34 grams (± 0.057 grams). Seed yield varied from 2.52 tonnes ha$^{-1}$ to 4.40 tonnes ha$^{-1}$, with an average seed yield of 3.34 tonnes ha$^{-1}$ (± 0.024 tonnes ha$^{-1}$) and seed sucrose concentration varied from 4.95% to 6.75%, with an average sucrose concentration of 5.84% (± 0.014%). Transgressive segregation was noted for seed yield and seed sucrose concentration in both populations. While some RILs exhibited transgressive segregation in individual environments for seed weight, this was not observed when the combined environment data was considered (Supplementary Table 2-4).

Our previous study revealed significant differences (p < 0.01) in genotype, environment, and genotype x environment treatments for protein concentrationand yield in these populations [41], which indicates the important role of genetic factors on the performance of these target traits. High heritability was noted for protein concentration and 100-seed weight ($H^2$ = 0.93-0.95 and 0.87-0.89, respectively; Table 1). Moderate heritability was observed for

84    sucrose concentration ($H^2$ = 0.70-0.81; Table 1), and low heritability was observed for seed yield ($H^2$ = 0.22-0.36)

85    (Table 1).

86

87    **Table 1** Broad-sense heritability of protein concentration, sucrose concentration, seed weight and seed yield in two

88    RIL populations evaluated in five environments (CHA15, CHA16, MER15, MER16 and PAL16)

89
90

|        | Protein | Yield  | Seed Weight | Sucrose |
|--------|---------|--------|-------------|---------|
| **POPn_1** | 0.9275  | 0.3603 | 0.8648      | 0.7035  |
| **POPn_2** | 0.9501  | 0.2180 | 0.8924      | 0.8132  |

91

92    **Relationships between Traits**

93    Pearson's correlation coefficients were used to determine the relationship between seed protein

94    concentration and sucrose concentration, seed weight and yield. Large, significant negative correlations were

95    observed between seed protein and sucrose concentration in both populations (POPn_1: r = -0.47; POPn_2: r = -

96    0.70; Fig. 2). InPOPn_1, seed protein concentration and seed weight were positively correlated (POPn_1: r = 0.53),

97    and seed weight and sucrose concentration were negatively correlated (POPn_1: r = -0.29). Interestingly,

98    nosignificant relationships were noted between seed protein concentration and seed yield in either population

99    (POPn_1: r = 0.09; POPn_2: r = -0.06) (Fig. 1; Fig. 2).

100

101   **SNP Mapping of the Soybean Genome**

102   Linkage maps were constructed from polymorphic SNP markers in each population. In POPn_1, a linkage

103   map was created using 807 SNP markers that were divided into 39 linkage groups. A linkage map consisting of

104   1,406 SNP markers on 40 linkage groups was created on POPn_2. All 20 chromosomes in the soybean genome were

105   represented, with most chromosomes consisting of two or more linkage groups. The linkage maps were 2,385 and

106   2,690 cM in length for POPn_1 and POPn_2, respectively. The number of linkage groups was attributed to a lack of

107   polymorphic markers between the parental genotypes distributed over large chromosomal regions, as elite Canadian

108   soybean cultivars may share similar pedigrees.

109

110   **QTL Associated with Seed Protein Concentration**

111    In total, from the analysis of both populations, fourteen large-effect QTL affecting protein content were

112    identified on Chromosomes 1, 2, 4, 5, 6, 8, 12, 13, 15 and 18. The fourteen QTL explained between 10.4% and

113    21.9% of the observed phenotypic variation (Table 2). Six of these QTL – *qProt_Gm01-2*, *qProt_Gm04-3*,

114    *qProt_Gm06-1*, *qProt_Gm06-3*, *qPro_Gm12-3*, and *qPro-Gm12-4* – carried the beneficial alleles from 'S18-R6' or

115    'S23-T5', while the remaining eight QTL – *qProt_Gm02-3*, *qProt_Gm04-4*, *qPro-Gm05-2*, *qPro-Gm06-6*, *qPro-

116    Gm08-2*, *qPro-Gm13-4*, *qPro_Gm15-3,* and *qProt_Gm18-3* – carried the favorable alleles from 'AC X790P'.

117    Positive protein-related QTL alleles in different genetic backgrounds suggests that it may be possible to stack

118    favorable alleles to develop superior high-protein progeny.

119    Nine putative QTL – *qPro_Gm01-2* (R2 = 10.4%), *qPro-Gm04-4* (R2 = 13.7%), *qPro-Gm05-2* (R2 =

120    14.2%), *qPro_Gm06-1*(R2 = 21.9%), *qPro_Gm06-3* (R2 = 12.6%), *qPro_Gm08-2* (R2 = 12.3%), *qPro-Gm12-3* (R2

121    = 11.6%), *qPro-Gm12-4* (R2 = 12%), and *qPro_Gm13-4* (R2 = 11.6%) – identified in this study were previously

122    unreported (Table 2; 26]. Four of these QTL were identifiedin both mapping populations (Table 2).  The five QTL

123    associated with seed protein concentration that co-localized with previously identified protein-related QTL on

124    SoyBaseare listed in Table 2; Supplementary Table 6.

125

126    **QTL Associated with Additional Value-Added Traits**

127    Genomic regions harboring putative large-effect QTL associated with seed protein concentration were

128    evaluated for their associations with seed yield, sucrose concentration and seed weight (Table 3; Supplementary

129    Table 5). Of the fourteen protein-related QTL, eight QTL were co-localized with QTL associated with other traits.

130    Three protein-related QTL – *qPro_Gm01-2*, *qPro_Gm02-3*, and *qPro_Gm12-4* –were co-localized with QTL

131    associated with seed sucrose concentration (Table 3). The favorable alleles were inherited from opposing parental

132    sources for each of these genomic regions, which supports the significant negative relationship observed between

133    seed protein and sucrose concentration in this study. (Table 3; Fig. 3). The remaining five protein-related QTL were

134    associated with seed weight, with positive associations noted for three of these regions (Table 3; Fig. 3). Favourable

135    alleles were donated by each parental cultivar for all traits-of-interest. Protein-related QTL were not co-localized

136    with significant regions for seed yield, consistent with the non-significant relationship between seed protein

137    concentration and seed yield in both populations. SoyBase associated seven of our protein-related QTL with

138  previously identified QTL for seed weight (nine QTL), seed oil concentration (five QTL) and seed yield (two QTL)

139  (Supplementary Table 6; 26].

140

141  **Candidate Genes**

142  A list of candidate genes was compiled using the Glyma 2.0 Assembly of Williams 82 on SoyBase

143  (Wm82.a2.v1) according to their functionknowledge [26]. The number of genes in each QTL flanking region varied

144  from four to seventy-four. In the flanking region corresponding to *qPro_Gm13-4* (spanning 26 kb), five genes were

145  identified. These genes include Glyma.13G167800 and Glyma.13G167900, which are located 6 and 9 kb

146  downstream of the SNP peak (28246299) and are annotated as a ribosomal protein and a ribosome biogenesis

147  regulatory protein, respectively (Table 4). These genes have an indirect role in protein synthesis. Gene expression

148  data provided by Severin et al. [42] noted that Glyma.13G167800 is expressed in the seed from 10 to 21 day after

149  flowering (DAF). Glyma.13G167900 is also expressed in the seed albeit at a lower level compared to

150  Glyma.13G167800. Two candidate genes, Glyma.06G004500 and Glyma.06G001800, underlying *qPro_Gm06-*

151  *1*were identified. These genes, located in 74 kb upstream and 148 kb downstream of the QTL peak, respectively,

152  encode transmembrane amino acid transporter proteins and ribosomal family proteins and (Table 4). Previous

153  transcriptomic analyses noted increased expression of Glyma.06G004500 in the seed at 14 to 17, and 21 DAF [42].

154  Glyma.04G212500 and Glyma.04G214500 were identified under *qPro_Gm04-4*intervals. These genes are

155  associated with the cupin superfamily and ribosomal protein family, respectively (Table 4). The cupin superfamily is

156  involved in seed storage protein [43], while ribosomal protein family genes are associated with mRNA translation.

157  In addition, candidate gene Glyma.04212500 are located exactly in the SNP peak position, which support the role of

158  cupin associated with seed protein concentration. Glyma.06G113700, Glyma.06G116400, and Glyma.06G119700

159  were located in*qPro_Gm06-3*region (Table 4). Glyma.06G113700 encodes a potential structural constituent of 40S

160  ribosomal protein. Glyma.06G116400 and Glyma.06G119700 were associated with a transmembrane amino acid

161  transporter protein and an intracellular transport protein, respectively (Table 4).

162  Three candidate genes, Glyma.15G129800, Glyma.15G130000, and Glyma.15G134800, were identified

163  from *qPro_Gm15-3*which are involved in structural constituents of the ribosome (Table 4). Moreover,

164  Glyma.06G225600 andGlyma.06G225700, which were annotated as translation initiation factor proteins were

165  identified under *qPro_Gm06-6* intervals (Table 4). Glyma.02G220000 and Glyma.02G221500, which contribute to

166    the structural integrity of the ribosome and play a role in translation were locatedin *qPro_Gm02-3*region (Table 4).

167    Based on previous transcriptomic analyses, Glyma.02G220000 is expressed in the seed 14 to 17, 21, 25, 28 and 35

168    DAF [42].

169          Candidate genes were also postulated for sucrose- and seed weight-related QTL that co-localized with

170    protein-related regions. Four candidate genes were identified: Glyma.06G004400 and Glyma.06G007900, which

171    were located under *qPro_Gm06-1* and *qWt_Gm06-1 region*, and Glyma.15G133600 and Glyma.15G133800 that

172    were located under *qPro_Gm15-3* and *qWt_Gm15-4 region*. All four genes are involved in carbohydrate metabolism

173    (GO:0005975) (Table 5).

174

175    **Discussion**

176          Soy-based food manufacturers require specific physical and chemical characteristics of the soybean seed to

177    maintain their production practices. For example, optimal tofu production requires high concentrations of both

178    protein and sucrose in the soybean seed. However, protein and sucrose concentration have a negative relationship

179    [27, 44-47]. These significant negative relationships between seed protein concentration and other value-added traits

180    have been major deterrents to the development of competitive food-grade soybean cultivars through conventional

181    breeding methods[14-23, 48]. The identification of protein-related QTL that has no effect on sucrose or has a

182    positive impact on other value-added traits would be of major benefit.The relationship between seed protein

183    concentration, seed weight and yield in our study indicated that both current populations are desirable for the

184    selection of optimal protein concentration with competitive yield and large seed size. On the other hand, negative

185    relationship between seed protein and sucrose concentration indicated the selection for protein concentration may

186    occur at the expense of seed sucrose concentration (and vice versa). These relationships could be attributed to tightly

187    linked loci governing these traits separately, or to pleiotropic effects of specific loci [19].

188          Broad-sense heritability estimations in current study confirmed that a large proportion of the observed

189    phenotypic variation for seed protein concentration, seed sucrose concentration, and seed weight are attributed to

190    genotype. Therefore, phenotypic selection may be a successful tool to increase genetic gain for these traits. This is

191    consistent with previous studies, in which moderate to high heritability estimates have been reported for seed protein

192    concentration ($H^2 = 0.81$-$0.92$; 16,49], seed sucrose concentration ($H^2 = 0.46$-$0.86$;45,50] and seed weight ($H^2 =$

193    $0.73$-$0.89$; 49] across different genetic backgrounds and environments.

194    It is possible to 'stack' desirable QTL for multiple traits of interest using MAS, which allows breeders to

195    screen early generation material for optimal trait combinations. This approach has been utilized breeding programs,

196    especially for breeding disease resistance cultivars [51-53]. Maroof et al. [54] discussed the value of pyramiding

197    race-specific soybean mosaic virus resistance genes using MAS, which involved the curation of specific genetic

198    combinations for optimal multiple resistance. This approach increased the ability of the breeding program to select

199    homozygous plants with multiple resistance, as the epistatic interactions among disease resistance genes made the

200    phenotypic screening of disease reaction unreliable [54]. This strategy was also utilized by Jiang et al. [55] where

201    the pyramiding of positive alleles from different parental sources was shown to increase seed protein filling rate and

202    overall seed quality in soybean.

203    In this study, fourteen large-effect QTL associated with seed protein concentration were identified, with the

204    positive alleles derived fromeach of the parental sources. This may be attributed to the unique mapping populations

205    utilized in this study. Previous QTL studies have used mapping populations that were derived from exotic

206    germplasm or parental cultivars with large phenotypic differences for the desired trait-of-interest [40]. However,

207    many modern elite soybean cultivars already possess high protein concentrations (approximately 40%, dry basis)

208    and may be fixed for the large-effect QTL identified in diverse populations. In the current study, the utilization of

209    moderate and highprotein elite parental cultivars facilitated the identification of novel QTL that may have been

210    masked in other populations [49,56,57]. For instance, we did not detect any major QTL in chromosomes 15 and 20

211    that are frequently reported to be important genomic regions associated with seed protein content. Due mainly to

212    limited number of polymorphic markers between the parents, in this study, resulted in having two or more linkage

213    groups for most of the chromosomes and probably some missing regions. The elimination of these regions may have

214    also restricted the full scope of QTL interactions in these populations, and exaggerated the influence of the identified

215    QTL on the traits-of-interest [56,58,59].Additionally, many QTL mapping procedures have difficulty with the

216    identification of small and intermediate effect QTL. Thesesmall and intermediate QTL are primarily associated with

217    quantitative traits, such as seed protein concentration [60,61]. The Beavis effect suggests that estimates of

218    phenotypic variance may be greatly overestimated in smaller mapping populations (<1000 progeny; 60), which may

219    have further exaggerated the influence of the identified QTL in this study.

220    Recently, Hagely et al. [62] utilized direct molecular-assisted selection to improve the carbohydrate

221    composition of soybean seeds. A natural variant of the raffinose synthase 3 gene (*rs3 snp5*) was associated with an

222    ultra-low raffinose family oligosaccharide (UL RFO) carbohydrate profile, which improved the sucrose

223    concentration and available metabolized energy of the soybean meal [63,64]. The reduction in raffinose and

224    stachyose was attributed to a specific genetic combination – *rs2 W331 + rs3 snp5/rs3 snp 6* haplotype C – that

225    results from a defect in the RS3 gene. Molecular marker assays were developed to detect these variants, which

226    streamlined their introgression into elite soybean cultivars [62].

227          In an effort to further understand the underlying mechanisms of protein concentration in the soybean seed,

228    candidate genes were identified from the flanking regions of our protein-related QTL and screened for their

229    functional role in protein accumulation. In this study, 491 genes were identified and grouped using their biological

230    process and functional annotation in SoyBase (www.soybase.org;65]. Numerous putative candidate genes were

231    identified (Table 5) through GO annotation, including sixteen genes were associated with protein translation

232    processes (GO:0006412, GO:0015171, GO:0006413, GO:0042254, GO:0006886, AT6G61750, and PF01490).

233    Eight genes were found associated with carbohydrate metabolism (GO:0005975), three genes were associated with

234    lipid metabolism (GO:0006629), and the remaining genes were involved in signal transduction, transport,

235    biosynthetic processes, nucleic acid metabolism, photosynthesis, and numerous other functions. The significant

236    relationships between protein, oil and sucrose[27,44,46,47] support the role of genes associated with lipid and

237    carbohydrate metabolism, which were also identified in the flanking region of these protein-related QTL.

238          Transcriptome analysis data provided by Severin et al., [42] showed Glyma.13G167800 (ribosome

239    biogenesis), Glyma.13G167900 (ribosome biogenesis), Glyma.06G004500 (transmembrane amino acid transporter

240    protein) and Glyma.02G220000 (60S ribosomal protein) are expressed in the seed, which supports their role in

241    soybean seed protein accumulation. Glyma.04G212500 was associated with the cupin superfamily, which includes

242    the 11S (glycine) and 7S (ß-conglycinin) seed storage proteins. 11S and 7S seed storage proteins account for ~70%

243    of storage proteins within the soybean seed [43,66]. Therefore, Glyma.04G212500 may have a strong association

244    with seed protein accumulation in soybean. Zhang et al. [67] identified thirteen candidate genes with putative roles

245    in protein biosynthesis on Chromosome 15 and 20, with functional annotation of a structural constituent of

246    ribosome, 60S ribosomal protein, amino acid transmembrane transport, and translation initiation factor 3. These

247    annotations were also associated with seven candidate genes in our study, which strongly supports their role in

248    protein accumulation in our populations. Zhang et al. [67] also conducted gene expression analyses of ribosomal,

249    translation initiation factor 3 and amino acid transmembrane transport genes, which showed significant up-

250 regulation of expression in the high-protein parent during the reproductive growth stage in the pod. This is

251 consistent with their role in protein accumulation in soybean seeds [67]. Li et al. [68] also found a candidate gene in

252 the flanking region of a protein QTL on chromosome 9, which was annotated as an amino acid transporter gene. In

253 another study, the overexpression of one amino acid transporter gene in *Vicianarbonensis* and pearesulted in

254 significant increasesin seed protein concentration[69]. Further exploration of these candidate genes and their

255 possible variants would further our understanding of protein accumulation pathways in the soybean seed and may

256 lead to improved marker- or molecular-assisted breeding techniques for the improvement of soybean seed

257 composition traits.

258

259 **Conclusion**

260     In summary, nine of the protein-related QTL identified in this study were validated in both populations and

261 may be suitable for marker-assisted selection. Some of these QTL were collocated with other value-added traits and

262 can be used for simultaneous improvement of multiple traits. Their value will be dictated by the objective of the

263 breeding program. For example, *qPro_Gm06-1, qPro_Gm06-6, qPro_Gm08-2, and qPro_Gm15-3* were positively

264 associated with seed weight QTL. These QTL may be unsuitable for a natto breeding program, which would favour

265 smaller seed size. In this case, *qPro_Gm05-2* – a protein-related QTL inversely associated with seed weight – would

266 be preferable. A curated panel of multiple-trait QTL may allow breeders to screen early-generation germplasm for

267 the specific physical and chemical characteristics required by soy-food processors.

268     Future studies may look to consider the impact of protein biosynthesis, storage and metabolism on seed

269 protein concentration in soybean, as suggested by the postulated candidate gene functions noted in this study, to

270 foster a better understanding of protein accumulation pathways in the soybean seed. Breeders may also wish to dive

271 deeper and explore the potential variants of these candidate genes, and their role in plant metabolism. The QTL

272 identified this study can be used for marker-assisted selection and as a starting point for the discovery of variants in

273 the protein biosynthesis pathway.

274

275 **Abbreviations**

276 QTL: Quantitative trait loci

277 MAS: Marker-assisted selection

278    RIL: Recombinant inbred line

279    DAF: Day after flowering

280    MG: Maturity group

281    NNA: Nearest neighbour analysis

282    NIR: Near infrared reflectance

283    ANOVA: Analysis of variance

284    SNP: Single-nucleotide polymorphisms

285    CIM: Composite interval mapping

286    MQM: Multiple QTL mapping

287    LOD: Likelihood of odd

288    SMA: Single marker analysis

289    eFP: electronic Fluorescent Pictograph

290

291    **Methods**

292    **Mapping Populations**

293    Two populations of $F_4$-derived recombinant inbred lines (RILs) were used to identify putative quantitative

294    trait loci (QTL) for seed composition traits and yield. The first population (POPn_1) consisted of 190 RILs derived

295    from a cross between 'AC X790P' and 'S18-R6'. 'AC X790P' is a 2.2 relative maturity group (MG) cultivar

296    developed by Agriculture and Agri-Food Canada in Harrow, Ontario, with a high, stable seed protein concentration

297    (48.6%, dry weight basis; 70]. The seeds were obtained from The Harrow Research and Development Centre

298    (Harrow RDC) located in Harrow, Ontario. 'S18-R6' is a 1.8 MG commercial cultivar with a moderate seed protein

299    concentration (40.4%), developed by Syngenta Canada, Inc. in Arva, Ontario [71], where the seeds were obtained.

300    The second population (POPn_2) was comprised of 193 RILs from a cross between 'S23-T5' and 'AC

301    X790P'. 'S23-T5' is a high-yielding 2.3 MG elite cultivar with moderate seed protein (41.3%) developed by

302    Syngenta Seeds, Inc. in Owatonna, Minnesota [72]. The seeds were obtained form Syngenta Canada, Inc. in Arva,

303    Ontario. Parental cultivars were considered high yielding when compared to the historical yield for southwestern

304    Ontario [73]. Both RIL populations were developed at the University of Guelph, Ridgetown Campus.

305

**Experimental Design**

The RIL populations were grown in five environments across southwestern Ontario in 2015 and 2016: Chatham 2015 (CHA15), Merlin 2015 (MER15), Chatham 2016 (CHA16), Merlin 2016 (MER16) and Palmyra 2016 (PAL16). Field trials were planted using randomized complete block designs with two replications, in which the plot performance was adjusted for spatial variability through nearest neighbour analysis (NNA) using information from the immediate neighbouring plots in each of the five environments [74]. Plots consisted of five 4-m rows with 43-cm row spacing and were trimmed to 3.8-m in length following emergence. Plots were seeded at a rate of 69 seeds/m$^2$ or 500 seeds per plot. Trials were maintained using standard tillage and cultural practices, and the three center rows of each plot were harvested for seed yield estimation and post-harvest evaluations.

**Phenotypic Data Collection**

Seed protein and sucrose concentrations were determined for each harvested plot using near infrared reflectance (NIR) with a DA 7250 NIR analyzer (Perten Instruments Canada, Winnipeg, MB) with calibrations provided by Perten Instruments. NIR measurements were an average of three technical replications. Seed yield (tonnes ha$^{-1}$) and seed weight (grams per 100 seeds) were also recorded for each harvested plot.

**Statistical Analyses**

Statistical analyses were performed using SAS 9.4 (SAS Institute Inc., Cary, NC). An analysis of variance (ANOVA) was conducted and PROC MIXED was used to generate LSMEANS for each environment with 'genotype' as a fixed effect and 'block' as a random effect. PROC MIXED was also used to perform combined ANOVAs for seed weight, and protein and sucrose concentrations using the model:

$$Y_{ij} = \mu + \alpha_i + \beta_j + \alpha\beta_{ij} + \varepsilon_{ij}, \qquad j = 1, \dots, n; i = 1, \dots, k$$

where $Y_{ij}$ represented the trait of interest (seed protein accumulation, seed sucrose accumulation, seed weight or seed yield), $\alpha_i$ represents the 'genotype' effect, $\beta_j$ represents the 'environment' effect, $\alpha\beta_{ij}$ represents the 'genotype-by-environment' effect and $\varepsilon_{ij}$ represented the residual effect. 'Genotype', 'environment' and 'genotype-by-environment' were considered fixed effects and 'block(environment)' was considered a random effect. PROC CORR was used to examine the relationships between entry trait estimates.

**Genotypic Data Collection**

Young trifoliate leaf tissue was collected from the first replicate block of each population at the Palmyra 2016 location. Leaf tissue for each RIL was sampled from multiple plants in each plot and stored in 2mL screw cap tubes. The samples were freeze-dried for 72-hours using a Savant ModulyoDThermoquest (Savant Instruments, Holbrook, NY), and then stored at -80°C for future use. Genomic DNA was extracted from the freeze-dried tissue samples using a modified procedure from the Sigma GenElute™ DNA Extraction Kit (SIGMA®, Saint Louis, MO) methodology. DNA quality was verified using electrophoresis with 1% agarose gels, while quantity was verified using a Qubit® 2.0 fluorometer (Invitrogen, Carlsbad, CA).

DNA samples (30μl of 10ng μl$^{-1}$ DNA) were transferred to Plate-formeD'analysesGénomiques at Université Laval (Laval, Quebec, Canada) for genotyping-by-sequencing (GBS), using the Fast-GBS pipeline with the *Gmax_275_v2* reference genome [75]. The Fast-GBS pipeline identified 24,738 high-quality single-nucleotide polymorphisms (SNPs). Heterozygous SNPs were considered missing data. SNPs with >20% missing data or a minimum minor allele frequency less than 0.3 were discarded prior to imputation with Beagle [76].

**Linkage Map Construction and QTL Mapping**

JoinMap 5.0 software was used to construct genetic linkage maps for each population [77]. SNP markers with significant levels of segregation distortion that differed from the expected 1:1 ratio based on a chi-square test (α = 0.01) were removed from further analysis. Markers that segregated identically within the population were reduced to a single marker for linkage map construction. Markers were grouped into linkage groups within each chromosome using a minimum likelihood of odds (LOD) ≥ 3, and Kosambi's mapping function was used to calculate genetic distances. Thereafter, the genetic position of these markers was anchored on physical position.

Composite interval mapping (CIM) was performed for the traits of interest using the multiple QTL mapping (MQM) algorithm in MapQTL® 6 [78]. The empirical LOD threshold values were calculated through a permutation test with 1,000 iterations and a Type I error rate of 0.05. The automatic cofactor selection function was used to identify significant cofactors for MQM. Graphic representations of significant QTL were created using MapChart 2.32 [79].

Putative QTL regions associated with seed protein concentration were also screened for significant QTL associated with seed weight, seed yield and seed sucrose concentration. SoyBase was used to compare the putative

362 QTL to published genomic regions related to seed protein concentration [26]. Putative QTL were also confirmed in

363 the alternate population using single marker analysis (SMA) in SAS 9.4 (SAS Institute Inc., Cary, NC). PROC GLM

364 was used to identify significant single marker effects ($\alpha < 0.0001$) with LSMEAN estimates as the dependent

365 variable and SNP marker as the independent variable. The SNP positions from genotype-by-sequencing were used

366 to denote marker names in MQM and SMA.

367

368 **Candidate Gene Search**

369 The flanking markers of each QTL were chosen based on the LOD values surrounding each peak

370 marker.To ensure that the actual QTL was located within the range selected, the first marker below the LOD

371 threshold on each side of the QTL peak was selected as the flanking marker. For each of the protein-related QTL,

372 the regions between the flanking markers were used to identify candidate genes according to their function. A total

373 of 491 genes were extracted from the flanking regions using the SoyBase Soybean Genetic Map. The functional

374 annotation of each gene was identified from TAIR (www.arabidopsis.org/), GO (http://geneontology.org/), PFAM

375 (http://pfam.xfam.org/), and PANTHER (http://www.pantherdb.org/) through SoyBase (https://soybase.org/). This

376 functional knowledge used to reduce number of genes and identify putative candidate genes.

377 The Electronic Fluorescent Pictograph (eFP) browser for soybean (www.bar.utoronto.ca) was used to

378 generate additional information about the candidate genes, such as tissue- and developmental-stage dependent

379 expression (based on transcriptomic data from Severine et al. [42]). Pfam, a comprehensive collection of protein

380 domains and families, and NCBI were used to obtain additional information about candidate genes.

384 **Availability of data and materials**

385 All datasets will be freely available upon request.

386 **Authors' contributions**

387  ME designed the project. RW performed the experiments, collected and analyzed the data. ST mined the candidate

388  genes. RW and ST wrote the manuscript. ME and LL assisted to analysis and revised the manuscript.All authors

389  read and approved the final manuscript.

390

391  **Ethics declarations**

392  **Ethics approval and consent to participate**

393  Not applicable.

394  **Consent for publication**

395  All authors agreed to publish this manuscript.

396  **Competing interests**

397  The authors declare that they have no competing interests.

402

403  **References**

404  1.Wang HL, Hesseltine CW. Coagulation conditions in tofu processing. Process Biochem. 1983; 17:7–12.

405  2. Shen CF, DeMan L, Buzzell RI, DeMan JM. Yield and quality of tofu as affected by soybean and soymilk

406  characteristics: glucono-δ-lactone coagulant. J Food Sci. 1991; 56:109–12.

407  3. Schaefer MJ, Love J. Relationships between soybean components and tofu texture. J Food Qual. 1992; 15:53–66.

408  4. Cai T, Chang KC. Processing effect on soybean storage proteins and their relationship with tofu quality. J Agric

409  Food Chem. 1999; 47:720–7.

410    5. Poysa V, Woodrow L. Stability of soybean seed composition and its effect on soymilk and tofu yield and quality. Food Res Inst. 2002; 35:337–45.

412    6. Kim Y, Wicker L. Soybean cultivars impact quality and function of soymilk and tofu. J Sci Food Agric. 2005; 85:2514–8.

414    7. Stanojevic SP, Barac MB, Pesic MB, Vucelic-Radovic B V. Assessment of soy genotype and processing method on quality of soybean tofu. J Agric Food Chem. 2011; 59:7368–76.

416    8. Carver BF, Burton JW, Carter TE, Wilson RF. Response to environmental variation of soybean lines selected for altered unsaturated fatty acid composition. Crop Sci. 1986; 26:1176–81.

418    9. Vollmann J, Fritz CN, Wagentristl H, Ruckenbauer P. Environmental and genetic variation of soybean seed protein content under Central European growing conditions. J Sci. 2000; 1306:1300–6.

420    10. Sudaric A, Simic D, Vrataric M. Characterization of genotype by environment interactions in soybean breeding programmes of southeast Europe. Plant Breed. 2006; 125:191–4.

422    11. Akond M, Liu S, Boney M, Kantartzi SK, Meksem K, Bellaloui N, et al. Identification of quantitative trait loci (QTL) underlying protein, oil, and five major fatty acids contents in soybean. Am J Plant Sci. 2014; 5:158–67.

424    12. Chaudhary J, Patil GB, Sonah H, Deshmukh RK, Vuong TD, Valliyodan B, et al. Expanding omics resources for improvement of soybean seed composition traits. Front Plant Sci. 2015; 6:1021.

426    13. Ma Y, Kan G, Zhang X, Wang Y, Zhang W, Du H, et al. Quantitative trait loci (QTL) mapping for glycinin and beta-conglycinin contents in soybean (Glycine max L. Merr). J Agric Food Chem. 2016; 64:3473–83.

428    14. Shannon G, Wilcox JR, Probst AH. Estimated gains from selection for protein and yield in the F4 generation of six soybean populations. Crop Sci. 1972; 12:824–6.

430    15. Burton JW. Breeding soybeans for improved protein quantity and quality. In: Shibles R, editor. 3rd Soybean Research Conference, Ames, IA, 12-17 Aug 1984. Boulder, CO: Westview Press; 1985. p. 361–7.

432    16. Burton JW. Quantitative genetics: results relevant to soybean breeding. In: Wilcox JR, editor. Soybeans: improvement, production, and uses. Madison, WI: American Society of Agronomy; 1987. p. 211–47.

434    17. Wilcox JR, Cavins JF. Backcrossing high seed protein to a soybean cultivar. Crop Sci. 1995; 35:1036–41.

435    18. Helms TC, Orf JH. Protein, oil, and yield of soybean lines selected for increased protein. Crop Sci. 1998;

436    38:707–11.

437    19. Chung J, Babka H, Graef G, Staswick P, Lee D, Cregan P, et al. The seed protein, oil, and yield QTL on soybean

438    linkage group I. Crop Sci. 2003; 43:1053–67.

439    20. Cui Z, James AT, Mizazaki S, Wilson RF, Carter TE. Breeding specialty soybeans for traditional and new

440    soyfoods. In: Liu K, editor. Soybeans as functional foods and ingredients. Champaign, IL: AOCS Press; 2004. p.

441    264–322.

442    21. Yin X, Vyn TJ. Relationships of isoflavone, oil, and protein in seed with yield of soybean. Agron J. 2005;

443    97:1314–21.

444    22. Bandillo N, Jarquin D, Song Q, Nelson R, Cregan P, Specht J, et al. A population structure and genome-wide

445    association analysis on the USDA soybean germplasm collection. Plant Genome. 2015; 8:1–13.

446    23. Kim M, Schultz S, Nelson RL, Diers BW. Identification and fine mapping of a soybean seed protein QTL from

447    PI 407788A on Chromosome 15. Crop Sci. 2016; 56:219.

448    24. Zeng J, Chen P, Shi A, Wang D, Zhang B, Orazaly M, et al. Identification of quantitative trait loci for sucrose

449    content in soybean seed. Crop Sci. 2015; 54:554–64.

450    25. Hymowitz T, Collins FI, Panczner J, Walker WM. Relationship between the content of oil, protein, and sugar in

451    soybean seed. Agron J. 1972; 64:613–6.

452    26. SoyBase. SoyBase, the USDA-ARS soybean genetics and genomics database. 2019. https://soybase.org/.

453    Accessed 5 Apr 2019.

454    27. Patil G, Mian R, Vuong T, Pantalone V, Song Q, Chen P, et al. Molecular mapping and genomics of soybean

455    seed protein: a review and perspective for the future. TheorAppl Genet. 2017; 130:1975–91.

456    28. Qi Z, Zhang Z, Wang Z, Yu J, Qin H, Mao X, et al. Meta-analysis and transcriptome profiling reveal hub genes

457    for soybean seed storage composition during seed development. Plant Cell Environ. 2018; 41:2109–27.

458    29. Panthee DR, Pantalone V, West DR, Saxton AM, Sams CE. Quantitative trait loci for seed protein and oil

459    concentration, and seed size in soybean. Crop Sci. 2005; 45:2015–22.

460    30. Bernardo R. Breeding for quantitative traits in plants. 2nd edition. Woodbury, MN: Stemma Press; 2010.

461    31. Qi Z, Wu Q, Han X, Sun Y, Du X, Liu C, et al. Soybean oil content QTL mapping and integrating with mate-

462    analysis method for mining genes. Euphytica. 2011; 179:499–514.

463    32. Kadam S, Vuong TD, Qiu D, Meinhardt CG, Song L, Deshmukh R, et al. Genomic-assisted phylogenetic

464    analyses and marker development for next generation soybean cyst nematode resistance breeding. Plant Sci. 2015;

465    242:342–50.

466    33. Wang X, Jiang G, Song Q, Cregan P, Scott R, Zhang J, et al. Quantitative trait locus analysis of seed sulphur-

467    containing amino acids in two recombinant inbred line populations of soybean. Euphytica. 2015; 201:293–305.

468    34. Jofuku KD, Goldberg RB. Kunitz trypsin inhibitor genes are differentially expressed during the soybean life

469    cycle and in transformed tobacco plants. Plant Cell. 1989; 1:1079–93.

470    35. Walling L, Drews GN, Goldberg RB. Transcriptional and post-transcriptional regulation of soybean seed protein

471    mRNA levels. Proc Natl AcadSci USA. 1986; 83:2123–7.

472    36. Yeh KW, Chen JC, Lin MI, Chen YM, Lin CY. Functional activity of sporamin from sweet potato (Ipomoea

473    batatas Lam.): a tuber storage protein with trypsin inhibitory activity. Plant Mol Biol. 1997; 33:565–70.

474    37. Lotan T, Ohto MA, Yee KM, West MA, Lo R, Kwong RW, et al. Arabidopsis LEAFY COTYLEDON1 is

475    sufficient to induce embryo development in vegetative cells. Cell. 1998; 1:1195–205.

476    38. Soderman EM, Brocard IM, Lynch TJ, Finkelstein RR. Regulation and function of the Arabidopsis ABA-

477    insensitive4 gene in seed and abscisic acid response signaling networks. Plant Physiol. 2000; 124:1752–65.

478    39. Palomeque L, Li-Jun L, Li W, Hedges B, Cober ER, Rajcan I. QTL in mega-environments: II. Agronomic trait

479    QTL co-localized with seed yield QTL detected in a population derived from a cross of high-yielding adapted x

480    high-yielding exotic soybean lines. TheorAppl Genet. 2009; 119:429–36.

481    40. Hyten DL, Pantalone VR, Sams CE, Saxton AM, Landau-Ellis D, Stefaniak TR, et al. Seed quality QTL in a

482    prominent soybean population. TheorAppl Genet. 2004; 109:552–61.

483    41. Whaley R, Eskandari M. Genotypic main effect and genotype-by-environment interaction effect on seed protein

484    concentration and yield in food-grade soybeans (Glycine max (L.) Merrill). Euphytica. 2019; 215:33.

485    42. Severin AJ, Woody JL, Bolon YT, Joseph B, Diers BW, Farmer AD, et al. RNA-seq atlas of Glycine max: a

486    guide to the soybean transcriptome. BMC Plant Biol. 2010; 10:160.

487    43. Dunwell JM. Cupins: a new superfamily of functionally diverse proteins that include germins and plant storage

488    proteins. Biotechnol Genet Eng. 1998; 15:1–32.

489    44. Nichols D, Glover K, Carlson S, Specht J, Diers B. Fine mapping of a seed protein QTL on soybean linkage

490    group I and its correlated effects on agronomic traits. Crop Sci. 2006; 46:834.

491    45. Jaureguy LM, Chen P, Scaboo AM. Heritability and correlations among food-grade traits in soybean. Plant

492    Breed. 2011; 130:647–52.

493    46. Sonah H, O'Donoughue L, Cober E, Rajcan I, Belzile F. Identification of loci governing eight agronomic traits

494    using a GBS? WAS approach Valid by QTL Mapp soya bean. 2015; 12:211–21.

495    47. Patil G, Vuong TD, Kale S, Valliyodan B, Deshmukh R, Zhu C, et al. Dissecting genomic hotspots underlying

496    seed protein, oil, and sucrose content in an interspecific mapping population of soybean using high-density linkage

497    mapping. TheorAppl Genet. 2018; 16:1939–53.

498    48. Poysa V, Buzzell RI. AC X790P soybean. Can J Plant Sci. 2001; 81:447–8.

499    49. Eskandari M, Cober ER, Rajcan I. Genetic control of soybean seed oil: I. QTL and genes associated with seed

500    soil concentration in RIL populations derived from crossing moderately high-oil parents. TheorAppl Genet. 2013;

501    126:483–95.

502    50. Maughan PJ, Maroof MS, Buss GR. Identification of quantitative trait loci controlling sucrose concentration in

503    soybean (Glycine max). Mol Breed. 2000; 6:105–11.

504    51. Kelly JD, Afanador L, Haley SD. Pyramiding genes for resistance to bean common mosaic virus. Euphytica.

505    1995; 82:207–12.

506    52. Miklas PN, Delorme R, Stone V, Daly MJ, Stavely JR, Steadman JR, et al. Bacterial, fungal, and viral disease

507    resistance loci mapped in a recombinant inbred common bean population ('Dorado'/XAN 176). J Am SocHortic Sci.

508    2000; 176:476–81.

509    53. Singh S, Sidhu JS, Huang N, Vikal Y, Li Z, Brar DS, et al. Pyramiding three bacterial blight resistance genes

510    (xa5, xa13 and Xa21) using marker-assisted selection into indica rice cultivar PR106. TheorAppl Genet. 2001;

511    13:1011–5.

512    54. Maroof S, Jeong SC, Gunduz I, Tucker DM, Buss GR, Tolin SA. Pyramiding of soybean mosaic virus resistance

513    genes by marker-assisted selection. Crop Sci. 2008; 48:517–26.

514    55. Jiang Z, Han Y, Teng W, Zhang Z, Sun D, Li Y, et al. Identification of QTL underlying the filling rate of protein

515    at different developmental stages of soybean seed. Euphytica. 2010; 175:227–36.

516    56. Asins MJ. Present and future quantitative trait locus analysis in plant breeding. Plant Breed. 2002; 121:281–91.

517    57. Winter S, Shelp BJ, Anderson TR, Welacky TW, Rajcan I. QTL associated with horizontal resistance to soybean

518    cyst nematode in Glycine soja PI464925B. TheorAppl Genet. 2007; 114:461–72.

519    58. Hyne V, Kearsey MJ. QTL analysis: further uses of marker regression. TheorAppl Genet. 1995; 91:471–6.

520    59. Kearsey M, Farquhar A. QTL analysis in plants: where are we now? Heredity. 1998; 80:137–42.

521    60. Beavis WD. QTL analyses: Power, precision and accuracy. In: Paterson AH, editor. Molecular Dissection of

522    Complex Traits. Boca Raton, FL: CRC Press; 1998. p. 145–62.5. Bernardo R. Breeding for quantitative traits in

523    plants. 2nd edition. Woodbury, MN: Stemma Press; 2010.

524    61. Xu S. Theoretical basis of the Beavis effect. Genetics. 2003; 165:2226–59.

525    62. Hagely KB, Jo H, Kim JH, Hudson KA, Bilyeu K. Molecular-assisted breeding for improved carbohydrate

526    profiles in soybean seed. TheorAppl Genet. 2020; 133:1189–200.

527    63. Hagely KB, Palmquist D, Bilyeu KD. Classification of distinct seed carbohydrate profiles in soybean. J Agric

528    Food Chem. 2013; 61:1105–11.

529    64. Schillinger JA, Dierking EC, Bilyeu K. Soybeans having high germination rates and ultra-low raffinose and

530    stachyose content. 2013;8471.

531    65. Morales AM, O'Rourke JA, Scheider K, Bancroft T, Borem A, Nelson R, et al. Transcriptome analyses and

532    virus induced gene silencing identify genes in the RRpp4-mediated Asian soybean rust resistance pathway. Funct

533    Plant Biol. 2013; 4:1029–47.

534 66. Yaklich RW, Helm RM, Cockrell G, Herman EM. Analysis of the distribution of the major soybean seed

535 allergens in a core collection of Glycine max accessions. Crop Sci. 1999; 39:1444–7.

536 67. Zhang T, Wu T, Wang L, Jiang B, Zhen C, Yuan S, et al. A combined linkage and GWAS analysis identifies

537 QTLs linked to soybean seed protein and oil content. Int J Mol Sci. 2019; 20:5915.

538 68. Li X, Shao Z, Tian R, Zhang H, Du H, Kong Y, et al. Mining QTLs and candidate genes for seed protein and oil

539 contents across multiple environments and backgrounds in soybean. Mol Breed. 2019; 39:139.

540 69. Rolletschek H, Hosein F, Miranda M, Heim U, Gotz KP, Schlereth A, et al. Ectopic expression of an amino acid

541 transporter (VfAAP1) in seeds of Vicianarbonensis and pea increases storage proteins. Plant Physiol. 2005; 1:1236–

542 49.

543 70. Poysa V, Buzzell RI. AC X790P soybean. Can J Plant Sci. 2001; 81:447–8.

544 71. Canadian Food Inspection Agency (CFIA). Crop Reports: S18-R6. 2011.

545 http://www.inspection.gc.ca/english/plaveg/pbrpov/cropreport/soy/app00006462e.shtml. Accessed 13 Oct 2016.

546 72. Canadian Food Inspection Agency (CFIA). Crop Reports: S23-T5. 2011.

547 http://www.inspection.gc.ca/english/plaveg/pbrpov/cropreport/soy/app00007153e.shtml. Accessed 13 Oct 2016.

548 73. Ontario Ministry of Agriculture Food and Rural Affairs (OMAFRA). Provincial field crop production and

549 prices. 2016. http://www.omafra.gov.on.ca. Accessed 21 Dec 2017.

550 74. Whaley R, Eskandari M. Genotypic main effect and genotype-by-environment interaction effect on seed protein

551 concentration and yield in food-grade soybeans (Glycine max (L.) Merrill). Euphytica. 2019; 215:33.

552 75. Torkamaneh D, Laroche J, Bastien M, Abed A, Belzile F. Fast-GBS: a new pipeline for the efficient and highly

553 accurate calling of SNPs from genotype-by-sequencing data. BMC Bioinformatics. 2017; 18:5.

554 76. Browning BL, Browning SR. Genotype imputation with millions of reference samples. Am J Hum Genet. 2016;

555 98:116–26.

556 77. van Ooijen JW. JoinMap 4.0: Software for the calculation of genetic linkage maps in experimental populations.

557 2006.

558    78. van Ooijen JW. MapQTL® 6, Software for the mapping of quantitative trait loci in experimental populations of

559    diploid species. 2009.

560    79. Voorrips RE. MapChart: software for the graphical presentation of linkage maps and QTLs. 2002; 93:77–8.

561

562    **Figure Legends**

563    **Fig. 1** Relationship between average protein and sucrose concentrations (%, dry basis), seed weight (grams per 100

564    seeds) and seed yield (tonnes ha$^{-1}$) in RIL populations derived from (a) 'AC X790P' x 'S18-R6' and (b) 'AC

565    X790P' x 'S23-T5' examined under combined Ontario environments in 2015 and 2016. Trendlines depict the linear

566    regression between protein concentration and each trait. Pearson correlation coefficients are also noted (** denotes p

567    < 0.05; $^{ns}$ denotes a non-significant relationship

568

569    **Fig. 2** Distribution of LSMEANs and Pearson correlation coefficients among important seed quality traits in two

570    RIL populations examined under combined Ontario environments in 2015 and 2016: (a) 'AC X790P' x 'S18-R6'

571    and (b) 'AC X790P' x 'S23-T5'

572

573    **Fig. 3** Graphical representation of putative QTL identified using multiple QTL mapping (MQM) algorithms for seed

574    protein and sucrose concentrations, and seed weight in the two RIL populations: 'AC X790P' x 'S18-R6' and 'AC

575    X790P' x 'S23-T5'. Positive allele source is denoted by block pattern: 'AC X790P' is represented by a solid pattern,

576    while 'S18-R6' and 'S23-T5' are represented by a striped pattern. Traits of interest are denoted by colour: seed

577    protein concentration (red), seed sucrose concentration (navy) and seed weight (black)

578

579

580

581

582

583

584

585 Tables

586

587 **Table 2** Major putative QTL ($R^2 > 10.0\%$) associated with soybean seed protein concentration identified by multiple
588 QTL mapping (MQM) in the two RIL populations ('AC X790P x S18-R6' and 'AC X790P x S23-T5')evaluated in
589 five environments (CHA15, CHA16, MER15, MER16 and PAL16)
590

| QTL Name[z] | Chr. | POPn | Flanking Markers | | Size (cM) | LOD[y] | A[x] | $R^2$ (%) | Source | References[w] |
|---|---|---|---|---|---|---|---|---|---|---|
| *qPro_Gm01-2* | 1 | 2 | S01_42371693 | S01_42555910 | 2.19 | 4.56 | 0.4578 | 10.4 | S23-T5 | - |
| *qPro_Gm02-3* | 2 | 1 | S02_40793724 | S02_41072417 | 4.58 | 5.16 | 0.4115 | 10.4 | AC X790P | VAL[SMA]; 1,2 |
| *qPro_Gm04-3* | 4 | 2 | S04_44592458 | S04_45008840 | 1.64 | 5.25 | 0.4931 | 11.0 | S23-T5 | 2,3, 11 |
| *qPro_Gm04-4* | 4 | 1 | S04_48435528 | S04_49024162 | 14.21 | 6.03 | 0.3570 | 13.7 | AC X790P | - |
| ***qPro_Gm05-2*** | **5** | **1** | **S05_38330071** | **S05_38993543** | **12.31** | **6.80** | **0.4132** | **14.2** | **AC X790P** | **VAL[SMA]** |
| *qPro_Gm06-1* | 6 | 1 | S06_19074 | S06_699413 | 1.68 | 10.19 | 0.4408 | 21.9 | S18-R6 | - |
| ***qPro_Gm06-3*** | **6** | **1** | **S06_9128442** | **S06_11029737** | **19.08** | **5.51** | **0.3339** | **12.6** | **S18-R6** | **VAL[SMA]** |
| *qPro_Gm06-6* | 6 | 1 | S06_30639643 | S06_33589987 | 0.28 | 5.80 | 0.3046 | 13.2 | AC X790P | 2, 5, 6, 7 |
| ***qPro_Gm08-2*** | **8** | **1** | **S08_43864875** | **S08_43896183** | **2.25** | **5.38** | **0.3936** | **12.3** | **AC X790P** | **VAL[SMA]** |
| *qPro_Gm12-3* | 12 | 1 | S12_924424 | S12_1147989 | 11.46 | 6.45 | 0.4943 | 11.6 | S18-R6 | - |
| *qPro_Gm12-4* | 12 | 1 | S12_3518939 | S12_3666689 | 7.64 | 6.63 | 0.4757 | 12.0 | S18-R6 | - |
| ***qPro_Gm13-4*** | **13** | **2** | **S13_28227783** | **S13_28254683** | **4.46** | **8.54** | **2.2804** | **11.6** | **AC X790P** | **VAL[SMA]** |
| *qPro_Gm15-3* | 15 | 2 | S15_10218629 | S15_10877491 | 1.64 | 5.63 | 0.6925 | 11.5 | AC X790P | VAL[SMA]; 4,8,9,10 |
| *qPro_Gm18-4* | 18 | 1 | S18_52660341 | S18_53019901 | 18.54 | 4.50 | 0.2713 | 10.4 | AC X790P | VAL[SMA]; 2 |

[z]QTL for the same trait detected in all individual environments (CHA15, CHA16, MER15, MER16 and PAL16) and the combined environment (GMET) with the same or overlapping marker interval was designated as one QTL. QTL highlighted in bold are novel QTL and were validated in the other RIL population.
[y]LOD thresholds were calculated through a permutation test with 1,000 iterations and a Type I error rate of 0.001.
[x]Additive effects calculated as the absolute value of half the subtraction of the mean of genotypes with the 'S18-R6' ('POPn_1') or 'S23-T5' (POPn_2) allele (negative effect) from the mean of genotypes with the 'AC X790P' allele (positive allele).
[w]Indicating that the QTL was confirmed in the other RIL population through multiple QTL mapping (VAL[MQM]), single marker analysis (VAL[SMA]), and/or has been reported previously in the reference(s): 1. Qi et al. (2014); 2. Mao et al. (2013); 3. Stombaugh et al. (2004); 4. Lee et al. (1996) ;5. Rossi et al. (2013); 6. Liang et al. (2010); 7. Palomeque et al. (2009b); 8. Brummer et al. (1997); 9Warrington et al. 2015; 10. Fasoula et al., 2004; 11. Wang et al., 2014.

591
592

24

593
594
595
596
597
598
599

**Table 3** Putative QTL for additional food-grade traits of interest (seed yield, seed weight and sucrose concentration) associated with major seed protein concentration QTL identified by multiple QTL mapping (MQM) in a RIL population derived from 'AC X790P x S18-R6' and 'AC X790P x S23-T5' examined under combined Ontario environments from 2015 and 2016

| Protein QTL | QTL Name[z] | Chr. | POPn | Flanking Markers | | Size (cM) | LOD[y] | A[x] | R² (%) | Source | Relationship |
|---|---|---|---|---|---|---|---|---|---|---|---|
| qPro_Gm01-2 | qSuc_Gm01-2 | 1 | 2 | S01_42371693 | S01_42555910 | 2.19 | 6.67 | 0.1472 | 14.5 | AC X790P | Inverse |
| qPro_Gm02-3 | qSuc_Gm02-3 | 2 | 2 | S02_40716331 | S02_42411031 | 11.17 | 5.46 | 0.1993 | 10.7 | S23-T5 | Inverse |
| qPro_Gm05-2 | qWt_Gm5-2 | 5 | 2 | S05_38273700 | S05_38764985 | 1.94 | 3.98 | 1.2482 | 8.1 | S23-T5 | Inverse |
| qPro_Gm06-1 | qWt_Gm6-1 | 6 | 1 | S06_19074 | S06_798961 | 2.24 | 4.46 | 0.3927 | 10.3 | S18-R6 | Positive |
| qPro_Gm06-6 | qWt_Gm6-3 | 6 | 1 | S06_30639643 | S06_33589987 | 0.28 | 4.20 | 0.3754 | 9.4 | AC X790P | Positive |
| qPro_Gm08-2 | qWt_Gm8-2 | 8 | 1 | S08_43325761 | S08_43864912 | 17.39 | 4.29 | 0.5042 | 9.6 | AC X790P | Positive |
| qPro_Gm12-4 | qSuc_Gm12-1 | 12 | 1 | S12_3518939 | S12_3666689 | 7.64 | 5.49 | 0.1495 | 12.4 | AC X790P | Inverse |
| qPro_Gm15-3 | qWt_Gm15-4 | 15 | 2 | S15_10731054 | S15_11188445 | 3.33 | 2.78 | 0.8428 | 5.3 | AC X790P | Positive |

[z]QTL for the same trait detected in all individual environments (CHA15, CHA16, MER15, MER16 and PAL16) and the combined environment (GMET) with the same or overlapping marker interval was designated as one QTL.
[y]LOD thresholds were calculated through a permutation test with 1,000 iterations and a Type I error rate of 0.001.
[x]Additive effects calculated as the absolute value of half the subtraction of the mean of genotypes with the 'S18-R6' ('POPn_1') or 'S23-T5' (POPn_2) allele (negative effect) from the mean of genotypes with the 'AC X790P' allele (positive allele).

600
601
602
603
604
605
606

**Table 4** Major putative QTL (R² > 10.0%) and candidate genes identified in confidence intervals of QTL associated with soybean seed protein concentration in the two RIL populations ('AC X790P x S18-R6' and 'AC X790P x S23-T5')

| QTL Name[z] | Chr. | Flanking Markers | | | Candidate ID | Annotation | Type | Description | Position |
|---|---|---|---|---|---|---|---|---|---|
| *qPro_Gm02-3* | 2 | S02_40793724 | - | S02_41072417 | Glyma.02g220000 | GO:0006412 | GO-bp | 60S Ribosomal protein L16p/L10e | 40794106..40795066 |
| | | | | | Glyma.02g221500 | GO:0006412 | GO-bp | 30S Ribosomal protein S2 | 40921208..40921756 |
| *qPro_Gm04-4* | 4 | S04_48435528 | - | S04_49024162 | Glyma.04g212500 | AT5G61750 | AT | Cupin | 48435108..48435965 |
| | | | | | Glyma.04g214500 | GO:0006412 | GO-bp | Ribosomal protein L17 family protein | |
| *qPro_Gm06-1* | 6 | S06_19074 | - | S06_699413 | Glyma.06g004500 | GO:0015171 | GO-mf | Transmembrane amino acid transporter protein | 393722..398436 |
| | | | | | Glyma.06g001800 | GO:0006412 | GO-bp | Ribosomal protein L3 family | 171462..172334 |

25

| Protein QTL | Chr | Flanking Markers | | Candidate ID | Annotation | | Description | Position |
|---|---|---|---|---|---|---|---|---|
| | | 0 | | | | 2 | protein/Translation protein | |
| qPro_Gm06-3 | 6 | S06_9128442 | - | S06_11029737 | Glyma.06g113700 | GO:0006412 | GO-bp | 40S ribosomal protein S3a-like | 9225152..9227191 |
| | | | | | Glyma.06g116400 | PF01490 | PFAM | Transmembrane amino acid transporter protein | 9472699..9476835 |
| | | | | | Glyma.06g119700 | GO:0006886 | GO-bp | Intracellular protein transport | 9737256..9743653 |
| qPro_Gm06-6 | 6 | S06_30639643 | - | S06_33589987 | Glyma.06g225600 | GO:0006413 | GO-bp | Translation initiation | 31131372..31133932 |
| | | | | | Glyma.06g225700 | GO:0006412 | GO-bp | Translation initiation factor eIF-4F | 31209402..31216702 |
| qPro_Gm13-4 | 13 | S13_28227783 | - | S13_28254683 | Glyma.13g167800 | GO:0042254 | GO-bp | Ribosome biogenesis | 28237788..28239022 |
| | | | | | Glyma.13g167900 | GO:0042254 | GO-bp | Ribosome biogenesis regulatory protein | 28240381..28243803 |
| qPro_Gm15-3 | 15 | S15_10218629 | - | S15_10877491 | Glyma.15g129800 | GO:0006412 | GO-bp | Ribosomal protein S27a/Ubiquitin family | 10430457..10431571 |
| | | | | | Glyma.15g130000 | GO:0006412 | GO-bp | Structural constituent of ribosome | 10439067..10440332 |
| | | | | | Glyma.15g134800 | GO:0006412 | GO-bp | Ribosomal protein L7/L12 C-terminal domain | 10831146..10833232 |

[z]QTL for the same trait detected in all individual environments (CHA15, CHA16, MER15, MER16 and PAL16) and the combined environment (GMET) with the same or overlapping marker interval was designated as one QTL.

**Table 5** Major putative QTL ($R^2 > 10.0\%$) and candidate genes identified in confidence intervals of QTL associated with soybean seed protein concentration which co-located with seed weight or sucrose concentration in the two RIL populations ('AC X790P x S18-R6' and 'AC X790P x S23-T5')

| Protein QTL | QTL Name | Chr. | Flanking Markers | | Candidate ID | Annotation | Description | Position |
|---|---|---|---|---|---|---|---|---|
| qPro_Gm06-1 | qWt_Gm6-1 | 6 | S06_19074 | - | S06_798961 | Glyma.06g004400 | GO:0005975 | Carbohydrate metabolism | 380973..384365 |
| | | | | | | Glyma.06g007900 | GO:0005975 | Carbohydrate metabolism | 613002..614426 |
| qPro_Gm15-3 | qWt_Gm15-4 | 15 | S15_10731054 | - | S15_11188445 | Glyma.15g133600 | GO:0005975 | Carbohydrate metabolism | 10739528..10743270 |
| | | | | | | Glyma.15g133800 | GO:0005975 | Carbohydrate metabolism | 10754838..10756823 |