



(RESEARCH ARTICLE)



Statistical Approaches for Identifying eQTLs (Expression Quantitative Trait Loci) in Plant and Human Genomes

Tahiru Mahama*

Department of Mathematical Sciences, The University of Texas at El Paso.

International Journal of Science and Research Archive, 2023, 10(02), 1429-1437

Publication history: Received on 20 October 2023; revised on 21 December 2023; accepted on 29 December 2023

Article DOI: <https://doi.org/10.30574/ijrsra.2023.10.2.0998>

Abstract

Expression quantitative trait loci, or eQTLs, are genetic regions that play a crucial role in influencing how genes are expressed, making them a vital tool for connecting genetic makeup to observable traits in both plants and humans. In this review, a thorough overview of the statistical methods used in eQTL mapping, covering everything from traditional linear regression to more sophisticated techniques like mixed linear models, Bayesian inference, hidden confounder correction, multivariate frameworks, and machine learning algorithms were provided. The unique biological and computational hurdles that eQTL studies face in plants, such as polyploidy and genotype-by-environment interactions compared to those in humans, which often grapple with issues like tissue specificity, cell-type diversity, and ethical considerations were pointed out. Emerging trends like integrative multi-omics (including mQTLs and chromQTLs), single-cell eQTL mapping, graph-based genome modeling, and causal inference methods (like Mendelian randomization), all of which are enhancing the resolution and interpretability of eQTL analysis were explored. Additionally, popular software tools such as Matrix eQTL, FastQTL, and TensorQTL, evaluating their scalability, power, and replicability were discussed. As eQTL studies evolve towards greater complexity and clinical relevance, strong statistical modeling will continue to be essential for unraveling regulatory variations across various genomic landscapes.

Keywords: eQTL; Gene expression; Statistical genetics; Plant genomics; Human genomics; Machine learning; Pan-genome; Causal inference; Regulatory variants

1 Introduction

Understanding the genetic roots of gene expression variation is key to unraveling the mechanisms that drive phenotypic diversity and complex traits in both plants and humans. Expression quantitative trait loci, or eQTLs, are specific genomic regions that affect the levels of mRNA transcripts, creating a vital connection between genotype and transcriptomic phenotype (Gilad, Rifkin, & Pritchard, 2008). The mapping of these eQTLs has become a fundamental aspect of systems genetics, helping to dissect regulatory networks and pinpoint the causal genes that contribute to traits important for agriculture and medicine (Albert & Kruglyak, 2015). In the realm of plants, eQTL mapping has transformed functional genomics and strategies for crop improvement by uncovering transcriptional regulators linked to quantitative traits like drought tolerance, flowering time, and disease resistance (Kremling et al., 2018). For humans, extensive eQTL studies, such as those conducted by the Genotype-Tissue Expression (GTEx) Consortium, have laid the groundwork for understanding tissue-specific gene regulation and the genetic framework of complex diseases (GTEx Consortium, 2020). Even with the progress made in high-throughput sequencing technologies, accurately identifying eQTLs continues to pose computational and statistical challenges due to the intricacies of transcriptional regulation, population structure, batch effects, and the often high dimensionality of omics data (Sarkar et al., 2019). To tackle these challenges, a variety of statistical methods have been developed, ranging from traditional linear regression models to more advanced mixed models, Bayesian frameworks, and machine learning techniques (Zhou & Stephens, 2012; Stegle et al., 2012). For

* Corresponding author: Tahiru Mahama

instance, linear models that adjust for covariates (like PEER factors or surrogate variable analysis) have been commonly employed to account for hidden confounders in expression data (Leek & Storey, 2007; Stegle et al., 2012).

Mixed linear models are a game-changer in the world of plant genetics, especially when it comes to dealing with complex family trees or population structures. They cleverly incorporate random effects to help us understand genetic relationships (Korte & Farlow, 2013). Lately, Bayesian methods have stepped up, allowing for more nuanced eQTL mapping across different tissues and conditions, making the results easier to interpret (Li et al., 2022). On another front, deep learning and other data-driven techniques are starting to show promise as powerful tools for predicting regulatory variants and understanding the intricate, non-linear relationships in gene expression data (Zeng et al., 2021). In this review, a closer look at the range of statistical methods currently used for eQTL detection in both plants and humans was undertaken. Examination of their underlying assumptions, computational frameworks, strengths, and weaknesses was carried out, which was aimed towards providing a comprehensive view of how these methods have developed and how they can be utilized in future integrative omics research. Special attention was paid to the latest advancements that combine various omics layers, like methylation and chromatin accessibility, as well as methods tailored for single-cell eQTL analysis, which is rapidly reshaping our understanding of how specific cell types regulate gene expression (Van der Wijst et al., 2020). Additionally, the key challenges and opportunities in statistical eQTL mapping, touching on issues like power, replication, trans-eQTL detection, and causal inference were discussed. As the field moves towards more complex genomic analyses that span multiple scales and populations, having robust and interpretable statistical tools will be crucial for turning expression variation into meaningful biological and clinical insights. This review was designed to be a helpful resource for researchers navigating the evolving landscape of eQTL mapping and its growing significance in functional genomics, plant breeding, and precision medicine.

2 Fundamentals of eQTL Mapping

2.1 Overview of eQTLs

cis vs. trans Expression quantitative trait loci, or eQTLs, are specific regions in the genome that help explain the differences in gene expression levels among individuals. They serve as a bridge, connecting genetic variants like single nucleotide polymorphisms (SNPs) to how much of a particular transcript is present, which allows researchers to pinpoint the regulatory elements that influence gene activity (Gilad, Rifkin, & Pritchard, 2008). eQTLs are generally categorized into two types based on their location relative to the gene they regulate: cis-eQTLs and trans-eQTLs. Cis-eQTLs have an impact on the expression of genes that are situated close to the variant, usually within 1 Mb of the transcription start site. These are often easier to identify because they tend to have stronger effects and face less of a challenge from multiple testing (Albert & Kruglyak, 2015). In contrast, trans-eQTLs influence the expression of genes that are located far away, sometimes on different chromosomes or quite distant from the SNP locus. They usually exhibit smaller effect sizes, are more susceptible to confounding factors, and require larger sample sizes for reliable detection. Cis-eQTLs typically target promoter regions, enhancers, or untranslated regions that play a role in initiating transcription, while trans-eQTLs often operate through regulatory intermediates like transcription factors, miRNAs, or chromatin remodeling complexes (GTEx Consortium, 2020). Understanding the differences between cis- and trans-regulatory effects is essential for grasping the intricate nature of gene regulation and for differentiating between direct and indirect regulatory variants.

2.2 Biological Significance of eQTLs in Gene Regulation

eQTLs play a vital role as intermediaries between our genetic makeup and observable traits, shedding light on how variations in our DNA can lead to complex traits and diseases. They help map out the genetic landscape by pinpointing regulatory hotspots and essential transcriptional regulators (Cookson et al., 2009). Numerous genome-wide association studies (GWAS) have shown that SNPs linked to traits are often found in non-coding regions, underscoring the significance of regulatory variants revealed through eQTL mapping (Maurano et al., 2012). In human research, eQTL studies have been key in uncovering the genes responsible for diseases like schizophrenia, diabetes, and heart conditions (Battle et al., 2017). The Genotype-Tissue Expression (GTEx) Project has highlighted that the effects of eQTLs can vary by tissue type, which underscores the importance of the cellular environment in regulating gene expression (GTEx Consortium, 2020). In the plant kingdom, eQTLs have helped clarify the genetic pathways that control flowering time, responses to stress, and traits related to yield. For example, Kremling et al. (2018) showed that transcriptome-wide association studies (TWAS) in maize could identify promising candidate genes beyond the usual GWAS loci, paving the way for innovative approaches in crop genomics and breeding. Additionally, distinguishing between cis- and trans-eQTLs aids in reconstructing causal networks and gaining a systems-level understanding of gene interactions. Cis-eQTLs are often more directly linked to changes in gene expression, while trans-eQTLs tend to highlight regulatory hubs that have broader effects (Westra et al., 2013).

2.3 Challenges in eQTL Detection

Across Species While eQTL mapping is incredibly useful, it comes with its fair share of methodological and biological hurdles, particularly when we try to compare across different species. For starters, the population structure and genetic diversity can vary widely between plants and humans. Many plant species, especially crops, have intricate genomes characterized by polyploidy, high levels of linkage disequilibrium, and significant structural variation, which all make statistical inference a bit tricky (Korte & Farlow, 2013). On the flip side, human studies face their own set of limitations, including ethical concerns, dependence on post-mortem tissues, and the variability of environmental exposures. These factors make it tough to create standardized conditions for expression profiling (Li et al., 2019). Plus, issues like sample size and statistical power keep popping up: while it's relatively straightforward to detect cis-eQTLs even in smaller samples, trans-eQTLs demand larger cohorts and more rigorous statistical controls to steer clear of false positives. Cross-species comparisons are further complicated by batch effects, technical variability, and hidden confounding factors, such as cell composition or unmeasured environmental influences. To tackle these challenges, researchers often turn to statistical correction methods like surrogate variable analysis (SVA) and probabilistic estimation of expression residuals (PEER) (Leek & Storey, 2007; Stegle et al., 2012). Moreover, gene orthology and conservation introduce biological challenges when it comes to transferring eQTL findings across species. Regulatory networks can diverge significantly, even among closely related species, which limits how broadly we can apply eQTL results (Brawand et al., 2011). Still, comparative eQTL mapping remains a valuable tool for gaining evolutionary insights into gene regulation.

3 Statistical Frameworks for eQTL Detection

3.1 Linear Models and Regression-Based Methods

When it comes to eQTL analysis, linear regression is the go-to method. It helps us understand the relationship between genotype usually represented as 0, 1, or 2 for allele dosage and gene expression levels. The standard model looks like this: $Y = \beta_0 + \beta_1 X + \epsilon$

In this equation, Y represents gene expression, X stands for genotype, β_1 indicates the effect size, and ϵ is the random error. This method is popular because it's straightforward and easy to interpret. For example, the GTEx Consortium (2020) used simple linear regression models for their initial cis-eQTL detection across various tissues. That said, this approach does come with some assumptions, like the idea that samples are independent and that it overlooks potential confounding factors such as population stratification or hidden batch effects. Even with these limitations, linear models are still quite effective for cis-eQTL mapping, especially when paired with permutation testing to manage false discovery rates (FDR) (Shabalina, 2012).

3.2 Mixed Linear Models and Variance Component Approaches

MLM is a powerful technique that combines random effects to solve the problem of confounding factors such as kinship, genetic relatedness, or population structure. This means that they are ideally suited to plant systems and structured human cohorts. The general formula is as follows:

$$Y = X\beta + Zu + \epsilon$$

Here, Z is a design matrix for the random effects u (for example, a kinship matrix), which are assumed to be distributed according to a multivariate normal distribution. This approach is highly significant in plant eQTL research, where the generation of breeding populations can result in the non-independence of samples (Korte & Farlow, 2013).

Among the most recognized are:

- EMMAX (Kang et al., 2010)
- GEMMA (Zhou & Stephens, 2012)

These methods not only improve power and accuracy but also solve the problem of genetic relatedness very well. This is very important for trans-eQTLs, where fine differences can be easily mixed.

3.3 Bayesian Inference and Hierarchical Modeling

The use of Bayesian models in eQTL mapping is gaining popularity because these methods allow researchers to incorporate prior information, measure uncertainty, and handle multiple levels of data with relative ease; as such, they are particularly useful for estimating the posterior distributions of effect sizes for multiple conditions and tissues (e.g.,

Flutre et al. 2013). Additionally, Bayesian methods assist in prioritizing candidate regulatory variants by calculating posterior inclusion probabilities (PIPs) (Flutre et al., 2013). One of the standout benefits of Bayesian frameworks is their capability to analyze multiple SNPs or tissues simultaneously, which helps alleviate the burden of multiple testing and boosts statistical power for detecting trans-eQTLs.

3.4 Hidden Confounder Correction Methods (e.g., PEER, SVA)

When it comes to gene expression data, hidden technical or biological factors can really throw a wrench in the works. We're talking about things like batch effects, variations in cell types, or even how samples are handled. If we don't address these issues, we risk getting false positives or missing out on real connections. That's where methods like Surrogate Variable Analysis (SVA) and PEER (Probabilistic Estimation of Expression Residuals) come into play. They're popular tools for identifying and correcting those sneaky confounders. SVA works by modeling hidden variables as surrogate factors, using techniques like eigen-decomposition or regression residuals (Leek & Storey, 2007). On the other hand, PEER takes a Bayesian factor analysis approach to uncover latent variables that account for the variance in expression (Stegle et al., 2012). These methods have become essential in eQTL pipelines, such as GTEx, boosting the reproducibility and reliability of the signals we detect (GTEx Consortium, 2020).

3.5 Multivariate and Multi-Tissue eQTL Mapping

Multivariate models are a game changer because they allow us to look at multiple expression traits—like those across different tissues or genes all at once. This not only enhances our detection power but also helps us pinpoint pleiotropic eQTLs or effects that are shared across tissues. Some notable examples include:

- **MultiPhen:** A generalized linear model that works across various phenotypes
- **MT-eQTL:** Bayesian joint models designed for multi-tissue expression data (Flutre et al., 2013) These methods are especially useful for transcriptome-wide association studies (TWAS) and in plant research, where different tissues might exhibit eQTL effects that vary with developmental stages. They also support cross-tissue regulatory mapping, which is crucial for uncovering tissue-specific disease mechanisms in humans.

3.6 Machine Learning and Deep Learning-Based Models

Machine learning (ML) models are becoming more popular in eQTL studies, where they help with feature selection, pattern recognition, and predicting regulatory effects. These models excel at managing high-dimensional data and capturing complex, non-linear relationships that traditional linear models often overlook. Techniques like random forests, support vector machines (SVMs), and regularized regression (LASSO) have been employed to identify functional SNPs that affect gene expression (Lee et al., 2015). On the other hand, deep learning architectures, such as convolutional neural networks (CNNs), are capable of predicting SNPs that modulate expression directly from raw genomic sequences (Zeng et al., 2016). For instance, tools like DeepSEA and ExPecto leverage deep learning to understand how variations in sequences influence chromatin accessibility and gene expression across different tissues (Zhou et al., 2012). In the realm of plants, ML models have been utilized to forecast how expression responds to stress based on SNPs and methylation features

4 Comparative Analysis: Statistical Methods in Plant vs. Human Systems

4.1 eQTL Mapping in Plants: Challenges of Polyploidy and Structure

When it comes to expression QTL mapping in plants, researchers face a unique set of statistical and biological hurdles, primarily due to the intricate structures of plant genomes, polyploidy, and significant genetic and environmental variability. Unlike the diploid genomes found in humans, many key agricultural plants, like wheat (*Triticum aestivum*), can be tetraploid or even hexaploid. This complexity makes it tricky to pinpoint genetic effects to specific homoeologs (Kaur et al., 2022). The redundancy of genes can obscure or lessen the impact of regulatory variants, especially in trans-eQTL analyses. Moreover, linkage disequilibrium (LD) often spans larger genomic regions in self-fertilizing or inbred plant populations, which can lower the resolution of eQTL mapping (Korte & Farlow, 2013). In structured mapping populations, such as recombinant inbred lines (RILs) or backcross populations, it's crucial to consider population structure and familial relationships. This is typically done using mixed linear models (MLMs) or variance component models, which have become the go-to approach in plant eQTL research (Zhou & Stephens, 2012). On top of that, environmental variability has a more significant impact on plants because they are directly affected by various abiotic and biotic stresses. This leads to genotype-by-environment interactions (GxE) that complicate the interpretation of eQTLs, making it necessary to employ multi-environment and longitudinal study designs (Cubillos et al., 2014). For

instance, in maize, Kremling et al. (2018) utilized TWAS and cis-eQTL data across different developmental stages to identify candidate genes that influence complex agronomic traits.

4.2 Human eQTL Studies: Tissue Specificity and Disease Linkage

Unlike plant systems, human eQTL studies face unique challenges such as limited access to tissues, diverse sample types, and ethical considerations. Since gene expression is often specific to certain tissues and contexts, one of the biggest hurdles in human research is capturing the intricate spatiotemporal dynamics of gene regulation (GTEx Consortium, 2020). The GTEx project, a groundbreaking initiative in multi-tissue transcriptomics, has shown that while many cis-eQTLs are common across different tissues, trans-eQTLs and splicing QTLs (sQTLs) are usually quite tissue-specific. This specificity makes them harder to detect and replicate (Li et al., 2018). Additionally, the composition of cell types within tissues can mask genetic influences, which may lead to an underestimation or misinterpretation of eQTLs. To tackle these challenges, researchers are now using methods like cell-type deconvolution and single-cell eQTL mapping (van der Wijst et al., 2020). It's crucial to note that eQTLs in humans are key to understanding genome-wide association study (GWAS) signals, many of which are located in non-coding regulatory regions. Integrative strategies such as colocalization analyses and Transcriptome-Wide Association Studies (TWAS) have effectively connected eQTLs to common diseases like schizophrenia, type 2 diabetes, and coronary artery disease (Gusev et al., 2018). For instance, the colocalization of eQTLs at the FTO locus with obesity traits has shed light on the mechanisms behind non-coding GWAS signals (Claussnitzer et al., 2015).

4.3 Shared Insights and Methodological Convergence

Even though plants and humans are biologically different, the studies on expression quantitative trait loci (eQTL) in both fields are starting to align in terms of the statistical methods and computational tools they use. Popular tools like Matrix eQTL (Shabalina, 2012), FastQTL (Ongen et al., 2016), and TensorQTL (Taylor-Weiner et al., 2019) have become go-to options for scalable eQTL analysis across various species. On top of that, both areas are adopting Bayesian hierarchical models, machine learning techniques, and integrative multi-omics approaches to tackle the complexities of gene regulation. For example, joint eQTL mapping across different tissues or environments, probabilistic fine-mapping, and causal mediation analysis are being utilized in both plant breeding and human precision medicine (Wen et al., 2016). Furthermore, researchers in both plant and human studies are acknowledging the significance of trans-eQTL hotspots, gene co-expression networks, and regulatory modules. This points to a shared framework of gene regulation that is influenced by master regulators or the structure of chromatin (Cubillos et al., 2014; Battle et al., 2017). Unified Insight: Whether the goal is to enhance crop yields or predict disease risks, the fundamental challenge remains the same: breaking down complex, high-dimensional, and noisy biological data to pinpoint regulatory loci that have real phenotypic implications.

5 Advances in Integrative and High-Dimensional eQTL Analyses

The field of eQTL mapping has seen a remarkable shift, moving from traditional single-layer transcriptomics to more integrative, high-dimensional methods. These new approaches provide us with enhanced resolution, deeper biological context, and valuable causal insights.

5.1 Integrating eQTLs with Other Omics Layers (e.g., mQTLs, chromQTLs)

Recent developments have highlighted the importance of combining eQTL data with other types of QTLs, such as methylation QTLs (mQTLs), chromatin accessibility QTLs (chromQTLs), and histone QTLs. This integration helps us better grasp the regulatory pathways that connect genotype to phenotype. For instance, eQTLs that overlap with chromQTLs in immune cells have shed light on the enhancer promoter interactions that contribute to the risk of autoimmune diseases (Kumasaka et al., 2019). Multi-omic platforms like ENCODE and BLUEPRINT serve as essential frameworks for these integrative analyses, allowing us to pinpoint causal variants within complex regulatory landscapes.

5.2 eQTL Mapping from Single-Cell RNA

seq Data Single-cell RNA sequencing (scRNA-seq) has revolutionized our ability to detect cell-type-specific eQTLs, effectively addressing the challenges posed by tissue heterogeneity. Unlike bulk data, single-cell eQTLs reveal dynamic and context-specific regulatory effects, even uncovering eQTLs that are active only in rare or fleeting cell states (van der Wijst et al., 2020). Despite facing statistical hurdles like sparsity and dropout events, innovative methods such as sc-eQTLGen and CellRegMap have shown great promise in clarifying the intricate relationships between expression and genotype.

5.3 Causal Inference: Mendelian Randomization and Mediation Models

To shift our understanding from mere correlation to actual causation, researchers are now embracing Mendelian Randomization (MR) and causal mediation frameworks. These tools help us explore whether gene expression plays a role in mediating the effects of genetic variants on complex traits (Zhu et al., 2016). Techniques like SMR (Summary-based Mendelian Randomization) and TWMR (Transcriptome-wide Mendelian Randomization) combine GWAS and eQTL data to highlight potential causal genes. Thanks to these innovative approaches, we've been able to identify regulatory variants linked to diseases such as schizophrenia, type 2 diabetes, and coronary artery disease (Gusev et al., 2016).

5.4 Trans-eQTLs: Detection, Interpretation, and Limitations

Detecting trans-eQTLs can be quite a challenge due to their small effect sizes and the complications that come with multiple testing. Yet, their biological significance is huge, they often act as regulatory hubs, like transcription factors or signaling molecules. Recent advancements in statistical methods, including meta-analysis, variance component models, and network propagation, have made it easier to spot these elusive elements (Westra et al., 2013). Still, one of the biggest hurdles we face is telling apart genuine trans-regulation from technical noise or artifacts caused by linkage disequilibrium, especially in complex or polyploid genomes.

6 Benchmarking and Evaluation of Statistical Methods

6.1 Simulation-Based Assessments

Simulation frameworks provide researchers with a way to evaluate eQTL mapping methods in controlled environments where the true associations are already known. By tweaking parameters like effect size, minor allele frequency (MAF), linkage disequilibrium (LD), and sample size, these simulations can measure method sensitivity, specificity, and type I error rates (Westra et al., 2013). Tools like simeQTL and custom R/Python scripts are frequently used to replicate realistic genotype-expression relationships, taking into account confounders, population structure, and sparse expression profiles. This kind of simulation-based benchmarking has revealed, for example, that linear models excel at detecting strong cis-eQTLs, while Bayesian or mixed models tend to be more effective for trans-eQTLs and structured populations (Zhou & Stephens, 2012).

6.2 Power, False Discovery, and Replicability

When it comes to eQTL detection, three key metrics stand out: statistical power, false discovery rate (FDR), and replicability across different cohorts or tissues. Power tends to increase with effect size, sample size, and MAF, but it often diminishes in trans-eQTLs due to their generally weaker effects. To control FDR, researchers typically use methods like the Benjamini-Hochberg procedure, permutation testing, or Bayesian false discovery models (Storey & Tibshirani, 2003). Replicability can be particularly challenging, especially in trans-eQTL studies or when comparing across tissues. Nowadays, cross-validation and meta-analysis across independent datasets are commonly employed to confirm eQTL signals (GTEx Consortium, 2020). Metrics such as π_1 statistics, reproducibility indices, and posterior probabilities are useful for quantifying replicability.

6.3 Available Software and Pipelines (e.g., Matrix eQTL, FastQTL, TensorQTL)

There's a solid lineup of powerful and scalable software tools ready for eQTL detection: -

- Matrix eQTL (Shabalín, 2012): This is a popular R-based tool that uses linear models and allows for optional covariate correction. It's designed to be fast and memory-efficient, especially when working with large datasets.
- FastQTL (Ongen et al., 2016): This tool features a speedy permutation-based cis-eQTL mapping pipeline, making it a great choice for RNA-seq data. It's particularly useful for integrating with GWAS.
- TensorQTL (Taylor-Weiner et al., 2019): A GPU-accelerated Python tool that can handle millions of eQTL tests across various tissues or cell types, supporting both cis- and trans-eQTL analyses. While these pipelines have different model assumptions, scalability, and flexibility, they all work well with standard formats like VCF, BED, and GTF. More and more, you'll find them integrated into Snakemake or Nextflow workflows, which helps ensure reproducibility and enables high-throughput analysis

7 Current Trends and Future Directions

The world of eQTL analysis is changing quickly, thanks to the integration of pan-genomics, graph-based methods, precision breeding applications, and ethical considerations in human studies. These new directions show a move away from traditional, linear models toward more dynamic, population-aware, and translational frameworks in regulatory genomics.

7.1 Pan-genome and Pan-transcriptome eQTL

Studies In the past, eQTL studies typically depended on a single reference genome, which limited discoveries in genetically diverse populations. However, the rise of pan-genome and pan-transcriptome frameworks allows researchers to include presence-absence variations (PAVs), structural variants, and new transcripts that aren't found in reference genomes (Golicz et al., 2016). In the realm of plants, where domestication and polyploidy lead to significant genomic variation, pan-genomic eQTL analyses have uncovered new regulatory loci that would go unnoticed in single-reference studies (Hirsch et al., 2014). For humans, initiatives like the Human Pangenome Reference Consortium (HPRC) are working to enhance variant detection and expression mapping in underrepresented populations (Ebbert et al., 2022).

7.2 Application of Graph-Based Genomics in eQTL

Analysis Graph-based genomics marks a significant shift from linear reference models to variation graphs, where alternative alleles and haplotypes are depicted in a graph structure. This method improves read alignment, variant calling, and genotype-expression associations, particularly in genomes that are structurally diverse or admixed (Garrison et al., 2018). Tools such as VG (Variation Graph Toolkit) and PanGenie facilitate graph-based eQTL analysis, enhancing the discovery of regulatory variants within complex loci like the MHC region. Graph eQTL mapping is especially useful for capturing regulatory diversity across pan-genomes and multi-ethnic cohorts (Eggertsson et al., 2019).

7.3 Ethical and Data-Sharing Considerations in Human

eQTL Studies As eQTL studies tap into the vast resources of biobank-scale data, they face a host of ethical dilemmas surrounding informed consent, privacy, and the representation of diverse ancestries. It's essential to ensure that transcriptomic datasets reflect diversity and equity, as this is key to making findings applicable to a broader population (Fatumo, 2020). Additionally, controlled-access policies, like those implemented by GTEx and the UK Biobank, need to strike a balance between promoting open science and safeguarding participant privacy. Guidelines from organizations such as the Global Alliance for Genomics and Health (GA4GH) advocate for data sharing that is interoperable, secure, and ethically sound (Knoppers, 2014).

7.4 Prospects for Translational and Breeding Applications

In the realm of human genomics, eQTL mapping is becoming a vital part of precision medicine, helping to identify drug targets, forecast therapeutic responses, and fine-tune polygenic risk scores. Techniques like TWAS and fine-mapped eQTLs have already been instrumental in validating targets for conditions such as Crohn's disease and breast cancer. When it comes to crop improvement, eQTLs are paving the way for marker-assisted selection, functional trait analysis, and gene editing approaches (like CRISPR-Cas9) by identifying regulatory variants that influence traits such as yield, stress resilience, and flowering time (Kremling et al., 2018). The incorporation of eQTL data into genomic selection models is set to speed up breeding cycles and improve the predictability of traits.

8 Conclusion

Expression quantitative trait loci (eQTL) mapping has become a fundamental aspect of functional genomics, serving as a crucial bridge between genetic variation and gene regulation in both plants and humans. This review highlights how progresses have been made from basic linear models to more sophisticated Bayesian frameworks, mixed-effect modeling, and machine learning approaches, significantly enhancing our ability to pinpoint regulatory variants with both accuracy and clarity. When the two fields are compared, we see that plant studies often deal with challenges like polyploidy and environmental factors, while human studies face issues related to tissue specificity and ethical considerations. Despite these differences, both areas are increasingly utilizing similar computational tools, integrating multi-omics data, and employing fine-mapping techniques. The rise of high-resolution technologies like single-cell RNA sequencing, graph-based reference genomes, and pan-genomic frameworks has broadened the scope of eQTL research, allowing for analyses that are specific to cell types and aware of population differences. At the same time, improvements

in causal inference methods, such as Mendelian randomization and mediation models, are shifting eQTLs from mere descriptive correlations to deeper mechanistic insights that hold promise for real-world applications. Looking ahead, the combination of innovative statistical methods, a variety of genomic resources, and a commitment to ethical practices will be essential for leveraging eQTLs in areas like trait enhancement, drug target identification, and precision medicine. Whether it's about developing crops that can withstand stress or fine-tuning polygenic risk scores, strong eQTL analytics will continue to play a pivotal role in unraveling the complexities of the regulatory genome.

References

- [1] Albert, F. W., & Kruglyak, L. (2015). The role of regulatory variation in complex traits and disease. *Nature Reviews Genetics*, 16(4), 197-212.
- [2] Brawand, D., Soumillon, M., Necsulea, A., Julien, P., Csárdi, G., Harrigan, P., ... & Kaessmann, H. (2011). The evolution of gene expression levels in mammalian organs. *Nature*, 478(7369), 343-348.
- [3] Claussnitzer, M., Dankel, S. N., Kim, K. H., Quon, G., Meuleman, W., Haugen, C., ... & Kellis, M. (2015). FTO obesity variant circuitry and adipocyte browning in humans. *New England Journal of Medicine*, 373(10), 895-907.
- [4] Cookson, W., Liang, L., Abecasis, G., Moffatt, M., & Lathrop, M. (2009). Mapping complex disease traits with global gene expression. *Nature Reviews Genetics*, 10(3), 184-194.
- [5] Cubillos, F. A., Coustham, V., & Loudet, O. (2012). Lessons from eQTL mapping studies: non-coding regions and their role behind natural phenotypic variation in plants. *Current Opinion in Plant Biology*, 15(2), 192-198.
- [6] Ebbert, M. T., Jensen, T. D., Jansen-West, K., Sens, J. P., Reddy, J. S., Ridge, P. G., ... & Fryer, J. D. (2019). Systematic analysis of dark and camouflaged genes reveals disease-relevant genes hiding in plain sight. *Genome biology*, 20(1), 97.
- [7] Eggertsson, H. P., Kristmundsdottir, S., Beyter, D., Jonsson, H., Skuladottir, A., Hardarson, M. T., ... & Melsted, P. (2019). GraphTyper2 enables population-scale genotyping of structural variation using pangenome graphs. *Nature communications*, 10(1), 5402.
- [8] Fatumo, S. (2020). The opportunity in African genome resource for precision medicine. *EBioMedicine*, 54.
- [9] Flutre, T., Wen, X., Pritchard, J., & Stephens, M. (2013). A statistical framework for joint eQTL analysis in multiple tissues. *PLoS genetics*, 9(5), e1003486.
- [10] Garrison, E., Sirén, J., Novak, A. M., Hickey, G., Eizenga, J. M., Dawson, E. T., ... & Durbin, R. (2018). Variation graph toolkit improves read mapping by representing genetic variation in the reference. *Nature biotechnology*, 36(9), 875-879.
- [11] Gilad, Y., Rifkin, S. A., & Pritchard, J. K. (2008). Revealing the architecture of gene regulation: the promise of eQTL studies. *Trends in genetics*, 24(8), 408-415.
- [12] Golicz, A. A., Bayer, P. E., Bhalla, P. L., Batley, J., & Edwards, D. (2020). Pangenomics comes of age: from bacteria to plant and animal applications. *Trends in Genetics*, 36(2), 132-145.
- [13] GTEx Consortium. (2020). The GTEx Consortium atlas of genetic regulatory effects across human tissues. *Science*, 369(6509), 1318-1330.
- [14] Gusev, A., Mancuso, N., Won, H., Kousi, M., Finucane, H. K., Reshef, Y., ... & Price, A. L. (2018). Transcriptome-wide association study of schizophrenia and chromatin activity yields mechanistic disease insights. *Nature genetics*, 50(4), 538-548.
- [15] Hirsch, C. N., Foerster, J. M., Johnson, J. M., Sekhon, R. S., Muttoni, G., Vaillancourt, B., ... & Buell, C. R. (2014). Insights into the maize pan-genome and pan-transcriptome. *The Plant Cell*, 26(1), 121-135.
- [16] Kang, H. M., Sul, J. H., Service, S. K., Zaitlen, N. A., Kong, S. Y., Freimer, N. B., ... & Eskin, E. (2010). Variance component model to account for sample structure in genome-wide association studies. *Nature genetics*, 42(4), 348-354.
- [17] Knoppers, B. M. (2014). Framework for responsible sharing of genomic and health-related data. *The HUGO journal*, 8(1), 3.
- [18] Korte, A., & Farlow, A. (2013). The advantages and limitations of trait analysis with GWAS: a review. *Plant methods*, 9, 1-9.

- [19] Kremling, K. A., Diepenbrock, C. H., Gore, M. A., Buckler, E. S., & Bandillo, N. B. (2019). Transcriptome-wide association supplements genome-wide association in *Zea mays*. *G3: Genes, Genomes, Genetics*, 9(9), 3023-3033.
- [20] Kumasaka, N., Knights, A. J., & Gaffney, D. J. (2019). High-resolution genetic mapping of putative causal interactions between regions of open chromatin. *Nature genetics*, 51(1), 128-137.
- [21] Laboratory, D. A., Fund, N. C., Site—NDRI, B. C. S., Site—RPCI, B. C. S., Resource—VARI, B. C., of Miami, B. B. R. U., ... & Statistical Methods groups—Analysis Working Group. (2017). Genetic effects on gene expression across human tissues. *Nature*, 550(7675), 204-213.
- [22] Lee, D., Gorkin, D. U., Baker, M., Strober, B. J., Asoni, A. L., McCallion, A. S., & Beer, M. A. (2015). A method to predict the impact of regulatory variants from DNA sequence. *Nature genetics*, 47(8), 955-961.
- [23] Leek, J. T., & Storey, J. D. (2007). Capturing heterogeneity in gene expression studies by surrogate variable analysis. *PLoS genetics*, 3(9), e161.
- [24] Li, Y. I., Knowles, D. A., Humphrey, J., Barbeira, A. N., Dickinson, S. P., Im, H. K., & Pritchard, J. K. (2018). Annotation-free quantification of RNA splicing using LeafCutter. *Nature genetics*, 50(1), 151-158.
- [25] Maurano, M. T., Humbert, R., Rynes, E., Thurman, R. E., Haugen, E., Wang, H., ... & Stamatoyannopoulos, J. A. (2012). Systematic localization of common disease-associated variation in regulatory DNA. *Science*, 337(6099), 1190-1195.
- [26] Ongen, H., Buil, A., Brown, A. A., Dermitzakis, E. T., & Delaneau, O. (2016). Fast and efficient QTL mapper for thousands of molecular phenotypes. *Bioinformatics*, 32(10), 1479-1485.
- [27] Sarkar, A. K., Tung, P. Y., Blischak, J. D., Burnett, J. E., Li, Y. I., Stephens, M., & Gilad, Y. (2019). Discovery and characterization of variance QTLs in human induced pluripotent stem cells. *PLoS genetics*, 15(4), e1008045.
- [28] Shabalin, A. A. (2012). Matrix eQTL: ultra fast eQTL analysis via large matrix operations. *Bioinformatics*, 28(10), 1353-1358.
- [29] Stegle, O., Parts, L., Durbin, R., & Winn, J. (2010). A Bayesian framework to account for complex non-genetic factors in gene expression levels greatly increases power in eQTL studies. *PLoS computational biology*, 6(5), e1000770.
- [30] Taylor-Weiner, A., Aguet, F., Haradhvala, N. J., Gosai, S., Anand, S., Kim, J., ... & Getz, G. (2019). Scaling computational genomics to millions of individuals with GPUs. *Genome biology*, 20(1), 228.
- [31] Van Der Wijst, M. G., Brugge, H., De Vries, D. H., Deelen, P., Swertz, M. A., LifeLines Cohort Study, ... & Franke, L. (2018). Single-cell RNA sequencing identifies celltype-specific cis-eQTLs and co-expression QTLs. *Nature genetics*, 50(4), 493-497.
- [32] Wen, X., Pique-Regi, R., & Luca, F. (2017). Integrating molecular QTL data into genome-wide genetic association analysis: Probabilistic assessment of enrichment and colocalization. *PLoS genetics*, 13(3), e1006646.
- [33] Westra, H. J., Peters, M. J., Esko, T., Yaghootkar, H., Schurmann, C., Kettunen, J., ... & Franke, L. (2013). Systematic identification of trans eQTLs as putative drivers of known disease associations. *Nature genetics*, 45(10), 1238-1243.
- [34] Zeng, H., Edwards, M. D., Liu, G., & Gifford, D. K. (2016). Convolutional neural network architectures for predicting DNA-protein binding. *Bioinformatics*, 32(12), i121-i127.
- [35] Zhou, X., & Stephens, M. (2012). Genome-wide efficient mixed-model analysis for association studies. *Nature genetics*, 44(7), 821-824.
- [36] Zhu, Z., Zhang, F., Hu, H., Bakshi, A., Robinson, M. R., Powell, J. E., ... & Yang, J. (2016). Integration of summary data from GWAS and eQTL studies predicts complex trait gene targets. *Nature genetics*, 48(5), 481-487.