# GigaScience

## GALLO: An R package for Genomic Annotation and integration of multiple data source in livestock for positional candidate LOci

### --Manuscript Draft--

| Manuscript Number: | GIGA-D-20-00265R2 |
|---|---|
| Full Title: | GALLO: An R package for Genomic Annotation and integration of multiple data source in livestock for positional candidate LOci |
| Article Type: | Technical Note |

| Abstract: | The development of high-throughput sequencing and genotyping methodologies and precision livestock farming allowed the identification of thousands of genomic regions associated with several complex traits. The integration of multiple sources of biological information is a crucial step to better understand patterns regulating the development of complex traits. Genomic Annotation in Livestock for positional candidate LOci (GALLO) is an R package, for the accurate annotation of genes and quantitative trati loci (QTLs) located in regions identified in the most common genomic analyses performed in livestock, such as Genome-Wide Association Studies and transcriptomics using RNA-Sequencing. Moreover, GALLO allows the graphical visualization of gene and QTL annotation results, data comparison among different grouping factors (e.g., methods, breeds, tissues, statistical models, studies, etc.), and QTL enrichment in different livestock species including cattle, pigs, sheep, chicken, etc. Consequently, GALLO is a useful package for annotation, identification of hidden patterns across datasets, datamining of previous reported associations, as well as the efficient scrutinization of the genetic architecture of complex traits in livestock. |
|---|---|

| Corresponding Author: | Pablo Augusto de Souza Fonseca<br>University of Guelph<br>Guelph, ON CANADA |
|---|---|
| Corresponding Author Secondary Information: | |
| Corresponding Author's Institution: | University of Guelph |
| Corresponding Author's Secondary Institution: | |
| First Author: | Pablo Augusto de Souza Fonseca |
| First Author Secondary Information: | |
| Order of Authors: | Pablo Augusto de Souza Fonseca |
| | Aroa Suárez-Vega |
| | Gabiele Marras |
| | Angela Cánovas |
| Order of Authors Secondary Information: | |

| Response to Reviewers: | Reviewer reports:<br>Reviewer #1: The authors have clarified the questions I had and included the clarification appropriately in the manuscript.<br><br>L298-302: this should be in the discussion not in the conclusion. One reviewer |
|---|---|

suggested to include more species than just livestock and the authors responded that it is possible as long as the species has a resource list like the AnimalQTL database. "We acknowledge this comment in the revised version of the manuscript and have included a sentence highlighting the applicability to other species (Lines 298-302)." This sentence is added to the conclusion, however, I urge discussing it in the discussion. It is a point of discussion and not a conclusion of the core manuscript. (If properly discussed it may be included in the conclusion.)

Answer: Thank you for your suggestion. We edited the current version of the manuscript and the information about applicability to other species was moved to the discussion section.

Some more textual errors arose in the newly written sections:
L277: remove one 'the'

Answer: Done.

L278: change 'sama' into 'same'

Answer: Done.

L282: change 'easy to be handle' into easy to handle of easy to be handled

Answer: Done.

L283: change 'have' into 'has'

Answer: Done.


Just a note for future revisions: For this comment & answer below (Reviewer 1) I didn't find the sections on the lines indicated, but elsewhere (141-153 & 277-285). Please make sure you refer to the correct line-numbers in the future to accommodate the reviewers.

The authors indicated that the R package is similar to BiomaRt, and gave performance differences in term of execution time of comparable commands. BiomaRt is a renowned package and was faster. It would be nice if the authors can indicate what benefits GALLO has over BiomaRt. Why was this package needed (e.g. what did you miss in biomaRt)?
Also it may be worthwhile to explicitly indicate why R is the appropriate language for this package. There are thing mentioned scattered over the paper, e.g. like visuals and no need for intermediate output files, please summarize them somewhere.

Answer: Thank you for the comment. The comparison between GALLO and other available tools is better discussed on lines 241-253 and 468-476 of the revised version of the manuscript.


Answer: Thank you very much for the comment. In the next submissions the authors will be awarded about this issue.

| Additional Information: | |
|---|---|
| Question | Response |
| Are you submitting this manuscript to a special series or article collection? | No |
| Experimental design and statistics | Yes |
| Full details of the experimental design and | |

| | |
|---|---|
| statistical methods used should be given in the Methods section, as detailed in our Minimum Standards Reporting Checklist. Information essential to interpreting the data presented should be made available in the figure legends.<br><br>Have you included all the information requested in your manuscript? | |
| **Resources**<br><br>A description of all resources used, including antibodies, cell lines, animals and software tools, with enough information to allow them to be uniquely identified, should be included in the Methods section. Authors are strongly encouraged to cite Research Resource Identifiers (RRIDs) for antibodies, model organisms and tools, where possible.<br><br>Have you included the information requested as detailed in our Minimum Standards Reporting Checklist? | Yes |
| **Availability of data and materials**<br><br>All datasets and code on which the conclusions of the paper rely must be either included in your submission or deposited in publicly available repositories (where available and ethically appropriate), referencing such data using a unique identifier in the references and in the "Availability of Data and Materials" section of your manuscript.<br><br>Have you have met the above requirement as detailed in our Minimum Standards Reporting Checklist? | Yes |

1    **GALLO: An R package for Genomic Annotation and integration of multiple**

2    **data sources in Livestock for positional candidate LOci**

3    Pablo A.S. Fonseca[1][*], Aroa Suárez-Vega[1], Gabriele Marras[1,2], and Ángela Cánovas[1][*]

4    [1]University of Guelph, Department of Animal Biosciences, Centre for Genetic Improvement of

5    Livestock, Guelph, N1G 2W1, Ontario, Canada.

6    [2]The Semex Alliance, Guelph N1G 3Z2, Ontario, Canada

7    Contact:

8    PASF: pfonseca@uoguelph.ca - https://orcid.org/0000-0002-6917-7475

9    ASV: asuarezv@uoguelph.ca - https://orcid.org/0000-0002-7726-4288

10   GM: gmarras@uoguelph.ca - https://orcid.org/0000-0001-5115-370X

11   AC: acanovas@uoguelph.ca -https://orcid.org/0000-0002-0036-0757

12   [*] Corresponding author

13

14    **Abstract**

15    **Background**: The development of high-throughput sequencing and genotyping methodologies

16    allowed the identification of thousands of genomic regions associated with several complex traits.

17    The integration of multiple sources of biological information is a crucial step required to better

18    understand patterns regulating the development of these traits. **Findings:** Genomic Annotation in

19    Livestock for positional candidate LOci (GALLO) is an R package developed for the accurate

20    annotation of genes and quantitative trait loci (QTLs) located in regions identified in common

21    genomic analyses performed in livestock, such as Genome-Wide Association Studies and

22    transcriptomics using RNA-Sequencing. Moreover, GALLO allows the graphical visualization of

23    gene and QTL annotation results, data comparison among different grouping factors (e.g.,

24    methods, breeds, tissues, statistical models, studies, etc.), and QTL enrichment in different

25    livestock species including cattle, pigs, sheep, and chickens, etc. **Conclusions:** Consequently,

26    GALLO is a useful package for the annotation, identification of hidden patterns across datasets,

27    datamining previously reported associations, as well as the efficient scrutinization of the genetic

28    architecture of complex traits in livestock.

29    **Keywords:** Multi-omics integration; QTL annotation; Gene annotation; Datamining; QTL

30    enrichment analysis; Livestock

31

32

33

34

**Background**

36   The identification of quantitative trait loci (QTLs), genomic regions linked to complex traits

37   through association tests using genetic markers and phenotypic traits, is a crucial step in the

38   improvement of genomic selection and economic profitability in livestock [1–4]. The development

39   of high-throughput methodologies (e.g., Genome-Wide Association Studies, Transcriptomics,

40   Metabolomics, Proteomics, etc.) for the study of the genetic architecture of complex traits allows

41   for the identification of potential candidate genes associated with economically relevant traits in

42   livestock. Taken together, these technologies can substantially improve the accuracy of detection

43   of candidate regions associated with economically important traits across the genome in livestock

44   species [5]. Consequently, the number of QTLs identified across the genome in livestock species

45   increased substantially in the last few years. As of October 2020, the Animal QTLdb can retrieve

46   information about QTLs previously identified in cattle (159,844), chickens (12,508), horses

47   (2,451), pigs (30,871), rainbow trout (584) and sheep (3,411) [6]. The proper integration of results

48   obtained from different methodologies and technologies is a crucial step for the accurate

49   identification of the biological processes regulating complex traits as well as, the identification of

50   potential functional candidate genes for each trait or those shared among traits [5,7–9]. The

51   integration of both structural and functional data can help scrutinize the genetic architecture of

52   economically relevant traits, and consequently, help to better understand complex biological

53   patterns regulating the expression of these traits, such as pleiotropic effects, epistasis, and genetic

54   hitchhiking, among others.

55   Despite the potential to improve the identification of functional candidate genes and/or QTLs

56   through the integration of multiple data sources, the current process poses limitations in the

57   pipelines and algorithms implemented in the tools available for livestock. Currently, there are

58  several tools that implement functions for gene (i.e., Biomart and BEDTools) and QTL annotation

59  (Animal QTLdb) [6,10,11]. However, these tools have limitations regarding the automatization

60  process to analyze results from multiple candidate regions (Biomart web application and the R

61  package and Animal QTLdb) or for the visualization of the results. Moreover, although

62  automatization is possible, there is no direct link between the candidate regions and/or markers

63  with the annotated genes and QTLs. Consequently, this gap forces the user to back solve the

64  overlap between the input and output files in order to perform the proper association between the

65  candidate region and/or markers and the annotated genes and/or positional co-localized QTLs. In

66  addition, there is still a need for customized QTL enrichment analyses in the available software

67  and databases. The Genomic Annotation in Livestock for positional candidate LOci (GALLO) is

68  an R package designed to provide an automatized and a straightforward environment for gene and

69  QTL annotation in multiple candidate regions, as well as the integration of data from multiple

70  sources. Additionally, the QTL enrichment analysis can be performed directly by GALLO using

71  the output obtained from the QTL annotation step. GALLO also provides a set of functions for

72  graphical visualization of the annotation, comparison, integration and QTL enrichment results. In

73  this context, the GALLO package was developed as an alternative tool: 1) to allow the integration

74  and simultaneous annotation of multiple datasets for genes and QTLs; 2) to provide graphical

75  visualization tools to visually integrate the annotation and similarity against datasets; 3) to perform

76  QTL enrichment analysis for the positional candidate genomic regions and/or markers associated

77  with economically relevant traits in livestock.

78  **Implementation**

79  The GALLO package was written in the R language [12]. The stable release is available as an R

80  package on CRAN (https://cran.r-project.org/web/packages/GALLO/index.html). The code was

22

81  extensively tested with several datasets from different sources and methodologies and reviewed to

82  ensure it meets the packages high quality standards. Additionally, the vignettes were created to be

83  comprehensive and to present practical examples in order to provide a user-friendly tutorial.

84  The GALLO package provides a useful set of functions that gives a straightforward approach to

85  data integration, comparison, gene and QTL annotation, and visualization of several data sources

86  and methodologies, such as variants from genome-wide association study (GWAS), RNA-

87  Sequencing, whole-genome sequencing, etc. (Figure 1 and Table 1). The main advantage to

88  perform an automated analysis from multiple datasets is the ability to handle the output using

89  different subsets (traits, populations, models, etc.) in the same environment without generating

90  multiple intermediate output files.

91  *Case study – Candidate regions for scrotal circumference and fertility in cattle*

92  The dataset used to present the basic usage and advantages of the GALLO package is composed

93  by the markers significantly associated with scrotal circumference in the Canchim breed [13] and

94  noncompensatory fertility in Holstein cattle [14]. These two studies were previously analyzed

95  together in a systematic review regarding male fertility in cattle [8]. Therefore, the data used herein

96  comprises a multi-study and multi-breed analysis. These candidate markers (527 single nucleotide

97  polymorphisms (SNPs)) are available in Supplementary Table 1. In addition to the candidate

98  markers, we presented as Supplementary Files 1 and 2, the annotation gff file containing the QTL

99  database     information     for     cattle     (obtained     from     the     Animal     QTLdb;

100  https://www.animalgenome.org/cgi-bin/QTLdb/BT/download?file=gffUMD_3.1) and the gtf file

101  containing    the    genes    annotated    in    the    cattle    genome    obtained    from    Ensembl

102  (ftp://ftp.ensembl.org/pub/release-94/gtf/bos_taurus/). The genomic coordinates of both files were

103    based on the bovine reference genome version UMD 3.1 due to the original coordinates used to

104    report the location of the candidate markers in the original studies. Here, the analysis performed

105    follows the same logical order to the one presented in the GALLO vignette

106    (https://rpubs.com/pablo_bio/GALLO_vignette). However, the dataset used in the user practical

107    tutorial is a subset of the data presented here, aiming to reduce the computational demand for the

108    user. The script with all the commands used to perform the analysis presented here are available

109    in Supplementary File 3. All the tests were performed using a desktop with a processor Intel Core

110    i5 2.4 GHz with 8 Gb of RAM memory.

111    *Importing datasets and annotating genes and QTLs around candidate markers*

112    The first step in the pipeline consists of importing the databases which will be used for the analysis

113    with the *import_gff_gtf()* function. In our specific example, we imported both cattle gene

114    annotation (gtf) and QTL (gff) databases. The *import_gff_gtf()* function receives the database file

115    (db_file) and the file type (*file_type= "gff" or "gtf"*) as arguments and creates a dataframe with

116    the respective information from each file. The system time taken to import the gtf and gff files

117    were 0.045 and 0.311 seconds, respectively, indicating an efficient importing process. The file

118    containing the candidate markers can be imported using any available function in the R

119    environment such as *read.table()* and *read.csv()*.

120    The main function of GALLO, *find_genes_qtls_around_markers()*, performs the annotation of

121    genes and/or co-localized QTLs within or nearby candidate markers or genomic regions (using the

122    user's defined interval/window). This function uses the information provided in the .gtf file (for

123    gene annotation) or .gff (for QTL annotation) to retrieve the requested information. The output

124    combines the information available in the input file provided by the user with the information

125    available for the genes and QTLs mapped in the candidate genomic regions. For example, for an

126    input file composed of three genomic coordinates where four genes are annotated in each of the

127    intervals determined by the user, the output file of *find_genes_qtls_around_markers()* will contain

128    12 rows.  The minimum information necessary for the gene and QTL annotation procedures is a

129    data frame with two columns containing the chromosome (CHR) and position in base pairs (BP)

130    in the case of the candidate SNPs input file. In the case of the candidate haplotypes, windows,

131    copy number variations (CNVs) or candidate regions; the input file is composed by three columns

132    corresponding to the chromosome (CHR), the start position in base pairs (BP1) and the end

133    position in base pairs (BP2). Data examples for the candidate markers and windows input files can

134    be obtained using the *data("QTLmarkers")* and *data("QTLwindows")* commands in R.

135    Additionally, examples of QTL and gene annotation results are accessible through the

136    *data("gtfGenes")* and *data("gffQTLs")* commands, respectively. These outputs can be easily

137    handled by summary functions in R, such as *table()*, to obtain information such as the total number

138    of genes and QTLs, the number of genes and QTLs annotated per variants, etc. The gene annotation

139    process was compared with the *getBM()* function from the biomaRt package.  The gene annotation

140    process on GALLO needed 0.424 seconds to completely annotate the genes in a 200 Kb interval

141    (upstream and downstream) from candidate markers, while the biomaRt function required 0.019

142    seconds. The QTL annotation on GALLO was compared with the Bedtools -wao -C command,

143    resulting in 0.851 and 0.12 seconds required for each approach, respectively. It is important to

144    highlight that for both gene and QTL annotation using biomaRt and Bedtools, respectively, a

145    posterior processing of the output file is required in order to match the candidate markers and the

146    genes and QTLs mapped within the candidate intervals. On the other hand, the output file from

147    *find_genes_qtls_around_markers()* function was designed to allow this match in an intuitive way,

148    combining the rows of both candidate markers file and database files (gff and gtf). Additionally,

149    GALLO allows the user to perform both annotations for genes and QTLs with a single software

150    and programming language. Consequently, GALLO obtains a more elaborate and informative

151    output without substantially compromising the computational demand required for the analysis.

152    The output files obtained in the gene and QTL annotation are available on Supplementary Tables

153    2 and 3, respectively.


154    *Comparing and visualizing the overlapping of genes and QTLs annotated within the candidate*

155    *regions*


156    The output file generated by the *find_genes_qtls_around_markers()* function can be used as an

157    input file for the other set of GALLO functions. An advantage from the output of

158    *find_genes_qtls_around_markers()* function is that any additional information present in the input

159    file will be retained in the output file. Consequently, this information can be used to compare the

160    retrieved information between groups of population, methodologies, statistical models, etc. For

161    example, the functions *overlapping_among_groups()* and *plot_overlapping()* can be used to create

162    matrices with the overlapping values among groups and to visualize this overlap. Figure 2 shows

163    the genes and QTLs overlapping between the positional markers obtained in the two selected

164    studies from the dataset of markers analyzed, Feugang et al. (2009) [14] and Buzanskas et al.

165    (2017) [13]. It is important to highlight that the overlapping matrix informing the percentage of

166    shared records is not symmetrical. The percentage of genes from study A shared with the study B,

167    and vice-versa, are calculated as a function of the total number of genes in A or B, respectively.

168    Briefly, this matrix is not symmetrical because GALLO calculates the percentage of records shared

169    as a function of the total number of records for each group. For example, groups A and B shared

170    5 records, where group A has 10 records in total and group B has 5 records. Consequently, the

171    percentage of shared records in A is 50% while the percentage of shared genes in B is 100%. In

172    the current example, it is possible to note that only a small percentage of the positional candidate

173    genes were shared between the studies. However, the analyses of overlapping QTLs (using the

174    trait name as reference ID) indicated a higher similarity between the studies, 46% of the QTLs

175    annotated in the candidate regions from Feugang et al. (2010) [14] were also present in Buzanskas

176    et al. (2017) [13] and 93% of the QTLs annotated in the candidate regions from Buzanskas et al.

177    (2017) were also present in Feugang et al. (2010) [13,14].


178    *Understanding the QTL context of the candidate regions*


179    A more precise investigation of the QTL representativeness and diversity can help to better

180    understand the genomic context of the candidate regions. The recurrent association of particular

181    genomic regions with multiple traits might suggest the presence of complex genetic mechanisms

182    regulating that region, such as pleiotropy, epistasis, hitchhiking effect, among others [15,16]. The

183    *plot_qtl_info()* function from GALLO allows for the graphical visualization of the summary of

184    QTL types and traits annotated. The percentage of each QTL type for cattle (i.e., milk, meat and

185    carcass, health, production, reproduction and exterior) annotated within the candidate regions is

186    presented in a pie plot through the use of the argument *qtl_plot="qtl_type"*, while the percentage

187    of each trait associated with a specific QTL type can be plotted using the argument

188    *qtl_plot="qtl_name"* and informing the additional argument *qtl_class* (that must receive the name

189    of the QTL class to be plotted). Figure 3 shows that for Feugang et al. (2009) [14] the two most

190    frequent QTL types were Milk (50.42%) and Reproduction (16.97%), while for Buzanskas et al.

191    (2017) [13] the most frequent QTL types were Reproduction (87.06%) and Meat and Carcass

192    (5.03%). An in-depth analyses can be performed for each QTL type in order to observe the

193    frequency of each trait associated with a specific QTL type. The most frequent traits related with

194    Reproduction QTLs were calving ease (>3%) and scrotal circumference (>60%) for Feugang et al.

195    (2009) and Buzanskas et al. (2017) [13,14], respectively (Figure 3). The comparison between the

196    frequency of traits related with Reproduction QTLs annotated in Feugang et al. (2009) and

197    Buzanskas et al. (2017) [13,14] indicated that among the top 10 most frequent QTLs, calving ease,

198    inhibin levels, stillbirth, interval to first estrus after calving, and birth index were shared between

199    the studies. The combined analysis (not filtering by study) indicated that the Reproduction and

200    Milk QTL types were the two most frequent classes with 76.99% and 10.62% of all QTL types,

201    respectively. In addition, scrotal circumference, inhibin level and calving ease were the most

202    frequent Reproduction QTL related traits in the combined analysis.

203    *QTL enrichment analysis*

204    In some cases, the biases produced with more research in certain areas/traits of higher relevance

205    to animal production (such as milk production related traits in the QTL database for cattle) may

206    result in a larger proportion of records for these traits in the QTL database. Consequently, the

207    simple investigation of the proportion of each QTL type might not be totally useful. The GALLO

208    package allows the user to perform a QTL enrichment analysis to test the significance of the QTL

209    representativeness. The QTL enrichment analysis function in the GALLO package is based on a

210    hypergeometric test approach, where the number of QTLs annotated within the candidate regions

211    for each QTL type or trait, is compared with the observed number of QTLs in the reference

212    database. Briefly, using an enrichment for individual traits in a chromosome-wide approach as an

213    example, the number of traits per chromosome annotated within the candidate regions and the total

214    number of each individual trait in the QTL database are computed. Subsequently, this information

215    is integrated into a hypergeometric test in order to estimate if the number of observed records, for

216    a specific trait, in a chromosome is larger than expected by chance.  The *qtl_enrich()* function

217    allows the user to perform the QTL enrichment analysis for both QTL types and traits (*qtl_type=*

218    *"QTL_type"* or *"Name"*), for the whole genome or chromosome-wide (*enrich_type= "genome"*

219    or *"chromosome"*) and for all the annotated chromosomes or a subset (*chr.subset= NULL* or the

220    object with the subset of chromosomes). The use of a chromosome-wide enrichment analysis

221    might help to detect specific regions across the genome with a high number of QTLs for a specific

222    trait, i.e. BTA14 in cattle for milk production [17]. A total of 161 unique pairs of traits and

223    chromosomes were tested for the enrichment using the annotated QTLs from both studies. The

224    system time required to perform the enrichment analysis was 5.32 seconds, suggesting efficient

225    processing. The top 10 enriched QTLs (False Discovery Rate (FDR) < 0.05) for the combined

226    analysis is shown in Table 2 and the enrichment results for all the annotated QTLs is shown in

227    Supplementary Table 4.   Additionally, GALLO also allows the user to obtain a graphical

228    visualization, in a bubble plot, of the enrichment results using the *QTLenrich_plot()* function. This

229    function receives the enriched table obtained from *qtl_enrich()*, the name of the column with the

230    trait names to be plotted and the name of the column with the p-values to be plotted as arguments.

231    A total of 28 pairs of traits and chromosomes were found to be enriched in the combined analysis,

232    with scrotal circumference (BTA 5, 18, 9, and 21), milk glycosylated kappa-casein percentage

233    (BTA 6 and 16), inhibin level (BTA 5), triglyceride level (BTA 5), milk kappa-casein percentage

234    (BTA 6) and milk iron content (BTA 23) in the list of top 10 most enriched traits. Figure 4 shows

235    the top 5 enriched QTLs identified in this analysis.

236    *Relationship between studies and enriched QTLs*

237    An interesting functionality of GALLO is the graphical visualization of the relationship between

238    groups using a chord plot. The *relationship_plot()* function receives as arguments a dataframe (it

239    can use the gene or QTL annotation results, the QTL enrichment, or any other table with two

240   groups of information to be compared), the two groups to be compared (arguments x and y) and

241   the graphical arguments to set the size, color and gap between the sector in the chord plot. Figure

242   5 shows the chord plot obtained using a subset of the QTL annotation dataframe composed only

243   by the top 10 enriched traits and the studies which these traits were annotated. This plot indicates

244   that only inhibin levels and scrotal circumference on BTA5 are shared between Feugang et al.

245   (2009) and Buzanskas et al. (2017) [13,14]. Additionally, milk glycosylated kappa-casein

246   percentage (BTA 6 and 16), milk kappa-casein percentage (BTA 6) and milk iron content (BTA

247   23) were annotated only in Feugang et al. (2009) [14] and scrotal circumference (BTA 9, 18, 21)

248   and triglyceride level (BTA 5) were annotated only in Buzanskas et al. (2017) [13]. Inhibin is

249   produced by the Sertoli cells and can be used as a biomarker for sexual development [18]. In

250   addition, the inhibin levels were already associated with both scrotal circumference and sperm

251   quality traits in several studies, suggesting an important role in male fertility [19–23]. The results

252   obtained here through the integration of the GWAS results from two independent studies followed

253   by QTL annotation reinforces this association. Additionally, QTLs not associated with

254   reproductive phenotypes were identified in the enrichment analysis, suggesting the presence of

255   complex biological mechanisms such as a pleiotropic effect, epistasis and genetic hitchhiking

256   effect. Previous studies have highlighted the possible role of genomic regions with these kinds of

257   processes in the cattle genome [24,25]. An additional integration of the QTL annotation and

258   enrichment analysis performed here with the gene annotation and prospection for functional

259   candidate genes can be a powerful tool to better understand the genetic architecture and the

260   relationship among complex traits.

261

262

263 **Discussion**

264 The GALLO package is composed of a group of functions designed to perform an efficient and

265 direct downstream analysis for the gene and QTL annotation for candidate markers/SNPs,

266 haplotypes, genomic windows, runs of homozygosity, CNVs, etc. The functions implemented in

267 GALLO were designed to allow the integration of multiple datasets simultaneously. A brief

268 summary of these functions is shown in Table 1. For example, GWAS results from multiple traits

269 and/or populations or breeds can be analyzed together and compared or, individually analyzed in

270 the downstream analysis. This can be easily performed by adding an extra column in the input file

271 with the grouping factors to classify each dataset. These input files can be easily adapted from the

272 output of commonly used softwares to analyze high-throughput genomic data, such as PLINK,

273 BLUPF90, DESeq2, etc. [26–28]. In addition, GALLO provides a set of functions designed for

274 the visualization of the annotation results, overlap among groups, relationship between groups

275 (i.e., markers and candidate genes, datasets and QTLs, models and positional candidate genes,

276 etc.), and QTL enrichment results. This set of functions provides the capability of integrating

277 several results from multiple sources including different methodologies (GWAS, RNA-

278 sequencing, proteomics, etc.), populations (breeds, time-points, etc.), traits or the different

279 combination of these groups or others. Taken together, this set of functions provide to the

280 possibility to perform all the steps of gene/QTL annotation, comparison and summary in the same

281 environment. Additionally, the output obtained using GALLO was designed to allow a direct

282 connection between the candidate genomic regions and the genes/QTLs which overlap those

283 regions. Therefore, compared with outputs provided by other tools, such as biomaRt and Bedtools,

284 the interpretation of the output provided by GALLO is straightforward and easy to handle. Finally,

285 the QTL enrichment analysis available on GALLO is a useful and new approach that has the

286 potential to better understand the relationship between candidate genomic regions and the target

287 phenotype. It is important to highlight that despite the fact that GALLO was primarily designed

288 for livestock species, the package can perform gene annotation and data comparison for any other

289 species without any additional alterations to the input files. Regarding the QTL annotation and the

290 respective graphical visualization, the user should provide the gff file from the QTL database in a

291 format matching the gff files available on Animal QTLdb.

292 A summary of usage examples and output descriptions for all the functions available on GALLO

293 can be found in the reference manual (Supplementary File 4). It is important to highlight that the

294 two studies used as an example here are also part of the bovine QTL database. Consequently, the

295 results obtained here for annotation and enrichment would be expected, once the candidate regions

296 from the example file are present in the database used for the annotation. This approach was used

297 as a proof of concept of the methodology and indicates a precise annotation of the candidate

298 regions.

299 **Conclusion**

300 The integration of multiple datasets for gene and QTL annotation is one of the major bottlenecks

301 for the automatization of functional analysis of the results obtained using high-throughput

302 methodologies. The GALLO package provides a user-friendly and straightforward environment to

303 perform gene and QTL annotation, visualization, data comparison and QTL enrichment for

304 functional studies in livestock species. Consequently, the use of GALLO in the analyses of data

305 generated from high-throughput methodologies may improve the identification of hidden patterns

306 across datasets, datamining of previously reported associations, as well as efficiency in the

307 scrutinization of the genetic architecture of complex traits in livestock.

308 **Availability and requirements**

309 Project name: Genomic Annotation in Livestock for positional candidate LOci (GALLO)

310 Project home page: https://github.com/pablobio/GALLO

311 Operating system(s): Platform independent

312 Programming language: R

313 Other requirements: Depends: R (>= 3.5.0)

314 License: GPL-3

315 SciCrunch: SCR_019212

316 Bio.tools: biotools:genomic_annotation_in_livestock_for_positional_candidate_loci_gallo

317 **Availability of supporting data**

318 All of the data analyzed in the present study can be accessed in the public repository hosting the

319 R package (https://github.com/pablobio/GALLO). The input files and results used as examples in

320 the manuscript preparation are available in the supplementary Tables 1-4. A manual including

321 usage examples and output descriptions for all the functions available on GALLO can be found in

322 the package vignette (https://cran.r-project.org/web/packages/GALLO/vignettes/GALLO.html).

323 An archival copy of the code and supporting data is available via the GigaScience repository,

324 GigaDB [28].

325 **Declarations**

326 *List of abbreviations*

327 BP: position in base pairs; BP1: start position in base pairs; BP2: end position in base pairs; CHR:

328 Chromosome; CNV: Copy Number Variation; GALLO: Genomic Annotation in Livestock for

329    positional candidate Loci; GWAS: Genome-Wide Association Study; QTL: Quantitative trait loci;

330    SNP: Single Nucleotide Polymorphism.

344    *Authors' contributions*

345    PASF and AC were responsible for the conceptualization. PASF, ASV and AC were responsible

346    for the data processing and review of the codes. PASF and ASV were responsible for data curation.

347  PASF and GM were responsible for the implementation of the bioinformatic pipeline, integration

348  of datasets, and the coding. AC was responsible for funding acquisition.

349  *Acknowledgements*

350  Not applicable.

351  **References**

352  1. Ron M, Weller JI. From QTL to QTN identification in livestock - Winning by points rather

353  than knock-out: A review. Anim. Genet. 2007.

354  2. Ernst CW, Steibel JP. Molecular advances in QTL discovery and application in pig breeding.

355  Trends Genet. 2013.

356  3. Miglior F, Fleming A, Malchiodi F, Brito LF, Martin P, Baes CF. A 100-Year Review:

357  Identification and genetic selection of economically important traits in dairy cattle. J Dairy Sci.

358  2017;

359  4. Pértille F, Guerrero-Bosagna C, Silva VH Da, Boschiero C, Nunes JDRDS, Ledur MC, et al.

360  High-throughput and Cost-effective Chicken Genotyping Using Next-Generation Sequencing.

361  Sci Rep. 2016;

362  5. Cánovas A, Reverter A, DeAtley KL, Ashley RL, Colgrave ML, Fortes MRS, et al. Multi-

363  tissue omics analyses reveal molecular regulatory networks for puberty in composite beef cattle.

364  PLoS One. 2014;

365  6. Hu ZL, Park CA, Reecy JM. Building a livestock genetic and genomic information

366  knowledgebase through integrative developments of Animal QTLdb and CorrDB. Nucleic Acids

367  Res. 2019;

368  7. De Souza Fonseca PA, Id-Lahoucine S, Reverter A, Medrano JF, Fortes MS, Casellas J, et al.

369  Combining multi-OMICs information to identify key-regulator genes for pleiotropic effect on

370  fertility and production traits in beef cattle. PLoS One. 2018;13:1–22.

371    8. Fonseca PA de S, dos Santos FC, Lam S, Suárez-Vega A, Miglior F, Schenkel FS, et al.

372    Genetic mechanisms underlying spermatic and testicular traits within and among cattle breeds:

373    Systematic review and prioritization of GWAS results. J Anim Sci. 2018;

374    9. Suárez-Vega A, Gutiérrez-Gil B, Benavides J, Perez V, Tosser-Klopp G, Klopp C, et al.

375    Combining GWAS and RNA-Seq approaches for detection of the causal mutation for hereditary

376    junctional epidermolysis bullosa in sheep. PLoS One. 2015;

377    10. Durinck S, Moreau Y, Kasprzyk A, Davis S, De Moor B, Brazma A, et al. BioMart and

378    Bioconductor: A powerful link between biological databases and microarray data analysis.

379    Bioinformatics. 2005;

380    11. Quinlan AR, Hall IM. BEDTools: A flexible suite of utilities for comparing genomic

381    features. Bioinformatics. 2010;

382    12. R Core Team (2019). R: A language and environment for statistical computing. Accessed 1st

383    April 2019. 2019;

384    13. Buzanskas ME, Grossi D do A, Ventura RV, Schenkel FS, Chud TCS, Stafuzza NB, et al.

385    Candidate genes for male and female reproductive traits in Canchim beef cattle. J Anim Sci

386    Biotechnol. 2017;

387    14. Feugang JM, Kaya A, Page GP, Chen L, Mehta T, Hirani K, et al. Two-stage genome-wide

388    association study identifies integrin beta 5 as having potential role in bull fertility. BMC

389    Genomics. 2009;

390    15. Hackinger S, Zeggini E. Statistical methods to detect pleiotropy in human complex traits.

391    Open Biol. 2017.

392    16. Id-Lahoucine S, Molina A, Cánovas A, Casellas J. Screening for epistatic selection

393    signatures: A simulation study. Sci Rep. 2019;

394    17. Kühn C, Thaller G, Winter A, Bininda-Emonds ORP, Kaupe B, Erhardt G, et al. Evidence

395    for multiple alleles at the DGAT1 locus better explains a quantitative trait locus with major

396    effect on milk fat content in cattle. Genetics. 2004;

397  18. Phillips DJ. Activins, inhibins and follistatins in the large domestic species. Domest. Anim.
398  Endocrinol. 2005.

399  19. Fortes MRS, Reverter A, Kelly M, Mcculloch R, Lehnert SA. Genome-wide association
400  study for inhibin, luteinizing hormone, insulin-like growth factor 1, testicular size and semen
401  traits in bovine species. Andrology. 2013;

402  20. Fortes MRS, Reverter A, Hawken RJ, Bolormaa S, Lehnert S a. Candidate genes associated
403  with testicular development, sperm quality, and hormone levels of inhibin, luteinizing hormone,
404  and insulin-like growth factor 1 in Brahman bulls. Biol Reprod [Internet]. 2012 [cited 2013 Sep
405  6];87:58. Available from: http://www.ncbi.nlm.nih.gov/pubmed/22811567

406  21. Bame JH, Dalton JC, Degelos SD, Good TEM, Ireland JLH, Jimenez-Krassel F, et al. Effect
407  of Long-Term Immunization against Inhibin on Sperm Output in Bulls1. Biol Reprod. 1999;

408  22. Martin TL, Williams GL, Lunstra DD, Ireland JJ. Immunoneutralization of Inhibin Modifies
409  Hormone Secretion and Sperm Production in Bulls1. Biol Reprod. 1991;

410  23. Sato T, Kudo T, Ikehara Y, Ogawa H, Hirano T, Kiyohara K, et al. Chondroitin sulfate N-
411  acetylgalactosaminyltransferase 1 is necessary for normal endochondral ossification and
412  aggrecan metabolism. J Biol Chem. 2011;

413  24. Bolormaa S, Pryce JE, Reverter A, Zhang Y, Barendse W, Kemper K, et al. A Multi-Trait,
414  Meta-analysis for Detecting Pleiotropic Polymorphisms for Stature, Fatness and Reproduction in
415  Beef Cattle. PLoS Genet. 2014;10.

416  25. Love MI, Huber W, Anders S. Moderated estimation of fold change and dispersion for RNA-
417  seq data with DESeq2. Genome Biol. 2014;15:1–21.

418  26. Purcell S, Neale B, Todd-Brown K, Thomas L, Ferreira MAR, Bender D, et al. PLINK: A
419  tool set for whole-genome association and population-based linkage analyses. Am J Hum Genet.
420  2007;

421  27. Misztal I, Tsuruta S, Strabel T, Auvray B, Druet T, H LD. BLUPF90 and related programs
422  (BGF90). Proc 7th world Congr Genet Appl to Livest Prod. 2002. p. 743–4.

423    28. Fonseca PAS, Suárez-Vega A, Marras G, Cánovas A. GALLO: An R

424    package for Genomic Annotation and integration of multiple data source

425    in livestock for positional candidate LOci GigaScience Database. 2020.

426    http://dx.doi.org/10.5524/100834

427

428

429  **Tables**

430  Table 1: Description of the functions implemented in the GALLO package.

| Function | Description | Output |
|---|---|---|
| **Gene and QTL annotation** | | |
| import_gff_gtf | Import the gff and gtf files used for QTL and gene annotation, respectively | A dataframe composed by the information present in the gtf and gff files |
| find_genes_qtls_around_markers | Annotation of genes and QTLs around candidate regions | A data frame composed of the columns present in the input file and the genes or QTLs mapped within or around (if interval provided) the candidate regions |
| **Data visualization** | | |
| overlapping_among_groups | Overlap between grouping factors (such as different traits, statistical models, populations, studies, stc.) | A list with three matrices: 1) A matrix with the number of overlapping data; 2) A matrix with the percentage of overlap; 3) A matrix with the combination of the two previous ones |
| plot_overlapping | Plot overlap between data and grouping factors | A heatmap with the overlap between groups |
| plot_qtl_info | Plot QTL information from the gene or QTL annotation output | A pie plot (if QTL class is chosen) or a bar plot (if trait name is chosen) for the annotated QTLs |
| relationship_plot | Plot the relationship among the candidate regions or | A chord plot linking a grouping factor (genomic regions, traits, populations, etc.) with the annotated genes or QTLs |

| | | |
|---|---|---|
| | grouping factors with the annotated genes and QTLs | |

| **QTL enrichment** | | |
|---|---|---|
| qtl_enrich | Performs a QTL enrichment analysis based on a Bootstrap simulation for each QTL class or trait | A data frame composed of the enrichment results for QTL classes or traits present in the input file. 1) QTL: The QTL class or trait used for the enrichment; 2) CHR: The chromosome for that specific QTL or trait (if the option "chromosome" is informed to the argument enrich_type); 3) N_QTLs: Number of observed QTLs or traits in the dataset; 4) N_QTLs_db: Number of each annotated QTL in the qTL database; 5) Total_annotated_QTLs: Total number of annotated QTLs; 6) Total_QTLs_db: Total number of QTLs in the QTL database; 7) pvalue: P-value for the enrichment analysis; 8) adj.pval: The adjusted p-value based on the multiple test correction selected by the user; 9) QTL_type= The QTL type for each annotated trait. |
| QTLenrich_plot | Creates a bubble plot with the QTL enrichment results | A plot with the QTL enrichment results |

431

432

433

434

435

436

437    Table 2: Top 10 enriched QTLs for the combined analysis performed with the candidate regions from the two studies, Feugang et al.

438    (2009) and Buzanskas et al. (2017), used in the example dataset.

| QTL | CHR | # QTLs | # QTLs db | Total # QTLs | Total # QTLs db | p-value | FDR | QTL type |
|---|---|---|---|---|---|---|---|---|
| Scrotal circumference | 5 | 132 | 134 | 347 | 5942 | 1.56E-171 | 4.98E-169 | Reproduction |
| Scrotal circumference | 18 | 11 | 13 | 41 | 2147 | 2.20E-18 | 3.52E-16 | Reproduction |
| Scrotal circumference | 9 | 11 | 14 | 30 | 1395 | 2.04E-17 | 2.18E-15 | Reproduction |
| Milk glycosylated kappa-casein percentage | 6 | 71 | 1607 | 204 | 12158 | 1.86E-15 | 1.49E-13 | Milk |
| Inhibin level | 5 | 47 | 285 | 347 | 5942 | 3.38E-11 | 2.16E-09 | Reproduction |
| Scrotal circumference | 21 | 4 | 5 | 12 | 3606 | 3.51E-10 | 1.87E-08 | Reproduction |
| Milk kappa-casein percentage | 6 | 76 | 2637 | 204 | 12158 | 2.39E-07 | 1.01E-05 | Milk |
| Triglyceride level | 5 | 6 | 7 | 347 | 5942 | 2.53E-07 | 1.01E-05 | Health |
| Milk glycosylated kappa-casein percentage | 16 | 7 | 44 | 21 | 1440 | 1.29E-06 | 4.58E-05 | Milk |
| Milk iron content | 23 | 4 | 8 | 19 | 1159 | 3.48E-06 | 0.000111329 | Milk |

439

**Figure legends:**

**Figure 1:** Workflow explaining the main functions implemented on GALLO. The grey rectangles represent the functions, while the rounded and sharp rectangles represent the main goal of that respective function and its input, respectively.
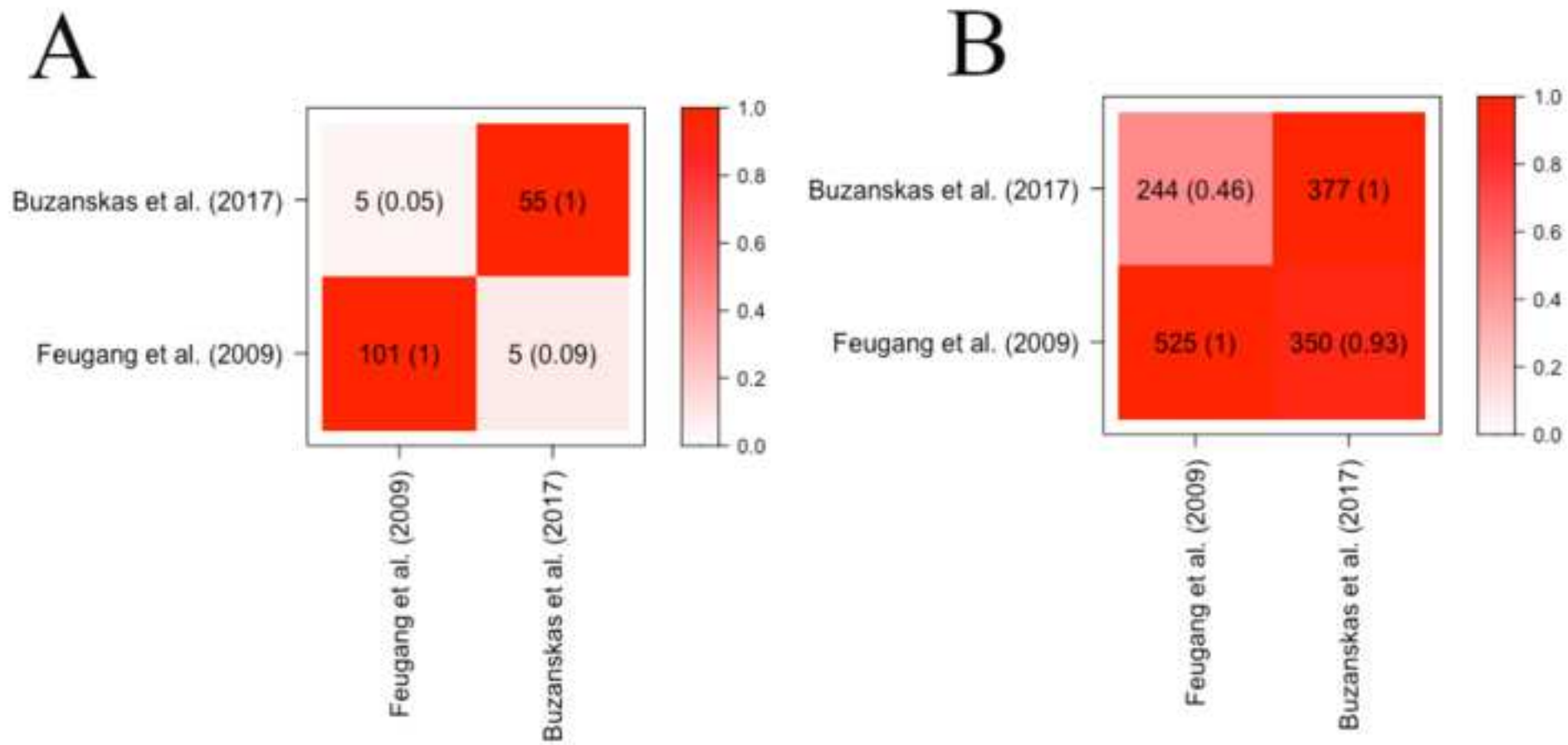
**Figure 2:** Overlapping between genes (A) and QTLs (B) annotated within the candidate regions (100 Kb downstream and upstream from the significant markers) from Feugang et al. (2009) and Buzanskas et al. (2017). The darker the color within the squares the higher the percentage of shared genes or QTLs.

**Figure 3:** Percentage of QTL type (pie plot) and trait related to Reproduction QTLs (barplots) for the QTL annotation results obtained for Feugang et al. (2009) (A), Buzanskas et al. (2017) (B) and the combined analysis (using both studies; C).
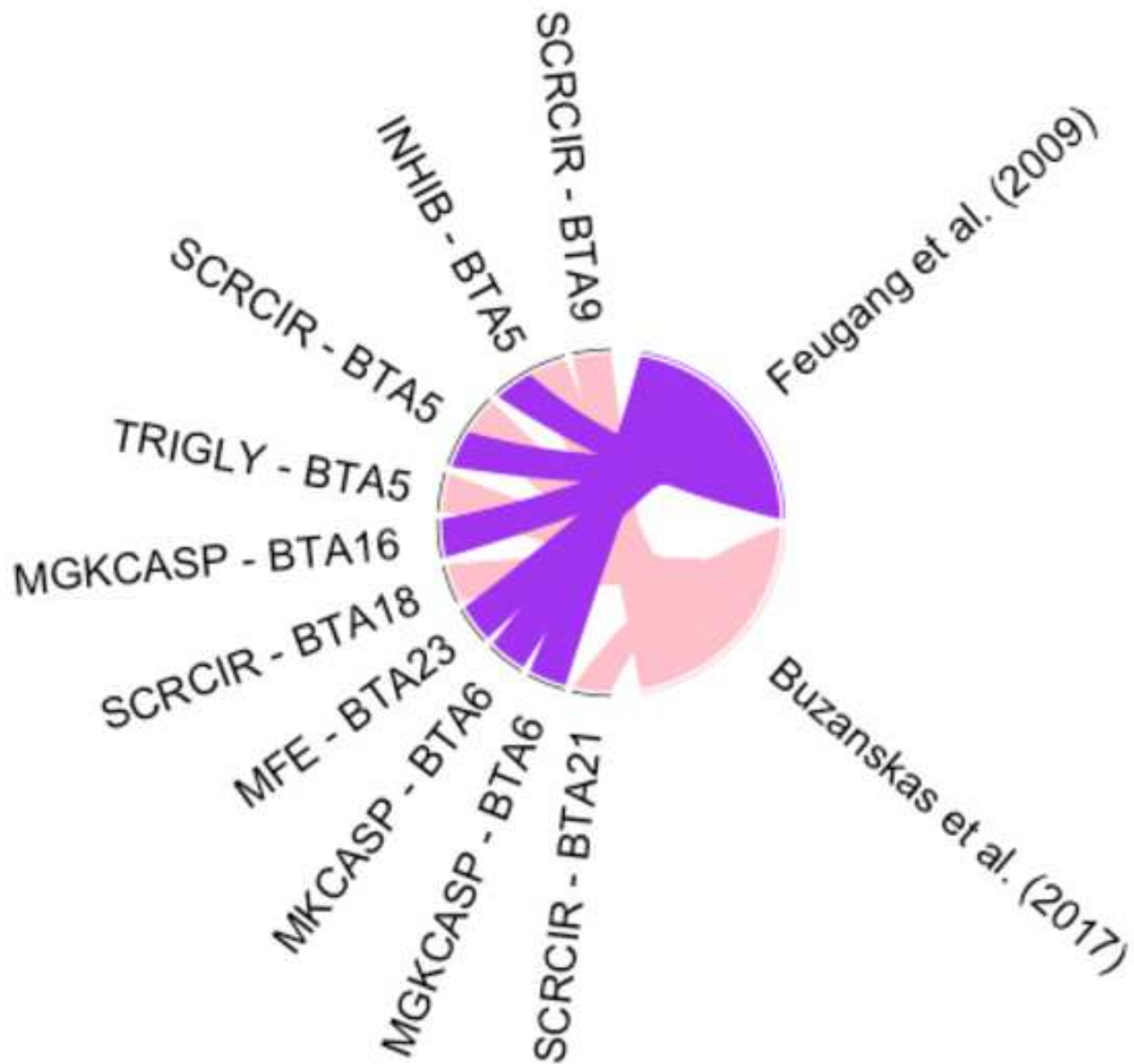
**Figure 4:** Bubble plot displaying the enrichment results for the top 5 enriched QTLs identified using the QTLs annotated within the candidate regions from Feugang et al. (2009) and Buzanskas et al. (2017). The darker the red shade in the circles, the more significant the enrichment. The area of the circles is proportional to the number of QTLs. The x-axis shows a richness factor obtained by the ratio of the number of QTLs annotated in the candidate regions and the total number of each QTL (and chromosome in the case of this plot) in the reference database.
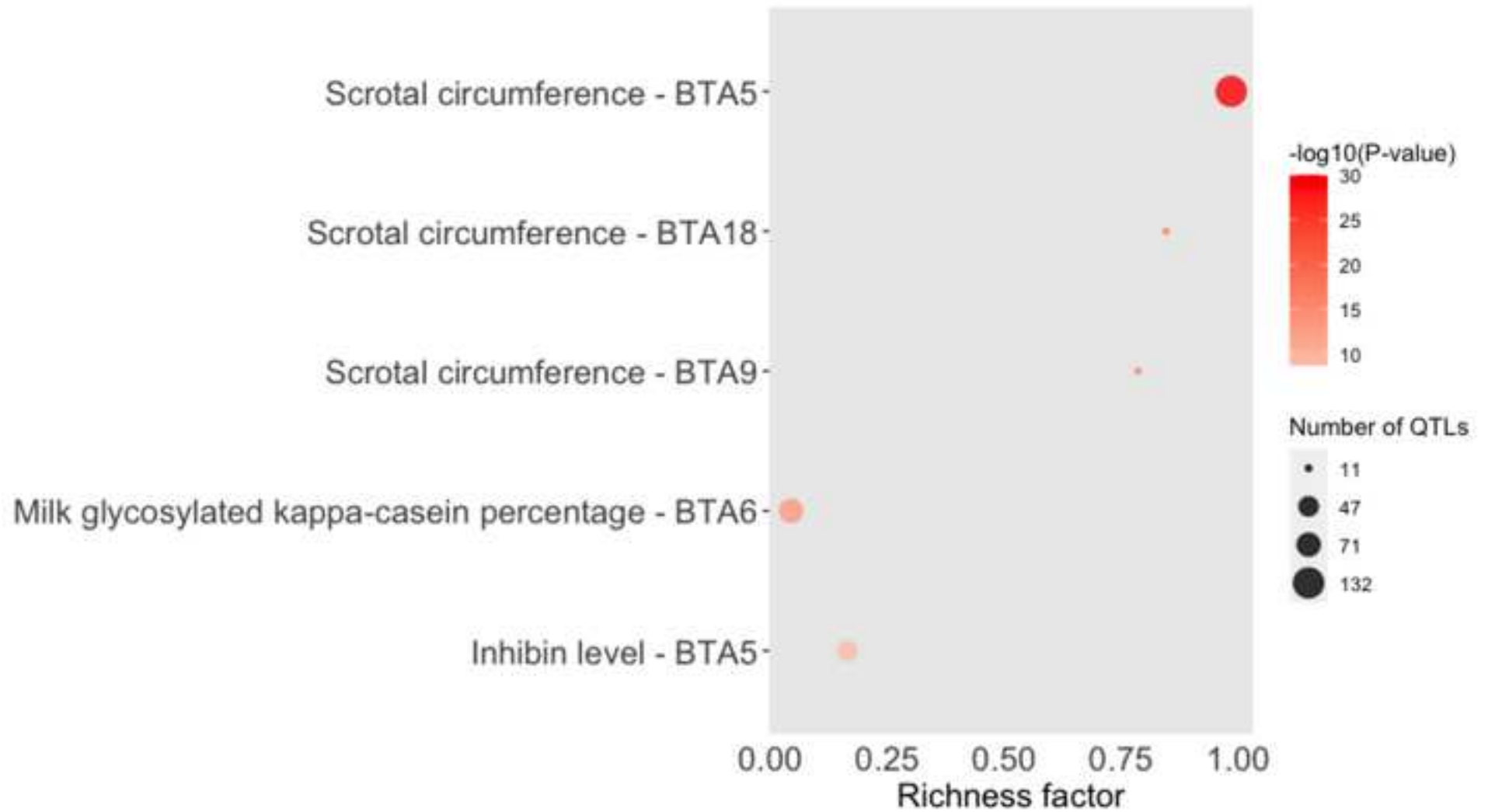
**Figure 5:** Chord plot showing the relationship between the top 10 enriched QTLs (Scrotal circumference – SCRCIR, Inhibin level – INHIB, Triglyceride level – TRIGLY, Milk glycosylated kappa-casein percentage – MGKCASP, Milk iron content – MFE, Milk kappa-casein percentage - MKCASP) and the studies (Feugang et al. (2009) in purple and Buzanskas et al. (2017) in pink).

Figure 1

Figure 2

Figure 3

Figure 3

Figure 5                                                                                      Click here to access/download;Figure;Figure5.png ±

Figure 4

Click here to access/download
**Supplementary Material**
Supplementary_File1.gtf

Click here to access/download
**Supplementary Material**
Supplementary_File2.gff.txt

Click here to access/download
**Supplementary Material**
Supplementary_file3.R

Click here to access/download
**Supplementary Material**
Supplementary_Table1.txt

Click here to access/download
**Supplementary Material**
Supplementary_Table2.txt

Click here to access/download

**Supplementary Material**

Supplementary_Table3.txt

Click here to access/download

**Supplementary Material**

Supplementary_Table4.txt

Guelph, September 1st, 2020

Dear Editorial Office,

We are pleased to re-submit the manuscript entitled "GALLO: An R package for Genomic Annotation and integration of multiple data source in livestock for positional candidate LOci" for consideration to publish it in the GigaScience. This is a resubmission of this manuscript after the inclusion of all the suggestion and considerations raised by the editor and the prior publication of the package in an official repository, in this case, the CRAN.

The present study introduces the applicability and the functionalities of GALLO package, developed in the R environment.

The identification of quantitative trait loci (QTLs) is a crucial step in the improvement of genomic selection and economic profitability in livestock. The development of high-throughput sequencing and genotyping methodologies and precision livestock farming allowed the identification of thousands of genomic regions associated with several complex traits. Consequently, the number of QTLs identified across the genome in livestock species increased substantially in the last years. Currently, in the Animal QTLdb it is possible to retrieve information about QTLs previously identified in cattle (127,191), chicken (11,340), horse (2,260), pig (29,865), rainbow trout (584) and sheep (3,001). The proper integration of the results obtained from different methodologies and technologies available is a crucial step for the accurate identification of the biological processes regulating the development of complex traits as well as the identification of potential functional candidate genes. However, currently, the integration of multiple data sources is not very straightforward due to limitations in the pipelines and algorithms implemented in the tools available for livestock. Moreover, although the automatization is possible, the direct link between the candidate regions and/or markers with the annotated genes and QTLs is missed. Consequently, this gap is forcing the user to back solve the overlap between the input and output files in order to perform the proper association between the candidate region and/or markers and the annotated genes and/or positional co-localized QTLs. In addition, nowadays there is still a lack of for customized QTL enrichment analyses in the available software and databases. Genomic Annotation in Livestock for positional candidate LOci (GALLO) is an R package, for the accurate annotation of genes and QTLs located in regions identified using the most common genomic analyses performed in livestock, such as Genome-Wide Association Studies and transcriptomics using RNA-Sequencing. Moreover, GALLO allows the graphical visualization of gene and QTL annotation results, data comparison among different grouping factors (e.g., methods, breeds, tissues, statistical models, studies, etc.), and QTL enrichment in different livestock species including cattle, pigs, sheep, chicken, etc. Consequently, GALLO is a useful package for annotation, identification of hidden patterns across datasets, datamining of previous reported associations, as well as the efficient scrutinization of the genetic architecture of complex traits in livestock.

We affirm that this manuscript has not been published elsewhere and is not under consideration by any other journal. All authors have approved the manuscript and agree with its submission to GigaScience.

The authors declare that they have no competing interests. With my best regards,

**Angela Cánovas, PhD**

Associate Professor Beef Genomics and Small Ruminants
University of Guelph
Department of Animal Biosciences
Centre for Genetics and Improvement of Livestock
Telephone: (519) 824-4129 ext. 56295
email: acanovas@uoguelph.ca