

Target as **Applications Note** to *Bioinformatics*

Genetics and population analysis

**FastQTLmapping: an ultra-fast package for mQTL-like analysis**

Xingjian Gao<sup>1,2</sup>, Jiarui Li<sup>3</sup>, Xinxuan Liu<sup>1,2,4</sup>, Qianqian Peng<sup>3</sup>, Han Jing<sup>3</sup>, Sijia Wang<sup>3,5,6</sup>, Fan Liu<sup>1,2,\*</sup>

<sup>1</sup> Key Laboratory of Genomic and Precision Medicine, Beijing Institute of Genomics, Chinese Academy of Sciences, Beijing, China.

<sup>2</sup> China National Center for Bioinformation, Beijing, China.

<sup>3</sup> CAS Key Laboratory of Computational Biology, Shanghai Institute of Nutrition and Health, University of Chinese Academy of Sciences, Chinese Academy of Sciences, Shanghai, China.

<sup>4</sup> School of Future Technology, University of Chinese Academy of Sciences, Beijing, China.

<sup>5</sup> Taizhou Institute of Health Sciences, Fudan University, Taizhou, Jiangsu, China.

<sup>6</sup> Center for Excellence in Animal Evolution and Genetics, Chinese Academy of Sciences, Kunming, China.

\*To whom correspondence should be addressed.

## Abstract

### Summary

Here we describe fastQTLmapping, a C++ package that is computationally efficient not only for mQTL-like analysis but as a generic solver also for conducting exhaustive linear regressions involving extraordinarily large numbers of dependent and explanatory variables allowing for covariates. Compared to the state-of-the-art MatrixEQTL, fastQTLmapping was an order of magnitude faster with much lower peak memory usage. In a large dataset consisting of 3,500 individuals, 8 million SNPs, 0.8 million CpGs and 20 covariates, fastQTLmapping completed the mQTL analysis in 7 hours with 230GB peak memory usage.

### Availability and Implementation

FastQTLmapping is implemented in C++ and released under a GPL license. FastQTLmapping can be downloaded from <https://github.com/TianTTL/xQTLmapping>.

### Contact

[liufan@big.ac.cn](mailto:liufan@big.ac.cn)

1 **INTRODUCTION**

2 Methylation quantitative trait loci (mQTL) analysis is to test association between a large  
3 number of genomic variants and a large number of CpG sites over the genome, ideally using  
4 well-sized population samples to obtain sufficient statistical power, with covariates to control  
5 for potential confounding effects, and in an exhaustive manner to maximize genome  
6 resolution. Such analysis is often highly computationally burdensome, easily involving  
7 trillions of multiple regressions. MatrixEQTL<sup>1</sup> represents the state-of-the-art in terms of  
8 computational efficiency, yet has room for improvement.

9 Here we describe fastQTLmapping, a computationally efficient, exact, and generic solver for  
10 exhaustive multiple regression analysis involving extraordinarily large numbers of dependent  
11 and explanatory variables with covariates, which is particularly helpful in mQTL-like  
12 analysis.

13

14 **IMPLEMENTATION**

15 FastQTLmapping accepts input files in text format and in Plink binary format. The output file  
16 is in text format and contains all test statistics for all regressions, with the ability to control  
17 the volume of the output at preset significance thresholds. Different thresholds can be  
18 specified according to physical distances between the markers under investigation, which  
19 facilitates the analysis of cis- and trans-mQTLs. Z- and rank-normalizations are optional for  
20 pre-processing certain or all input variables. In order to maximize variable retention,  
21 fastQTLmapping determines missing values when concatenating individual dependent and  
22 explanatory variables, which in turn is quality controlled at a user-specified threshold.

23 FastQTLmapping is deployed on Linux using MKL (<https://software.intel.com/tools/onemkl>)

24 and GSL (<http://www.gnu.org/software/gsl/>) library, and is run from the command line. C++

25 source code, an example run, and documentation are freely available at

26 <https://github.com/TianTTL/xQTLmapping>.

27 FastQTLmapping loads the entire input file into the memory as characters and then converts

28 them into floating point numbers. A floating-point parser is developed to effectively handle

29 different data types in parallel on a per-line basis.

30 Redundant calculations related to covariates are removed using an orthogonal analysis.

31 Consider a multiple regression,  $\mathbf{C}$  is the matrix of  $k$  covariates,

35 
$$y = \alpha + \beta x + \gamma \mathbf{C} + \epsilon$$

32 orthogonalizing  $x$  and  $y$  with respect to  $\mathbf{C}$ . Decompose  $\mathbf{C}$  using QR factorization  $\mathbf{C} =$

33  $\mathbf{Q}\mathbf{R}$ , then centralize  $\mathbf{Q}$  to  $\widehat{\mathbf{Q}}$ , where  $\widehat{\mathbf{Q}}$  is an orthonormal basis. Then subtract the projections

34 of  $x$  and  $y$  on  $\widehat{\mathbf{Q}}$  to get  $x'$  and  $y'$ .

36 
$$\widehat{\mathbf{Q}} = [q_1 \ q_2 \ \dots \ q_k]$$

37 
$$x' = x - \sum_{i=1}^k \langle x, q_i \rangle q_i$$

38 
$$y' = y - \sum_{i=1}^k \langle y, q_i \rangle q_i$$

39 The subsequent analyses only involve univariate regressions  $y' = \alpha' + \beta' x' + \epsilon'$  to test

40  $H_0: \beta' = 0$  with  $k$  less degrees of freedom. FastQTLmapping speeds up the pair-wise

41 univariate regressions by calculating Pearson product-moment correlations, obtains the p-

42 values based on the t-distribution, and produces other test statistics if a p-value satisfies a

43 preset significant threshold.

44 Mathematical operations are accelerated using the Math Kernel Library (MKL), serial

45 computations are parallelized using OpenMP, and peak memory consumption is controlled  
46 through data splitting and variable reuse. Omics data are divided into blocks and dynamically  
47 assigned to threads. The block sizes are self-adapted based on the data size to maximize the  
48 performance of Level-3 BLAS while controlling for memory consumption. Critical section is  
49 used to guard smooth writing from threads to hard disks. All intermediate variables are  
50 allocated as aligned memory buffer.

51

## 52 **PERFORMANCE**

53 We compared the performance of fastQTLmapping and MatrixEQTL. For a fair comparison,  
54 we made a parallel version of MatrixEQTL using R packages ‘doParallel’, linked MKL to R  
55 environment, and manually split data to feed MatrixEQTL to achieve its optimal performance.  
56 The test data was downloaded from the GEO database (<https://www.ncbi.nlm.nih.gov/geo/>).  
57 The SNP and CpG data are GSE79254 and GSE79144, respectively as described previously.<sup>2</sup>  
58 The original data contains 54 individuals,  $1.5 \times 10^6$  SNPs, and  $4.5 \times 10^5$  CpGs. For testing  
59 purpose, we individual-wise resampled the data with replacement and generated 4 datasets  
60 consisting of 100, 200, 300 and 400 individuals. MQTL analysis was conducted with single  
61 CPU thread and with 8 CPU threads without covariates. The testing environment is CPU:  
62 Xeon E5 2686 V4 (18 cores, 2.3GHz), RAM: 256 GB ECC DDR4, OS: CentOS Linux 7, g++  
63 version: 4.8.5, R version: 4.0.3, and MKL version: 2021.3.0. Both fastQTLmapping and  
64 MatrixEQTL produced the same and exact results in the absence of missing values under all  
65 investigated settings. In presence of missing values, both packages provided approximate  
66 results, while fastQTLmapping results were always closer to the exact results. The operating

67 time of both fastQTLmapping and MatrixEQTL was largely linear to the data size. For  
68 computation, fastQTLmapping was 9.0-19.2 times faster than MatrixEQTL under the single-  
69 thread setting and 9.2-16.9 times faster under the 8-threads setting (**Figure 1A**). For I/O,  
70 fastQTLmapping was 19.6-34.7 times faster than MatrixEQTL under the single-thread setting  
71 and 11.6-16.8 times faster under the 8-threads setting (**Figure 1B**). The peak memory  
72 consumption of fastQTLmapping (1.7-4.9 GB) was much smaller than that of MatrixEQTL  
73 (9.7–13.7 GB) under the single-thread setting. The increase of the peak memory consumption  
74 when parallelizing to 8 threads was slower for fastQTLmapping (1.5-2.5 folds) than  
75 MatrixEQTL (5.7-6.9 folds, **Figure 1C**).

76 To further examine the performance of fastQTLmapping in well-sized population studies using  
77 densely imputed SNPs and Illumina MethylationEPIC 850K BeadChip, we also generated a  
78 dataset consisting of 3,500 individuals,  $8 \times 10^6$  SNPs, and  $8 \times 10^5$  CpGs by individual-wise, SNP-  
79 wise, and CpG-wise resampling the original data with replacement, and added 20 normally  
80 distributed variables as covariates. With 20 threads running in parallel, fastQTLmapping  
81 completed the mQTL analysis including I/O in 7.0 hours with 230 GB peak memory.

82

### 83 **CONCLUSIONS**

84 FastQTLmapping is ultra-fast, easy-to-deploy, capable for conducting pair-wise regression  
85 analysis at extraordinarily large scales on regular servers, particularly helpful for well-sized  
86 mQTL studies.

87

### 88 **Acknowledgements**

89 This work was supported by the Strategic Priority Research Program of Chinese Academy of  
90 Sciences [Grant No. XDB38010400, XDB38010100, XDC01000000], Shanghai Municipal  
91 Science and Technology Major Project [Grant No. 2017SHZDZX01], National Natural  
92 Science Foundation of China (NSFC) [81930056, 91631307], Science and Technology  
93 Service Network Initiative of Chinese Academy of Sciences [KJF-STS-QYZD-2021-08-001,  
94 KJF-STS-ZDTP-079], the National Key Research and Development Project [Grant No.  
95 2018YFC0910403], the Max Planck-CAS Paul Gerson Unna Independent Research Group  
96 Leadership Award, and the CAS Youth Innovation Promotion Association [Grant No.  
97 2020276].

98 *Conflict of interest:* none declared.

99

100 **REFERENCES**

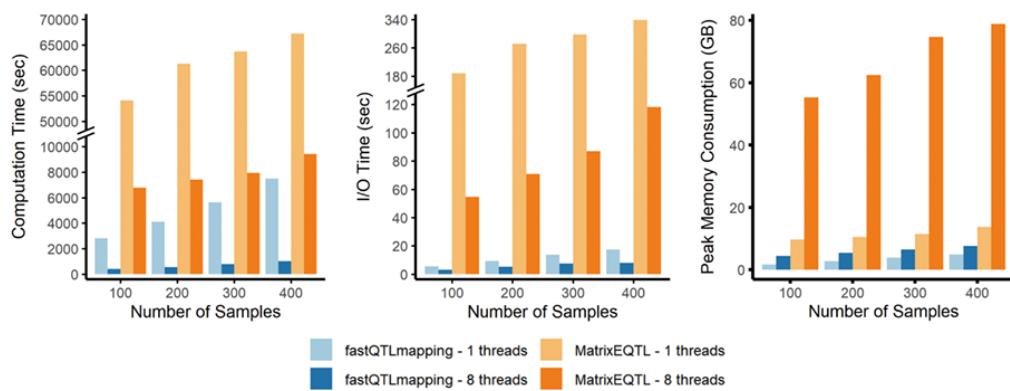
101 1. Shabalin, A.A. Matrix eQTL: ultra fast eQTL analysis via large matrix operations.  
102 *Bioinformatics* **28**, 1353-1358 (2012).

103  
104 2. Do, C. *et al.* Mechanisms and Disease Associations of Haplotype-Dependent Allele-  
105 Specific DNA Methylation. *Am J Hum Genet* **98**, 934-955 (2016).

106  
107

108 **Figure 1. Performance of fastQTLmapping and MatrixEQTL under various settings.**

109 A, Computation time; B, I/O time; C. Peak memory consumption.



110