Population-scale single-cell RNA seq empowers eQTL mapping **Cell transitions** B naive B memory **B** cells **Cell types** Genotype Environments, medical history, CD4+ T cells other phenotypes CD8+ Dendritic cells Rare migrating dendritic cells associated with disease **Cell states**

Cuomo et al, Nature Reviews Genetics 2023

Single-cell eQTL studies have gained prominence in the past 5 years



Cuomo et al, Nature Reviews Genetics 2023

Single-cell eQTL studies have gained prominence in the past 5 years



Cuomo et al, Nature Reviews Genetics 2023

OneK1K data: ~1,000 immune cells from each of ~1,000 individuals





Mem

48,023; 3.8%)

Imm, Naive

(n=82,068; 6.5%)

Classical

(n=38,233; 3.0%)

Non-classical

(n=15,166; 1.2%)

(n=8.690; 0.7%)

14 immune cell type:

CD4+

Naive, Eff, EM, CM

CD4+

Eff, CM

(n=61.786: 4.9%

CD4+

TGFB stimulated

CD8+

Eff

CD8+

Naive, Eff, Mem

CD8+

other

(n=34,528; 2.7%)

NK

mature

(n=9,677; 0.8%)

activated

(n=159.820:

Yazar et al, 2022

Typically, single-cell RNA-seq data are harmonized to **pseudo-bulk** for eQTL mapping





Cuomo et al, Genome Biology, 2021

1. Model repeat and complex data structure,

due to multiple cells per individual and relatedness between individuals



1. Model repeat and complex data structure,

due to multiple cells per individual and relatedness between individuals

2. Model discrete read counts



Expression levels of the SRGAP2 gene in induced pluripotent stem cells from the same ~100 individuals (Cuomo *et al*, Nature Reviews Genetics 2023)

1. Model repeat and complex data structure,

due to multiple cells per individual and relatedness between individuals

2. Model discrete read counts

3. Fast and scalable for large data, test 20k genes, 200 cell types, millions of cells, millions of variants



1. Model repeat and complex data structure,

due to multiple cells per individual and relatedness between individuals

2. Model discrete read counts

- **3. Fast and scalable for large data**, test 20k genes, 200 cell types, millions of cells, millions of variants
- **4. Test rare variations.** Single-variant test is underpowered



Motivation:

NO scalable computational tools exist to handle all current challenges for eQTL mapping with sc-RNA seq data

Mapping sc-eQTLs shares similar challenges with conducting GWAS for human diseases in large-scale biobanks

- 1. Model repeat and complex data structure
 - sc-eQTL: multiple cells per individual and relatedness between individuals
 - GWAS: population substructure and relatedness between individuals
- 2. Model phenotypes that are not continuous and Normally distributed
 - sc-eQTL: discrete read counts
 - GWAS: binary disease status, cases << controls
- 3. Fast and scalable for large data
 - sc-eQTL: test 20k genes, 200 cell types, millions of cells, millions of variants
 - GWAS: millions of genetic markers for thousands of human disease phenotypes
- 4. Test rare variations. Single-variant test is underpowered

SAIGE/SAIGE-GENE: R package suite developed for testing common and rare variant associations in large-scale biobanks

nature genetics

ANALYSIS https://doi.org/10.1038/s41588-018-0184-1

Efficiently controlling for case-control imbalance and sample relatedness in large-scale genetic association studies

Wei Zhou^{1,2}, Jonas B. Nielsen³, Lars G. Fritsche^{2,4,5}, Rounak Dey^{2,5}, Maiken E. Gabrielsen⁴, Brooke N. Wolford^{1,2}, Jonathon LeFaive^{2,5}, Peter VandeHaar^{2,5}, Sarah A. Gagliano^{2,5}, Aliya Gifford⁶, Lisa A. Bastarache⁶, Wei-Qi Wei⁶, Joshua C. Denny^{6,7}, Maoxuan Lin³, Kristian Hveem^{4,8}, Hyun Min Kang^{2,5}, Goncalo R. Abecasis^{2,5}, Cristen J. Willer⁵,^{13,9,10*} and Seunggeun Lee^{5,2,10*} TECHNICAL REPORT https://doi.org/10.1038/s41588-020-0621-6 nature genetics

Scalable generalized linear mixed model for region-based association tests in large biobanks and cohorts

Wei Zhou^{® 1,2,3,4,13} ⊠, Zhangchen Zhao^{® 1,5,13}, Jonas B. Nielsen[®]⁶, Lars G. Fritsche^{® 1,5}, Jonathon LeFaive^{® 1,5}, Sarah A. Gagliano Taliun^{® 1,5}, Wenjian Bi^{1,5}, Maiken E. Gabrielsen⁷, Mark J. Daly^{2,3,4,8}, Benjamin M. Neale^{® 2,3,4}, Kristian Hveem^{7,9}, Goncalo R. Abecasis^{1,5}, Cristen J. Willer^{® 6,10,11} and Seunggeun Lee^{® 1,5,12} ⊠

> nature oenetics

> > Check for updates

BRIEF COMMUNICATION https://doi.org/10.1038/s41588-022-01178-w

OPEN SAIGE-GENE+ improves the efficiency and accuracy of set-based rare variant association tests

Wei Zhou^{1,2,3,0}^{\infty}, Wenjian Bi^{4,5,6,0}^{\infty}, Zhangchen Zhao^{5,6,10}, Kushal K. Dey⁷, Karthik A. Jagadeesh⁷, Konrad J. Karczewski^{0,1,2,3}, Mark J. Daly^{1,2,3,8}, Benjamin M. Neale^{0,1,2,3} and Seunggeun Lee⁹

SAIGE-QTL: Generalized Linear Poisson Mixed-model for sc-eQTL mapping



1. Model repeat and complex data structure, due to multiple cells per individual and relatedness between individuals



- Mixed model with random effects 2. NIODEL DISCRETE READ COUNTS
- **3.** Generalized Poisson regression types, millions of cells, millions of variants

Approximations and computational speedups



★ MAE < 5

Set-based tests

SAIGE-QTL: Generalized Linear Poisson Mixed-model for sc-eQTL mapping

$$y \sim Poisson(\mu)$$
$$\mu = exp(\eta)$$
$$\eta = ZX_{ind}\alpha_1 + X_{cell}\alpha_2 + ZG\beta + Zb_{ind}$$
$$b_{ind} \sim N\left(0, \sum_{k=1}^{K} \tau_k \Psi_k\right)$$

- Z: the $N \times n$ design matrix indicating which cell is from which individual
- X_{ind}: covariates for individuals, e.g, age, sex, ancestry PCs
- X_{cell}: covariates for cells, e.g, pear factors
- **G**: genotype
- b_{ind} : $n \times 1$ vector for random effects
- *Ψ*: Genetic relationship matrix for related samples; Identity matrix for unrelated samples

Set-based association tests: grouping rare variants to test



 $S_j = \sum_{i=1}^n g_{ij}(y_i - \hat{\mu}_i)$ is the score statistic for the variant Under the null $H_0: \beta_i = 0$

Number of eGenes increases as number of cells increases with both methods







Using the same data, SAIGE-QTL identified the eQTL for the IBDassociated *ITGA4* locus that were missed in pseudo-bulk mapping



IBD: inflammatory bowel disease IBD GWAS: de Lange *et al*, Nature Genetics 2017 eQTL finempping: Kundu *et al*, Nature Genetics 2022 SAIGE-QTL identified the eQTL (only) in monocytes for the IBDassociated *ITGA4* locus (missed in pseudo-bulk mapping)

rs1375493

| Cell type | Beta | SE | P-value | |
|-----------|-------|------|---------|--|
| B_IN | 0 | 0.02 | 9.2e-01 | |
| B_Mem | -0.02 | 0.02 | 5.3e-01 | |
| CD4_ET | 0.03 | 0.02 | 6.7e-02 | |
| CD4_NC | 0.01 | 0.01 | 4.7e-01 | |
| CD4_SOX4 | 0.09 | 0.07 | 2.0e-01 | |
| CD8_ET | 0.01 | 0.01 | 5.4e-01 | |
| CD8_NC | 0.01 | 0.01 | 5.4e-01 | |
| CD8_S100B | 0.02 | 0.02 | 2.6e-01 | |
| DC | -0.01 | 0.04 | 8.3e-01 | |
| Mono_C | -0.2 | 0.03 | 4.9e-13 | |
| Mono_NC | -0.07 | 0.02 | 1.1e-03 | |
| NK | 0.02 | 0.01 | 2.7e-01 | — |
| NK_R | 0.02 | 0.03 | 5.7e-01 | |
| Plasma | 0.01 | 0.04 | 8.2e-01 | |
| | | | | |
| | | | | -0.275 -0.2 -0.15 -0.1 -0.05 0 0.05 0.1 0.15 0.2 |

Beta(Effect size)

IBD: inflammatory bowel disease

Across all 14 cell types, SAIGE-QTL identified 17,218 eGenes across all cell types (FDR<5%, as described below and in Methods), representing 5,894 unique genes of which 2,447 (**42%**) were specific to one cell type only



Values >10 are included

SAIGE-QTL identified additional immune diseaseassociated loci mediated by gene expression



- Among eGenes identified by SAIGE-QTL and TensorQTL across 14 cell types, gene-cell type pairs that are identified by the Summary-data—based Mendelian randomization (SMR, Zhu et al., 2016) framework to be associated with four autoimmune disease risks. CD: Crohn's disease, IBD: inflammatory bowel disease, RA: rheumatoid arthritis, SLE: systemic lupus erythematosus.
- e.g. 9q21 was shown to affect the risk of CD and IBD through changes in CARD9 expression in monocytes (MonoC)

Trans-eQTL regions for all protein coding genes in Naive CD4⁺T cells



Trans-region:

- MAF >= 10%
- Outside a 2Mb window around the gene

Trans-eQTL regions for all protein coding genes in Naive CD4⁺T cells



Trans-region:

- MAF >= 10%
- Outside a 2Mb window around the gene

SAIGE-QTL identified trans-eQTLs through genome-wide eQTL scan



cis
 trans

SPNS1 encodes SPNS lysolipid transporter 1, with potential role in maintaining **mitochondrial**-lysosomal homeostasis and preserving muscle mass and function (Zhang et al., 2023)



MRPL32 encodes Mitochondrial Ribosomal Protein L32. Among its related pathways are **Mitochondrial** translation and Metabolism of proteins (from GeneCards)

Set-based association tests: grouping rare variants to test



 $S_j = \sum_{i=1}^n g_{ij}(y_i - \hat{\mu}_i)$ is the score statistic for the variant Under the null $H_0: \beta_i = 0$ 5,541 eGenes (2,317 unique) eGenes have rare/less frequent signals (MAF <= 5%), 483 (21%) are independent from common eQTLs in the same genes



Different weighting strategies in the set-based tests across all cell types • Equal weights

- Beta(MAF, 1,25) upweight rarer variants
- based on the distance of each variant from the transcription start site (dTSS).





Linfeng Hu

95 out of 104 genes are eGenes in OneK1K (B_IN+B_Mem)

SAIGE-QTL: Generalized Linear Poisson Mixed-model for sc-eQTL mapping

- 1. Models repeat and complex data structure
- 2. Models discrete read counts
- 3. Fast and scalable for large data
- 4. Tests common and rare variations





Code and Data Availability

SAIGE-QTL is implemented as an open-source R method available at: <u>https://github.com/weizhou0/qtl</u>

Tutorial:

https://weizhou0.github.io/SAIGE-QTL-doc/

Code to reproduce all analyses and figures included here can be found at:

https://github.com/annacuomo/SAIGE_QTL_analyses

Summary statistics for the eQTLs in this study can be found on Zenodo:

https://zenodo.org/records/10811106

Pre-print

https://www.medrxiv.org/content/10.1101/2024.05.15.24307317v1.full.pdf

scRNA-seq data allow for identifying cell-level dynamic eQTLs



- CellRegMap. Cuomo et al., Mol Syst Biol (2022)
- GASPACHO. Kumasaka et al., Nature Genetics (2023)

More slides



lasma

Viono_c Viono_{nc}



Focus on three cell types:

- Most abundant: CD4 NC (n=463,528)
- 2. Middle: Naïve B cells (n=82,068)
- 3. Least abundant: Plasma cells (n=3,625)

Computation cost in OneK1K

| | | | Step 1 for 20,000 genes | | Step 2 for 20,000 genes | |
|-----------|----------------|-----------|-------------------------|-------------|-------------------------|---------------------|
| | | | | | | Time (CPU hours) |
| | | | | | Time (CPU hours) | for genome-wide |
| cell type | Individuals(n) | cells (N) | Time (CPU hours) | Memory (Gb) | for cis (+- 1Mb) | scan for 8M markers |
| Plasma | 795 | 3,625 | 24.99 | 0.23 | 46 | 604 |
| CD4_SOX4 | 857 | 4,065 | 25.78 | 0.23 | 47 | |
| DC | 968 | 8,690 | 35.88 | 0.24 | 46 | |
| NK_R | 969 | 9,677 | 37.56 | 0.25 | 50 | |
| Mono_NC | 969 | 15,166 | 49.72 | 0.26 | 47 | |
| CD8_S100B | 981 | 34,528 | 94.14 | 0.34 | 46 | |
| Mono_C | 969 | 38,233 | 106.79 | 0.35 | 49 | |
| B_Mem | 982 | 48,023 | 145.98 | 0.38 | 43 | |
| CD4_ET | 982 | 61,786 | 157.67 | 0.45 | 48 | |
| B_IN | 982 | 82,068 | 290.98 | 0.53 | 47 | 608 |
| CD8_NC | 982 | 133,482 | 337.27 | 0.73 | 50 | |
| NK | 982 | 159,820 | 425.84 | 0.80 | 54 | |
| CD8_ET | 982 | 205,077 | 596.32 | 0.97 | 50 | |
| CD4_NC | 982 | 463,528 | 1424.79 | 1.88 | 70 | 846 |

Compared to the R function glmer used in *Nathan et al.* (2022) paper

- 145 times faster for *cis*-eQTL mapping, and even more for *trans* mapping
- Allows for rare-variant association tests
- Allows for accounting for sample relatedness via the genetic relationship matrix
- Allows for accounting for cell-cell correlation with user-specified covariance matrices

Genes can with different sparsity levels (proportion of zero read counts)

• Example: Naïve B cells



SAIGE-QTL identified ~48% more eGenes than previous pseudo-bulk analysis using TensorQTL (FDR<5%)



Total: 11,585

Total: **17,033**

* Variants with MAF > 5% were tested

* Two analyses have highly correlated effect size estimates of cis-eQTLs (Pearson's R² = 0.98) in naive B cells

SAIGE-QTL has improved power with higher chi-square statistics



Simulation studies to evaluate the calibration of SAIGE-QTL based on OneK1K data

- Permuted genotypes across donors for genetic variants with MAF > 5% on Chr 1 for 20 times
- Randomly selected 50 genes from each sparsity group from three cell types
- Run single-variant assoc tests for Chr 1 using SAIGE-QTL for each gene



SAIGE-QTL has well calibrated false positive rates

