# Implementation of a new data enrichment analysis tool in the Animal QTL Database

Zhi-Liang Hu, Carissa A. Park, and James M. Reecy

Department of Animal Science, Iowa State University, 2255 Kildee Hall, 806 Stange Road, Ames, IA 50011

https://www.animalgenome.org/QTLdb

## Abstract

The development of Animal QTLdb over the past decade has provided a valuable resource for researchers to utilize a wealth of quantitative trait locus (QTL) and SNP association data to help elucidate potential genetic mechanisms that underlie traits of interest. The increasing growth of reported livestock animal QTL/association data in recent years provides new opportunities for data mining, meta-analysis, and big data analysis. We report here a new data enrichment analysis tool recently implemented in the QTLdb, which allows users to identify regions of the genome and to determine if they may be overrepresented in reported QTL/associations for traits of interest, similar to GO-enrichment testing of gene expression data. For example an analysis of 3,827 "milk yield" QTL/associations representing seven related traits found on 30 cattle chromosomes are "enriched" on chromosome 14. While the initial implementation of the enrichment tool demonstrates its utility, there are additional opportunities whereby this utility can be expanded by including more factors/parameters for more complex networked information analysis.

## Background

Gene Ontology (GO) [1, 2] enrichment analysis has been used on large gene expression datasets to effectively identify implicated biological processes underlying experimental outcomes [3]. The method has been useful and well described [4]. In simplified terms, both a gene and a QTL represent a region of the genome. Likewise, Gene Ontology terms can be associated with genes just as phenotypes are associated with QTL. Therefore, we hypothesized that it would be possible to analyze genome regions for over-enrichment of a particular phenotype/trait.

## Developments

As a preliminary trial, we assessed simple procedures to evaluate the enrichment of QTL/association data curated into the QTLdb using Chi-square analysis of a two-way contingency table (traits by chromosomes). The currently implemented tool was designed to allow evaluation of all reported QTL/associations for selected traits throughout a genome, to determine if the trait or traits are overrepresented in one or more regions. The analysis relies on an underlying assumption that the traits of interest are related. For example, they may belong to the same overarching trait type, or may be from a given trait ontology branch.

Figure 1 (a) shows the output from our initial implementation, in which 7,686 "growth" QTL/associations, representing 47 related traits (Figure 1 (b)) found on 30 cattle chromosomes, are "enriched" in certain chromosomal locations. Chi-square analysis was used to estimate the overrepresentation of reported QTL/associations in terms of their frequencies by chromosome locations. The contingency p-value (p) estimates are used to indicate the degree of overrepresentation (enrichment) of QTL/associations. The Benjamini-Hochberg procedure [5] is used to estimate the false discovery rate (FDR). The sizes of Chi-squares in each contingency category are graphically plotted with horizontal bars of varying lengths to indicate chromosomes or chromosome locations where larger numbers of QTL/associations are found.

If users are interested in sub-regions of a particular chromosome, they have the option to select a chromosome and give a window size in terms of the length of the chromosomal region to be evaluated (Figure 1 (a)).

Efforts are ongoing to expand functionality to allow options for users to select a set of traits for analysis based on their ontology reasoning. When available, trait correlation data are also appended to the enrichment results (Figure 1 (c)) to provide supporting information for user evaluation. This provides evidence to show that new potentials exist for further enrichment analysis involving factors that may be networked or related.

## Access to the tool

The data enrichment tool can be accessed via multiple paths on the Animal QTLdb website. The direct access to the tool is the "Whole genome enrichment analysis of QTL/associations" link in the "Search and Analysis" menu from a QTLdb species home page. There are also access points on other QTLdb data analysis tools where data that are currently being viewed can be further examined using the data enrichment analysis tool (refer to the website):

- "Search and Analysis" >> Search to "View genome distribution data for a trait" >> Click a trait type on the bottom list of the trait types >> "Try a genome wide enrichment analysis"

- "Search and Analysis" >> "Find related data by trait ontology search" >> Click a trait name for trait details view >> Explore Further: "Try a data enrichment analysis"

- "Data Browse": Open a trait hierarchy (Figure 2(b)) >> Navigate to a trait type of interest (Figure 2(a)) >> Click on a link-out icon >> "Find all QTL /association for this trait type" >> Below the genome plot there is a link to "Try enrichment analysis on this trait type" which will use the displayed trait data to start an enrichment analysis.

This is only a highlight of access routes. Other web access points to the enrichment analysis tool will also be created as the web site is continually developed.
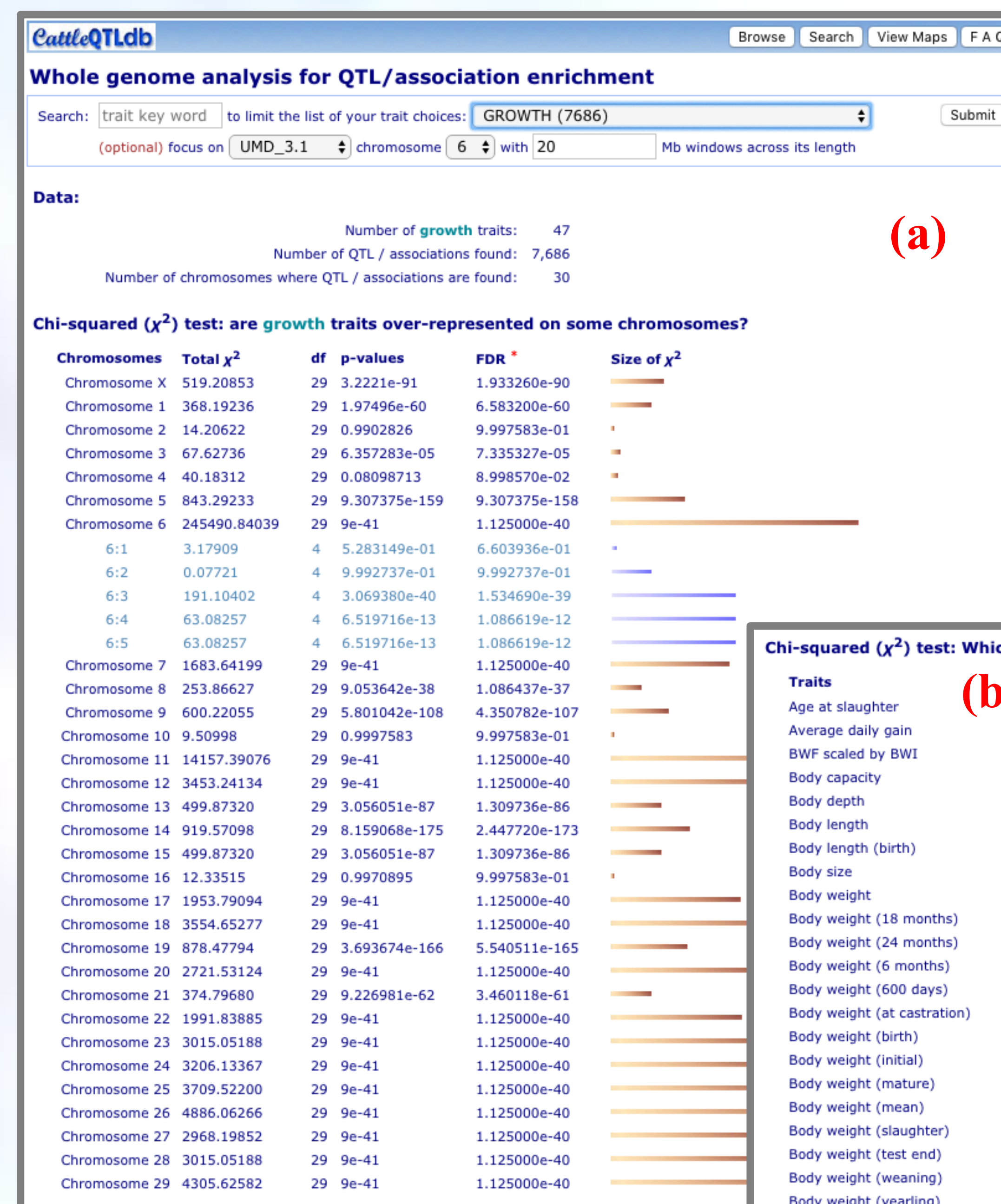
## Acknowledgments:

### Figure 1 (a, left; and b, below)

Screenshot of enrichment analysis on "growth" QTL/association data. Note that the screenshots of (a), (b), and (c) are all on one web page.

In (a) chromosome 6 was selected for detailed enrichment analysis with a 20-Mbp window.

In (b) is shown how each of the 47 "growth" traits contributes to the genome location-wise enrichment analysis outcome.

### Figure 1 (c, below)

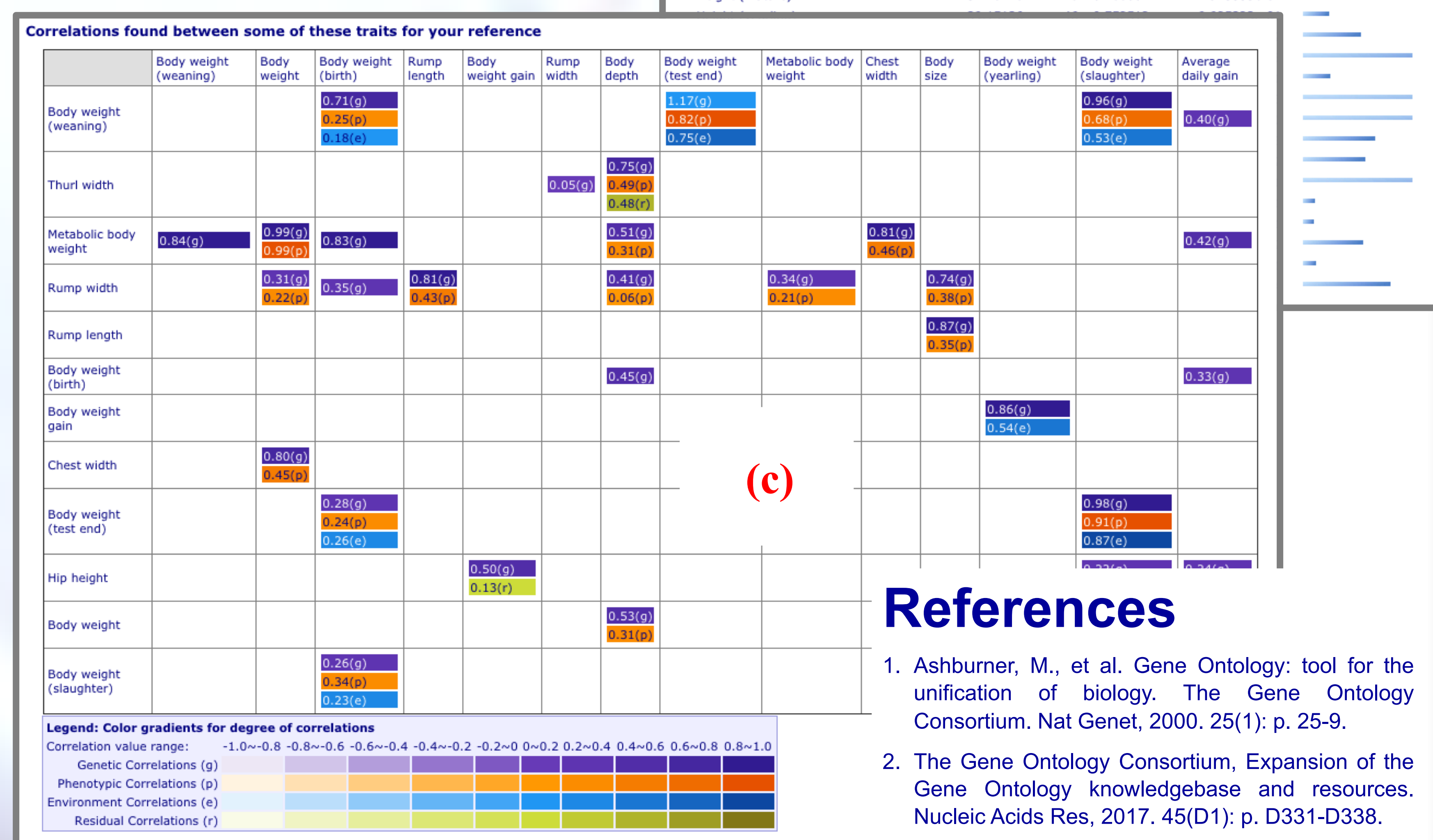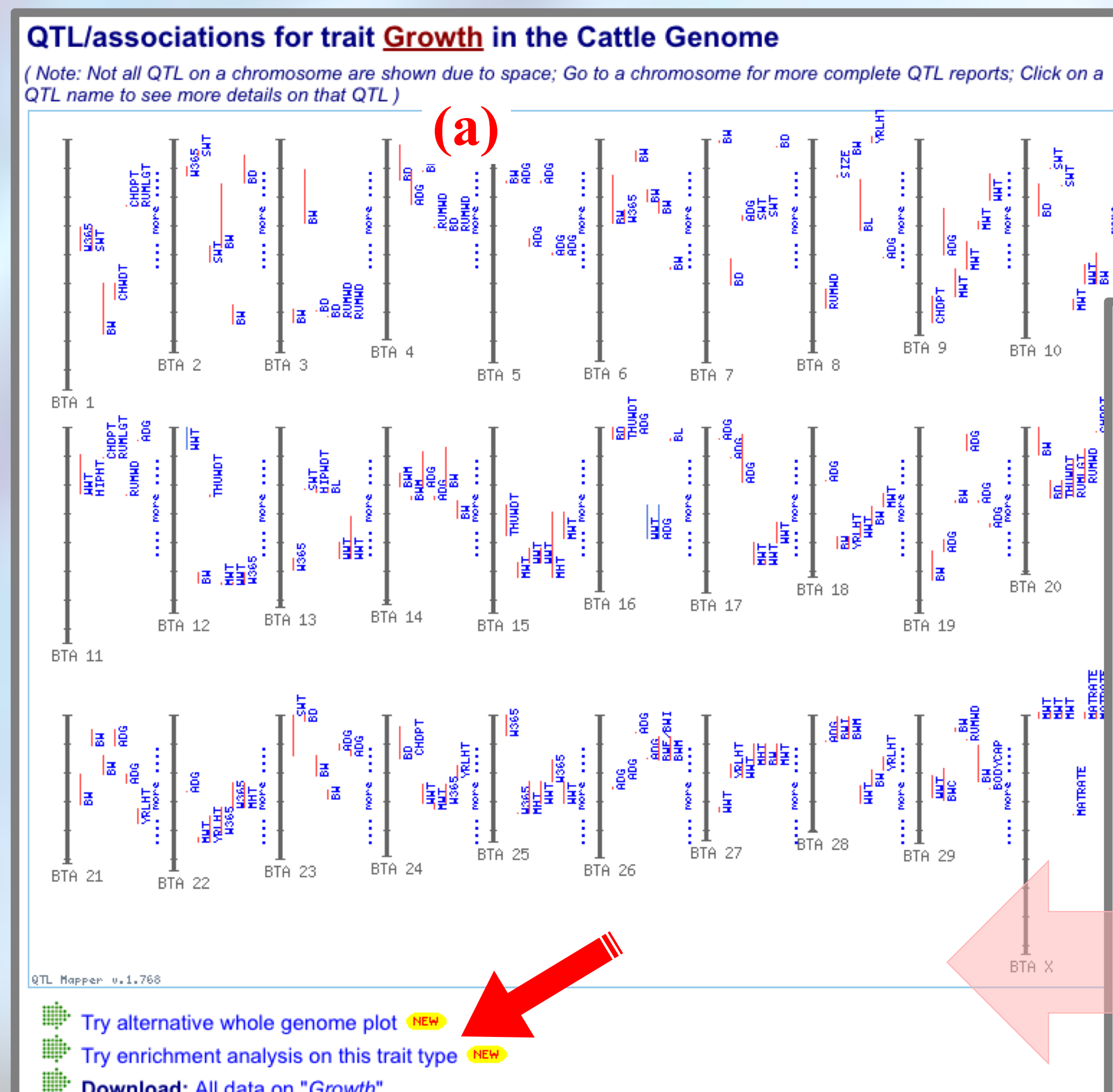Correlation data for the displayed traits is included when available.

### Figure 2

Links to the enrichment tool from other tools: from genome-wide data plot (a, below), and from the trait ontology browser (b, lower right).

## References

1. Ashburner, M., et al. Gene Ontology: tool for the unification of biology. The Gene Ontology Consortium. Nat Genet, 2000. 25(1): p. 25-9.

2. The Gene Ontology Consortium, Expansion of the Gene Ontology knowledgebase and resources. Nucleic Acids Res, 2017. 45(D1): p. D331-D338.

3. Huang, D.W., B.T. Sherman, and R.A. Lempicki, Bioinformatics enrichment tools: paths toward the comprehensive functional analysis of large gene lists. Nucleic Acids Res, 2009. 37(1): p. 1-13.

4. Subramanian, A., et al., Gene set enrichment analysis: a knowledge-based approach for interpreting genome-wide expression profiles. Proc Natl Acad Sci USA, 2005. 102(43): p. 15545-50.

5. Benjamini, Y. and Y. Hochberg, Controlling the false discovery rate: a practical and powerful approach to multiple testing. J R Stat Soc Series B Stat Methodol, 1995. 57(1).