

# Differentially Expressed Heterogeneous Overdispersion Genes Testing for Count Data\*

Yubai Yuan<sup>1</sup>, Qi Xu<sup>2</sup>, Agaz Wani<sup>3</sup>, Jan Dahrendor<sup>3</sup>, Chengqi Wang<sup>3</sup>, Janelle  
Donglasan<sup>3</sup>, Sarah Burgan<sup>3</sup>, Zachary Graham<sup>3</sup>, Monica Uddin<sup>3</sup>, Derek  
Wildman<sup>3</sup>, and Annie Qu<sup>2</sup>

<sup>1</sup>Department of Statistics, The Pennsylvania State University

<sup>2</sup>Department of Statistics, University of California Irvine

<sup>3</sup>College of Public Health, University of South Florida

## Abstract

The mRNA-seq data analysis is a powerful technology for inferring information from biological systems of interest. Specifically, the sequenced RNA fragments are aligned with genomic reference sequences, and we count the number of sequence fragments corresponding to each gene for each condition. A gene is identified as differentially expressed (DE) if the difference in its count numbers between conditions is statistically significant. Several statistical analysis methods have been developed to detect DE genes based on RNA-seq data. However, the existing methods could suffer decreasing power to identify DE genes arising from overdispersion and limited sample size. We propose a new differential expression analysis procedure: heterogeneous overdispersion genes testing (DEHOGT) based on heterogeneous overdispersion modeling and a post-hoc inference procedure. DEHOGT integrates sample information from all conditions and provides a more flexible and adaptive overdispersion modeling for the RNA-seq read count. DEHOGT adopts a gene-wise estimation scheme to enhance the detection power of differentially expressed genes. DEHOGT is tested on the synthetic RNA-seq read count data and outperforms two popular existing methods, DESeq and EdgeR, in detecting DE genes. We apply the proposed method to a test dataset using RNAseq data from microglial cells. DEHOGT tends to detect more differently expressed genes potentially related to microglial cells under different stress hormones treatments.

**Key words:** Differential expression, Gene expression, Generalized linear modeling, RNA-Seq data

---

\*This work is supported by the National Institutes of Health (R01MD011728).

# 1 Introduction

High-throughput sequencing of DNA fragments and mRNA-seq techniques are powerful tools based on next generation sequencing technologies [16] for monitoring RNA abundance to detect genetic variation. Specifically, for RNAseq, the sequenced RNA fragments are aligned with reference genome sequences, and the number of sequence fragments assigned to each gene is counted for each sample. Then we can compare read counts between different biological conditions or between different genetic variants to infer genetic information based on biological systems of interest [19]. In the analysis of RNA-seq data, read counts do not have a prior upper bound, thus regression models based on a binomial distribution with a pre-specified number of trials do not apply [34]. Linear regression is therefore not feasible as count data is always a non-negative integer. More importantly, RNA-seq data presents high overdispersion, implying that the variance of the count can be much larger than its mean. Given that the sample sizes are typically small for RNA-seq analysis due to the cost and other factors, statistical modeling needs to address the large variation from the data and to improve the power of detecting differential gene expressions.

One fundamental clinical interest of applying RNA-seq analysis is to understand the mechanism of post-traumatic stress disorder (PTSD) formulation. PTSD is a common severe psychiatric disorder that develops following exposure to a life-threatening or traumatic experience [31]. PTSD is known to cause negative effect on an individual's life quality via the PTSD condition itself or the relevant comorbidities. Previous works [12, 18] show that only a small proportion of individuals experience traumatic events will develop PTSD. Meanwhile, the majority of people exposed to trauma are resilient even after repeated exposures to trauma [35]. In addition, various risk factors of PTSD have been identified such as low socio-economic status, social support and gender [33, 6, 15].

Significant individual heterogeneity of either response to trauma or the PTSD development originates from the individual epigenetic variability. Specifically, previous studies reveal the connection between PTSD and immune system functioning, and several genes such as FKBP5 involved with the immune system are also found to be differentially expressed among PTSD individ-

uals [27, 32, 17]. In particular, previous work [13] has identified monocytes as a key cell type in differentiating male subjects with versus without lifetime PTSD. In addition, rodent studies have implicated peripheral monocytes in inducing anxiety-like behavior through trafficking of proinflammatory monocytes to the brain via activated microglia. Following this line of research, in this paper, we collect RNA-seq data from the well-designed lab experiments to investigate differential expression of genes in human microglia cells under different immune characteristic environments. This is an important step for understanding the role of microglia cells and immune-related genes in PTSD development.

The main challenge in analyzing microglial RNA-seq datasets lies in the high and heterogeneous overdispersion in the read counts. As an illustration, Figure 1 shows the histogram of the empirical RNA read counts from microglial data, where the read counts are highly spread out and the variance can be much larger than the mean. Several differential expression analysis methods have been developed to address the overdispersion issue in RNA-seq read counts. Among these methods, the DESeq [1] and EdgeR [23] are the most popular and are implemented and available using the R [8]. Specifically, the DESeq analyzes count data by using a shrinkage estimation for dispersions as well as fold changes to improve stability and interpretability of estimates. EdgeR is designed for the analysis of replicated count-based expression data, and is based on the method developed by Robinson and Smyth [25] using an overdispersed Poisson model to account for the read count variability. However, most existing methods adopt the shrinkage strategy when estimating the level of overdispersion by assuming that genes with similar expression strength have homogeneous dispersion levels. Although overdispersion regularization helps to increase the robustness of inference against the uncertainty due to limited sample size, it decreases the discriminative power in detecting differentially expressed genes with strong overdispersion effects at the population level.

In this paper, we propose a new differential expression analysis framework based on generalized linear modeling. Compared with other popular RNA-seq analysis methods such as DESeq and EdgeR, the main advantages of the proposed method for differentially expressed heteroge-

neous overdispersion genes testing (DEHOGT) are as follows. First, our method jointly estimates the fold change and overdispersion parameters over samples from all treatment conditions, which increases the effective sample size and leads to more accurate inference. Second, and more importantly, our model adopts a within-sample independent structure among genes without assuming that genes with similar expression strength have homogeneous dispersion levels. Therefore, our method can better account for the heterogeneity in count dispersion and select more relevant genes. Third, our method allows for fully independent gene-wise inference and hence can achieve computational scalability to handle large gene datasets by implementing parallel computing. Finally, the proposed method enjoys the flexibility of adapting different overdispersion patterns by allowing different count generating distributions in the inference procedure.

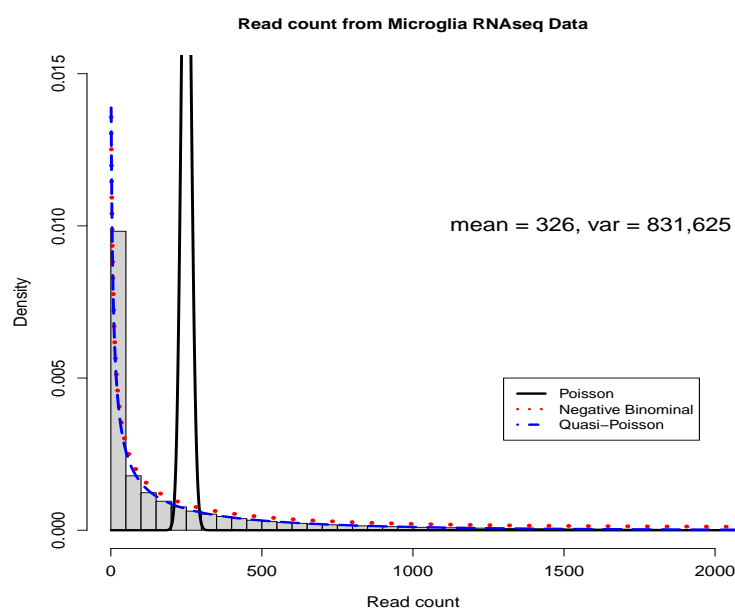


Figure 1: The overdispersion in the real RNA-seq count data

## 2 Methodology

We develop a new differentially expressed gene testing procedure to account for the heterogeneity in gene-wise overdispersion levels. Traditionally, Poisson and multinomial distributions are used to model count data with large variance. However, the variance of RNA sequence counts tends to



be much larger than that of the Poisson or multinomial distribution [30]. Overlooking the overdispersion could result in biased and misleading inference about gene association to the response of interest. To overcome this limitation, we first introduce adaptive distribution modeling in this paper to analyze the overdispersed RNA-seq count data. We utilize a quasi-Poisson distribution and a negative binomial distribution as the read count, thus generating a distribution similar to the overdispersion pattern which is based on empirical data. Specifically, we denote  $Y$  as the random count response, and the quasi-Poisson distribution satisfies:

$$E(Y) = \mu, ; \text{Var}(Y) = \theta\mu, \quad (1)$$

where  $\mu > 0$  is the mean of  $Y$ ,  $\theta \geq 1$  denotes the overdispersion parameter, and larger  $\theta$  indicates higher overdispersion level. Although  $\mu$  is larger than 0,  $Y$  can be any nonnegative integer. Note that Poisson model assumes that the variance is equal to the mean, e.g.,  $\theta = 1$ . In contrast, a quasi-Poisson distribution provides more flexibility to allow variance increases as a linear function of the mean. Accordingly, the quasi-Poisson regression generalizes the Poisson regression and is adopted to model an overdispersed count variable. The quasi-Poisson model is characterized by the first two moments, i.e., mean and variance. Besides the quasi-Poisson distribution, the negative binomial distribution can also be used to model overdispersed count data satisfying:

$$E(Y) = \mu, ; \text{Var}(Y) = \mu + \mu^2/\theta, \quad (2)$$

where  $\theta > 0$  is the overdispersion parameter, and smaller  $\theta$  indicates higher overdispersion level. Similar to the quasi-Poisson distribution, the negative binomial distribution is characterized by the mean and variance while modeling the variance as a quadratic function of the mean. In the following, we denote the quasi-Poisson distribution and negative distribution as  $\text{quasi-Poisson}(\mu, \theta^{QP})$  and  $\text{negative-binomial}(\mu, \theta^{NB})$ , respectively. In addition, we use  $NB$  and  $QP$  as the abbreviation of negative binomial and quasi-Poisson distribution. In Figure 1, we illustrate the distribution density functions by fitting the empirical read counts in our empirical data from microglia cells

(see Methods) with 1) Poisson distribution, 2) negative binomial distribution, and 3) quasi-Poisson distribution, respectively. Compared with the Poisson distribution, both the negative binomial and quasi-Poisson distributions provide better approximation by capturing the overdispersion in read counts.

In addition, the read counts of a gene can be affected by other factors in an experiment other than its expression level in the RNA-seq. Therefore, instead of directly modeling the raw count data  $Y$ , we first perform count normalization, which makes the expression levels of genes more comparable and accurate between samples. We utilize the Trimmed Mean of M-values normalization (TMM) [24] adopted by EdgeR to compute the normalization factors that correct sample-specific biases. TMM is recommended for most RNA-Seq data where most genes are not differentially expressed across any pairs of the samples. Specifically, we first calculate the normalization factors as the median ratio of gene counts relative to the geometric mean per gene within a specific sample. The normalization factors account for two main non-expression factors; e.g., sequencing depth and RNA composition before between-sample comparison [24]. Consequently, we divide raw counts by sample-specific size factors to yield the effective read count for cross-sample comparisons.

The proposed DEHOGT workflow combines the above ingredients to identify differentially expressed genes. Compared with the two popular RNA-seq analysis methods DESeq and EdgeR, the main difference of the proposed method is at the model fitting step of the above algorithm, where the overdispersion parameters  $\{\theta_i\}$  are estimated for each gene individually. The DESeq and EdgeR estimate the overdispersion parameters by pooling the samples from different genes under the assumption that genes with similar expression strength also share similar overdispersion levels. In contrast, the proposed method does not rely on the homogeneous dispersion assumption and can capture the heterogeneity in different genes' expression levels, especially when the overdispersion of gene is high. In addition, the proposed method allows one to choose different working distributions in Step 3 to model the RNA-seq count data to accommodate different associations between mean and variance presented in the empirical read count data. This provides us additional flexibility in modeling the overdispersion patterns to achieve more accurate read count

fitting. Consequently, correctly specified read count overdispersion patterns can lead to higher statistical power of post-hoc testing to detect differentially expressed genes.

We summarize the proposed method (DEHOGT) for the RNAseq read count for detecting differentially expressed (DE) genes as follows. Assume that there exists a total of  $R$  different treatments and  $S$  samples where each treatment has multiple samples as replicated measurements. We index the gene and sample measurements as  $g$  and  $s$  such that  $g = 1, 2, \dots, N$  and  $s = 1, 2, \dots, S$ . First, the read count data is modeled via one of the following generating distributions:

$$Y_{gs}/K_s \sim \text{quasi-Poisson}(\mu_{gs}, \theta_g^{QP}), Y_{gs}/K_s \sim \text{negative-binomial}(\mu_{gs}, \theta_g^{NB}),$$

where  $K_s$  denotes the normalization factor for the  $s$ th sample obtained by the TMM method. To determine the generating distribution, we check the overdispersion pattern between  $E_s(Y_{gs}/K_s)$  and  $\text{Var}_s(Y_{gs}/K_s)$  from the empirical data. A better quadratic function fitting leads to the choice of a negative binomial distribution and a better linear relation fitting leads to the quasi-Poisson. Here we assume that the gene-wise dispersion level is constant across all samples to estimate the quasi-Poisson distribution  $\theta^{QP}$ , or the negative binominal distribution  $\theta^{NB}$ , by utilizing information from samples under different treatments.

To differentiate genes' read counts under different treatments, we model the genewise read count mean via the following generalized linear model:

$$\log \mu_{gs} = \beta^{(1)}x_g + \beta_g^{(2)}\mathbf{T}_s^\top,$$

where  $\beta_g^{(2)} = (\beta_{g1}^{(2)}, \beta_{g2}^{(2)}, \dots, \beta_{gR}^{(2)})$  represents the fold change of the  $g$ th gene under  $R$  different treatments, and  $\mathbf{T}_s \in \{0, 1\}^R$  is the dummy coding for the treatment membership of the  $s$ th sample such that  $\mathbf{T}_{sr} = 1$  when the  $s$ th sample belongs to treatment  $r$ ,  $r = 1, \dots, R$ . In addition,  $x_g \in \mathbb{R}^p$  are the gene-wise covariates, so that our method can further adjust other non-expression factors to reduce the bias in inferring the genes' expression level, where  $R$  denotes a real number.

Given that  $\{\beta_g^{(2)}\}$  represents the gene-wise expression level under different treatments, we can

infer whether the  $g$ th gene is differentially expressed under the two treatments  $r_1$  and  $r_2$  based on the linear hypothesis testing  $H_0 : \beta_{gr_1}^{(2)} - \beta_{gr_2}^{(2)} = 0$ . Then we can identify the DE genes under the treatment comparison pair  $(r_1, r_2)$  when the corresponding p-value is smaller than a specific cutoff. To control the type-I error of simultaneously testing on multiple genes, we adopt the Benjamini-Hochberg procedure [3] to adjust the gene-wise p-value, and control the false discovery rate. In addition to the adjusted p-value, the magnitude of the logfold change is also suggested as another criterion for choosing DE genes with a logfold change of  $\log_2 |\beta_{gr_1}^{(2)} - \beta_{gr_2}^{(2)}|$  larger than 1.5 [20, 21]. Therefore, we combine these two criteria, and select the DE genes with an adjusted p-value smaller than 0.05 and an absolute logfold change larger than 1.5.

Our proposed DEHOGT algorithm is summarized as follows:

---

**Algorithm:** DEHOGT

---

1. (*Input*): For the  $i$ th gene  $i = 1, \dots, N$ , input read counts  $\{Y_{is}\}_{s=1}^S$  from  $S$  samples, the covariates  $x_i$  associated with the  $i$ th gene, and the treatment assignment for each sample  $T_s \in \{0, 1\}^R$  from each gene, ( $s = 1, \dots, S$  where  $R$  is the number of treatments). Specifying the working distribution indicator  $I$ :

$$I = \begin{cases} 1, & \text{choose the quasi-Poisson,} \\ 2, & \text{choose the negative binominal.} \end{cases}$$

2. (*Read count normalization*): Obtain normalization factor for the  $i$ th gene:  $K_i = \text{TMM}(\{Y_{is}\}_{s=1}^S)$  ( $i = 1, \dots, N$ ) where TMM denotes the Trimmed Mean of M-values normalization.
3. (*Fitting the generalized linear model*): For the  $i$ th gene, estimate the fold change parameter  $\beta_i^{(2)}$  and the overdispersion parameters  $\theta_i$ :

$$(\hat{\beta}_i^{(2)}, \hat{\theta}_i) = \underset{\beta, \theta}{\operatorname{argmax}} \prod_{s=1}^S f_I(Y_{is}/K_s, \mu_i(\beta), \theta_i), \quad (3)$$

$$\log \mu_i(\beta) = \beta^{(1)} x_i + \beta_i^{(2)} T_s^\top, \quad (4)$$

where  $f_I$  denotes the probability density function of the chosen working distribution.

4. (*Post-hoc testing*): For the  $i$ th gene and a specific interesting treatment pair  $(r_1, r_2)$ , perform

$$H_0 : \hat{\beta}_{gr_1}^{(2)} - \hat{\beta}_{gr_2}^{(2)} = 0,$$

and obtain  $p$ -value  $p_i$ .

5. (*DE gene filtering*): For the treatment pair  $(r_1, r_2)$ , obtain gene-wise adjusted  $p$  value, using Benjamini-Hochberg [3] adjusting for false positive discovery:

$$\{p_i^{adj}\}_{i=1}^N = \text{Benjamini-Hochberg}(\{p_i\}_{i=1}^N).$$

Select the  $i$ th gene if  $p_i^{adj} > 0.05$  and  $\log_2 |\beta_{ir_1}^{(2)} - \beta_{ir_2}^{(2)}| > 1.5$ .

6. (*Output*): Set of differentially expressed genes and the corresponding fold change estimation  $\hat{\beta}^{(2)}$ .

### 3 Simulation Studies

We compare the proposed DEHOGT method with two popular RNA-seq analysis methods DESeq [1] and EdgeR [23] in detecting differentially expressed genes on the simulated read count data and microglia cell RNA-seq data. In the first simulation setting, the discrepancy in expression level between the treatment and control group is weak for DE genes, while the average expression levels for both groups are high. In the second simulation setting, the expression discrepancy between the treatment and control group is strong for DE genes, while the average expression levels for both groups are low.

### 3.1 Read count with low discrepancy of expression level

In the first setting, we simulate the read count data following the negative binomial and the quasi-Poisson distribution:

$$Y_{gs} \sim QP \left( \text{mean} = \mu_{gs}, \text{var} = \mu_{gs} \theta_g^{\text{QP}} \right), \quad g = 1, \dots, N, \quad s = 1, \dots, |S|,$$

$$Y_{gs} \sim NB \left( \text{mean} = \mu_{gs}, \text{var} = \mu_{gs} (1 + \mu_{gs} / \theta_g^{\text{NB}}) \right), \quad g = 1, \dots, N, \quad s = 1, \dots, |S|,$$

where  $g \in \{1, \dots, N\}$  denotes gene indexes and the total number of genes  $N = 12,500$ . We use  $G^{\text{DE}} \subset \{1, \dots, N\}$  to denote the set of differentially expressed genes with  $|G^{\text{DE}}| = 2500$ . In addition,  $s \in S$  and  $|S| = 12$  denote the sample index with  $S = S_1 \cup S_2$ ,  $|S_1| = |S_2| = 6$ , where  $S_1$  and  $S_2$  indicate the samples in the control group and the treatment group, respectively. Here the mean parameters  $\mu_{gs}$  are similar to setting [29] in the RNA-seq data analysis. Specifically, the formulations are:

$$\mu_{gs} = E[Y_{gs}] = \begin{cases} M_{gs} + \eta_g, & g \in G^{\text{DE}}, s \in S_2 \\ M_{gs}, & o.w \end{cases},$$

where we sample  $M_{gs}$  from  $\text{Unif}[0, U_s]$ , and  $U_s \sim \text{Unif}[600, 800]$  is the sample-wise sequencing depth. Furthermore, we sample  $\eta_g$  from  $\exp(1/100)$  as the up-regulated signal of the differentially expressed genes. We consider three different overdispersion levels for the read counts from the quasi-Poisson distribution as

$$\theta_g^{\text{QP}} \sim \text{Unif}(1, 5), \theta_g^{\text{QP}} \sim \text{Unif}(5, 10), \theta_g^{\text{QP}} \sim \text{Unif}(10, 20),$$

where a larger  $\theta_g^{\text{QP}}$  indicates a greater overdispersion level. Similarly, we consider three overdispersion levels under the negative binomial read counts as

$$\theta_g^{\text{NB}} \sim \text{Unif}(0.1, 0.2), \theta_g^{\text{NB}} \sim \text{Unif}(0.2, 0.5), \theta_g^{\text{NB}} \sim \text{Unif}(0.5, 1),$$

where a smaller  $\theta_g^{\text{NB}}$  indicates a greater level of overdispersion.

We compare the performance of DESeq [1], EdgeR [23], and the proposed DEHOGT in identifying differentially expressed genes using an adjusted p value less than 0.05 and an absolute value of logfold change larger than 1.5. We first investigate the false negative rates from the comparison methods. The results under different data generations (quasi-Poisson or negative binomial) and different overdispersion levels are shown in Figures 2 and 3. Figures 2 and 3 suggest that the proposed DEHOGT method reaches the lowest false negative rate over competing methods under different overdispersion levels, indicating that most of the genes selected by the proposed method are differentially expressed. Note that the DEHOGT (NB) under the true negative binominal setting always achieves the lowest false negative rate when the cutoff of the adjusted p-value is set as 0.05. This is because the p-values from DEHOGT under NB tend to be smaller than for DEHOGT under QP. The better performance of DEHOGT under QP for the ROC and AUC (area under the ROC curve) implies that we can select a p-value cutoff larger than 0.05, under which the false negative rate of DEHOGT under QP can be smaller than the false negative rate of DEHOGT under NB.

We also investigate the overall DE gene discriminative power of different methods when the cutoff point of the adjusted p-value changes over the range from 0 to 1, as measured by the AUC (area under the ROC curve). Note that the AUC value is between 0 and 1, and a larger AUC value indicates that the algorithm can achieve an overall lower false positive rate and lower false negative rate simultaneously. The comparisons are shown in Figures 4 and 5, illustrating the AUC values for competing methods under different generating distributions and overdispersion levels.

The above results indicate that the proposed DEHOGT method outperforms both the DESeq and edgeR methods, and the proposed method can achieve the optimal AUC if the model is correctly specified. Specifically, DEHOGT under QP attains a higher AUC than DEHOGT under NB under varying  $\theta_g^{\text{QP}}$  when the read counts are generated from the quasi-Poisson distribution. Similarly, DEHOGT (NB) attains higher AUC than DEHOGT (QP) under varying  $\theta_g^{\text{NB}}$  if the read counts are generated from negative binomial distributions.

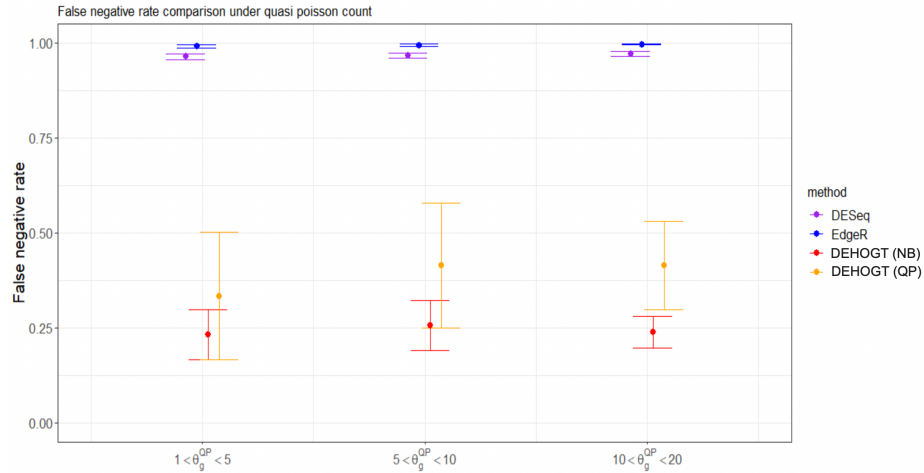


Figure 2: The false negative rate from different methods when the read counts follow the quasi-Poisson distribution with different overdispersion levels  $\theta_g^{QP}$  in simulation setting 1. The bars represents the standard deviation of the false negative rate over repeated experiments.

### 3.2 Read count with high discrepancy of expression level

In the second simulation setting, we simulate the read count data of the moderate overdispersion level in RNAseq read counts. Following the notations in simulation 1, we simulate the read count data from both the quasi-Poisson distribution and the negative binomial distribution as

$$Y_{gs} \sim QP \left( \text{mean} = \mu_{gs}, \text{var} = \mu_{gs} \theta_g^{QP} \right), \quad g = 1, \dots, N, \quad s = 1, \dots, |S|,$$

$$Y_{gs} \sim NB \left( \text{mean} = \mu_{gs}, \text{var} = \mu_{gs} (1 + \mu_{gs} / \theta_g^{NB}) \right), \quad g = 1, \dots, N, \quad s = 1, \dots, |S|,$$

where we choose  $N = 10,000$  and  $S = S_1 \cup S_2, |S_1| = |S_2| = 6$ . The GE genes are randomly selected and  $|G^{DE}| = 2,000$ . We consider three different overdispersion levels for the read counts from the quasi-Poisson distribution as

$$\theta_g^{QP} \sim \text{Unif}(50, 100), \quad \theta_g^{QP} \sim \text{Unif}(20, 50), \quad \theta_g^{QP} \sim \text{Unif}(10, 20).$$

Similarly, we also consider three overdispersion levels under negative binomial read counts as

$$\theta_g^{NB} \sim \text{Unif}(0.1, 0.2), \quad \theta_g^{NB} \sim \text{Unif}(0.2, 1), \quad \theta_g^{NB} \sim \text{Unif}(1, 2).$$



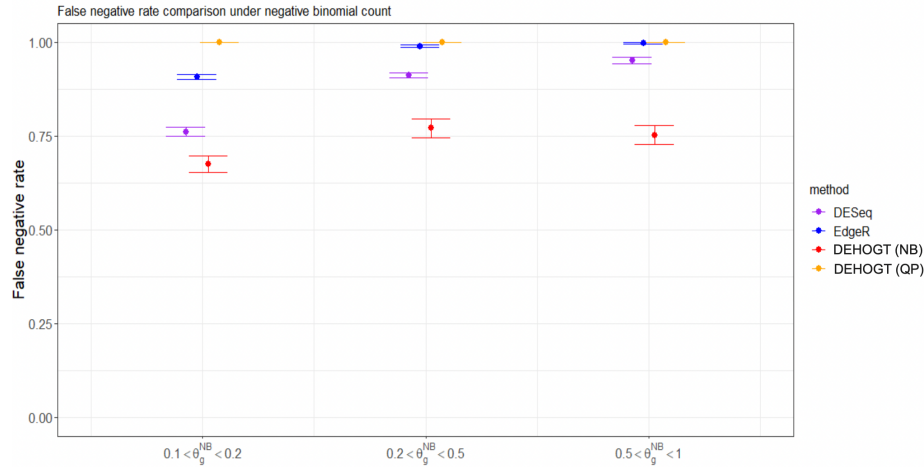


Figure 3: The false negative rate from different methods when the read counts follow the negative binomial distribution with different overdispersion levels  $\theta_g^{NB}$  in simulation setting 1. The bars represents the standard deviation of the false negative rate over repeated experiments.

We differentiate DE genes and non-DE genes with different sample means such that

$$\mu_{gs} \sim \text{Unif}(1, 500), s \in S, g \notin G^{DE},$$

$$\mu_{gs} \sim \begin{cases} \text{Unif}(1, 500), s \in S_1, g \in G^{DE} \\ \text{Unif}(1, 500) + \lceil 3.5 \times \bar{\mu}_{gS_1} \rceil, s \in S_2, g \in G^{DE} \end{cases}$$

where  $\lceil \cdot \rceil$  is the ceiling function, and  $\bar{\mu}_{gS_1} = \frac{1}{|S_1|} \sum_{s \in S_1} \mu_{gs}$ .

Notice that the expression discrepancy between the treatment and control group is strong for DE genes, while the average expression levels for both groups are low. To select the DE genes, we follow the selection criterion in the previous simulation such that the absolute value of log2fold change is larger than 1.5 and the adjusted p value is smaller than 0.05. We first investigate the false negative rates from different methods, and the results are illustrated in Figure 6 and 7.

The numerical results illustrates that the proposed method DEHOGT has a lower false negative rates than DESeq and EdgeR under different read count generation distributions and different overdispersion levels. Specifically, when the read count distribution is correctly specified, our method consistently achieves lower false negative rate than the EdgeR and DESeq. More importantly, the improvement from the DEHOGT increases as the degree of overdispersion in the read

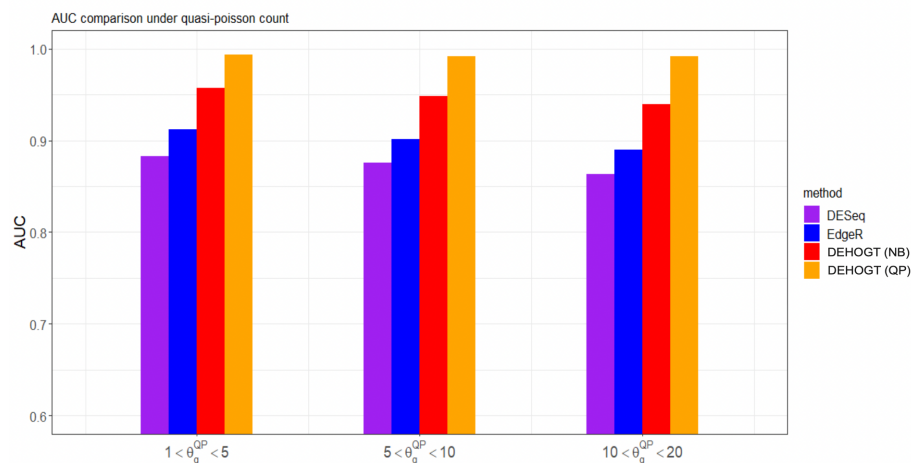


Figure 4: The area under the ROC curve from different methods when the read counts follow the quasi-Poisson distribution with different overdispersion levels  $\theta_g^{QP}$  in simulation setting 1.

count increases for both quasi-Poisson and negative binominal distributions.

We also investigated the overall discriminative power of the DE gene using different methods when the adjusted p-value cutoff varies between 0 and 1 instead of using 0.05. The overall classification performance is measured by the AUC. The Figures 9 and 8 illustrate the AUC from competing methods under different settings of read counts.

The above results show that the proposed DEHOGT method achieves a higher AUC in detecting the DE genes than the DESeq and EdgeR, indicating that our method offers a better balance between decreasing false positive rate and false negative rate. In addition, the improvement from our method is more significant as the overdispersion level increases, which is consistent with the aforementioned false negative rate comparison. A higher AUC from the DEHOGT method also implies that it can be more robust against the selection of different cutoff of p-value for DE genes.

We also illustrate the ROC curves in Figure 10 for two representative cases where read counts follow the negative binomial distribution with  $\theta^{NB} \in (1, 2)$ , and the quasi-Poisson distribution with  $\theta^{NB} \in (50, 100)$ , respectively.

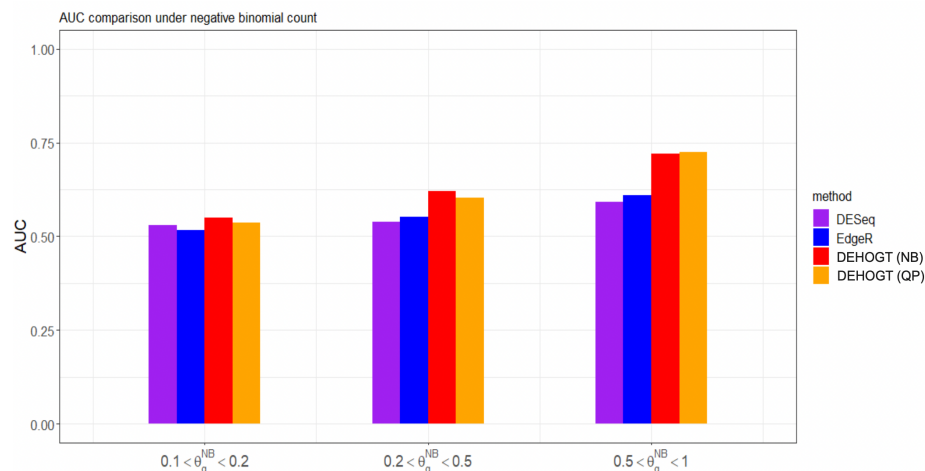


Figure 5: The area under the ROC curve from different methods when the read counts follow the negative binomial distribution with different overdispersion levels  $\theta_g^{NB}$  in simulation setting 1.

## 4 Application on Microglia RNA-seq read count data

In this section, we apply the proposed DEHOGT method, DESeq, and EdgeR in the study of post-traumatic stress disorder described in the Introduction section. Specifically, we aim to identify differentially expressed genes from microglia cells that are relevant to the PTSD progress. The RNA-seq data were collected by Uddin research team and Wildman lab at the University of the South Florida. The research performed in-vitro experiments on microglial cells which utilized stress hormones to imitate immune environments similar to PTSD. The function of stress hormones is to adjust the human interior environment, provide energy, and increase heart rate when experience stress [22]. The experiments exposed microglial cells to dexamethasone (dex) and hydrocortisone (cort) serving as stress hormones. The alcohol is also utilized as an additional control treatment to validate if changes in gene expressions are due to the exposure to stress hormones or just a random treatment (alcohol). Specifically, the experiments grew microglial cells under one of the four treatments: hydrocortisone, dexamethasone, alcohol (vehicle), or control. After exposure of three days, RNA-seq data was extracted from the cells on the third day and on the final day of the washout period (day 6), respectively. The goal of study is to identify the genes that are differentially expressed in microglia cells when exposed to different hormones and to determine if the dose of the hormone affects gene expression levels.

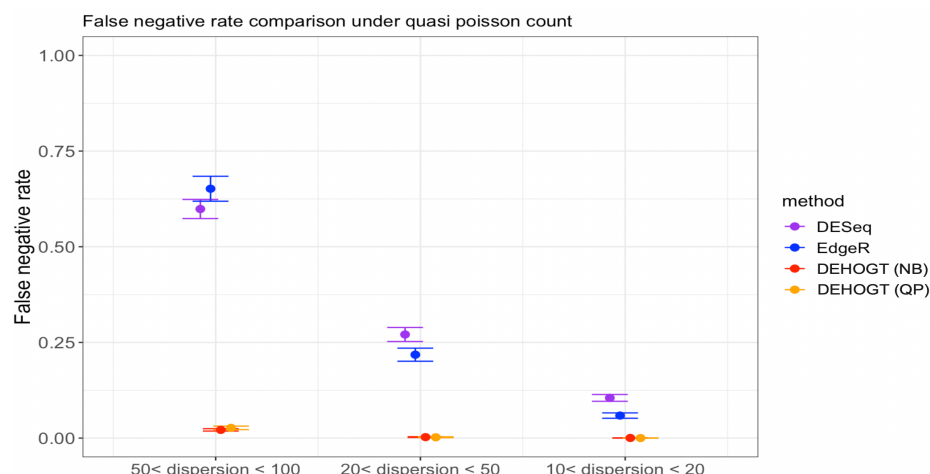


Figure 6: The false negative rate from different methods when the read count follow the quasi-Poisson distribution with different overdispersion levels  $\theta^{QP}$ . The variance of FNR obtained from repeated experiments is illustrated using the bars.

More specifically, there are a total of 20,052 expressed genes after quality control preprocessing. There is a total of 9 different treatments with the combination of media (dex, cort, vehicle, and control) and dosage (low and high): dex high, dex low, cort high, cort low, dex vehicle high, dex vehicle low, cort vehicle high, cort vehicle low, and control. On day 3 (time point 3), three repeated samples are collected under each treatment. On day 6 (time point 6), three repeated samples are collected under treatments dex high, dex low, cort high, and cort low, and one sample under dex vehicle high, dex vehicle low, cort vehicle high, and cort vehicle low.

We first investigated the level of empirical dispersion in the microglia RNA-seq read counts. Specifically, we examine the relation between sample count mean and sample count variance across all genes. Figure 11 illustrates a quadratic growth of count variance over count mean. In addition, we fit a quadratic regression on count variance over count mean, where an adjusted  $R^2$  coefficient reaches 0.66. Therefore, we choose to use a negative binomial distribution as the read counts generating process in the proposed DEHOGT method.

We utilize DEHOGT, DESeq, and EdgeR to select DE genes under the following 7 treatment comparison pairs: dex high at time point 3 and control (dexh3 vs control), dex high at time point 6 and control (dexh6 vs control), cort high at time point 3 and control (corth3 vs control), cort high at

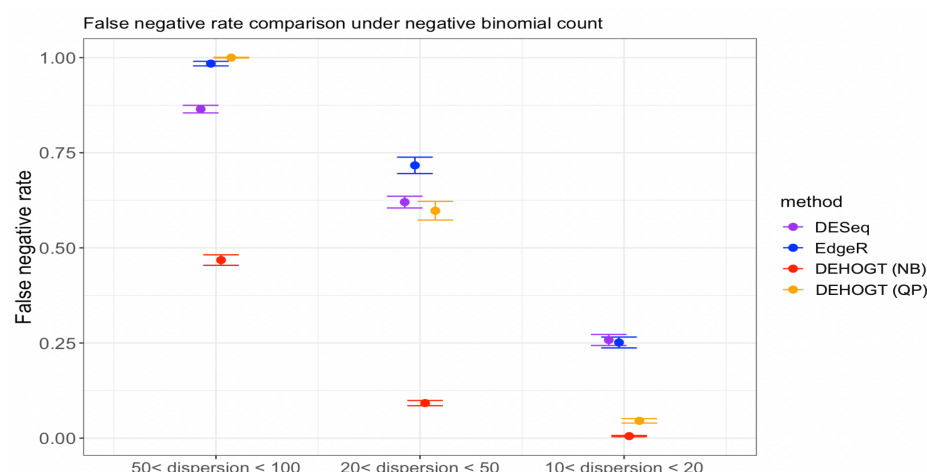


Figure 7: The false negative rate from different methods when the read count follows a negative binomial distribution with different overdispersion levels  $\theta^{NB}$ .

time point 6 and control (corth6 vs control), dex vehicle high and dex high at time point 3 (dexvh3 vs dexh3), dex vehicle low and dex low at time point 3 (dexvl3 vs dexl3), cort vehicle high and cort high at time point 3 (corth3 vs corth3). In selecting DE genes between the two treatments, we follow the criterion in Section 2 in that the adjusted p value is smaller than 0.05, and the log2fold change is larger than 1.5.

We first illustrate the number of DE genes selected by competing methods. Table 1 shows that the proposed method tends to select more genes than the other two methods, especially compared to the DESeq. In the exploratory stage, it is critical to include as many relevant genes as possible for the downstream analysis. The DEHOGT method is more effective in reducing the false negative rate in detecting PTSD-related genes by identifying a larger candidate pool of DE genes.

We conduct detailed analysis for the DE genes based on three methods for each treatment pair. In general, we investigate the overlapping in DE genes from three methods, where the findings are illustrated via the Venn diagram in Figure 12 to Figure 18. Notice that the proposed DEHOGT method selects more DE genes than DESeq and EdgeR for all pairwise comparisons between treatments except dexvh 3 vs dexh 3 and dexvl 3 vs dexl 3, demonstrating that the proposed method can identify more DE genes to reduce the potential risk of missing underlying relevant genes. In the following, we provide an interpretation for the treatment pair dexh6 versus control. The

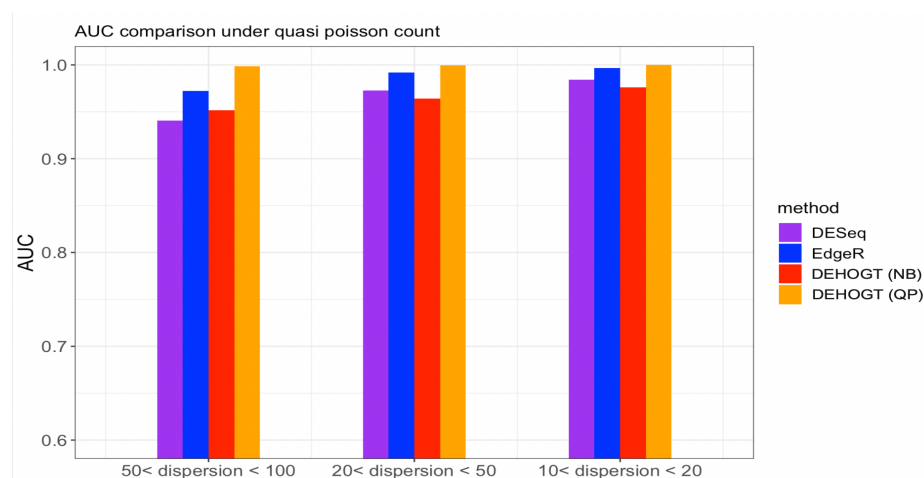


Figure 8: The AUC from different methods when the read count follows the quasi-Poisson distribution with different overdispersion levels  $\theta^{QP}$ . The variance of FNR obtained from repeated experiments is illustrated using the bars.

interpretation of other treatment pairs can be conducted similarly. The Venn diagram in Figure 13 shows that all the DE genes selected by the DESeq are also selected by EdgeR, and 86.7% of the DE genes selected by DESeq are also selected by DEHOGT. In addition, 61.7% of the DE genes selected by EdgeR are detected by DEHOGT.

The proposed method identifies three genes *CRISPLD2*, *TSC22D3*, and *PSGI* which are differentially expressed under the three treatment comparisons: dexvh 3 versus dexh 3, dexvl 3 versus dextl 3, and cortvh3 versus corth3. Specifically, the glucocorticoid-responsive gene *CRISPLD2* is found to be differentially expressed in read counts from an RNA-seq experiment with muscle cells exposed to dexamethasone [10]. The another glucocorticoid-responsive gene *TSC22D3* (*GILZ*) is found to be differentially expressed under gonorrhea or chlamydia exposure based on many animal and human gene studies that examine different cell types [7, 9]. These evidences support the fact that *TSC22D3* serves as a mediator for the anti-inflammatory activity of gonorrhea or chlamydia summarized in [26]. The gene *PSGI* is found to activate the underlying beta 1 (TGF- $\beta$ 1) known as transforming growth factor, which is an essential cytokine process in suppression and immunoregulation of inflammatory T cells [28, 5].

We also list the significant DE genes uniquely selected by the three methods in Table 2, which

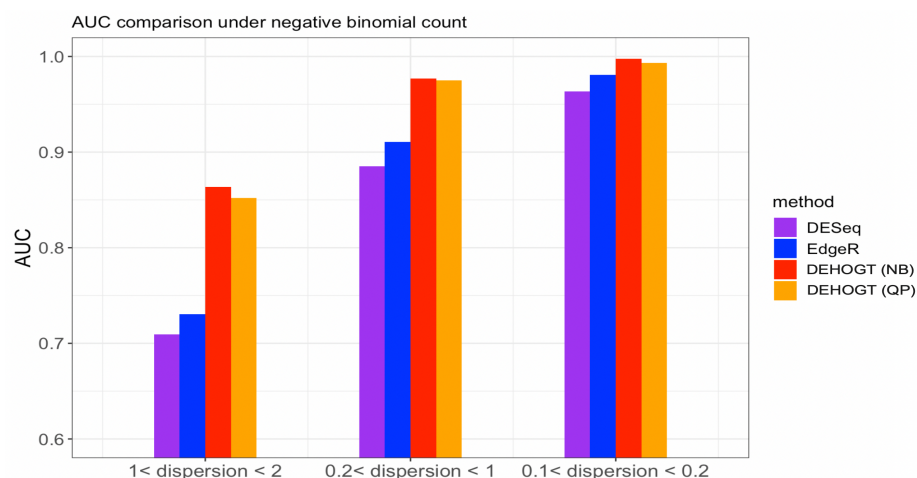


Figure 9: The AUC from different methods when the read counts follow the negative binomial distribution with different overdispersion levels  $\theta^{NB}$ .

demonstrates that most of the DE genes identified by DESeq are also selected by DEHOGT and EdgeR. Specifically, the gene *FKBP5* is identified by the proposed method but not identified by the other methods under the comparison dexh3 versus dexvh3. The gene *FKBP5* is a co-chaperone adjust the activity of glucocorticoid receptor. *FKBP5* is known as an important modulator of responding stress. In many studies using different cell types, the dysregulation phenomenon of *FKBP5* is found in many stress-related psychopathologies via investigating single nucleotide polymorphisms [4, 2], gene expression [11], and DNA methylation profiles [14].

In addition, we examine the most significant DE genes among the overlaps of the three methods in Figure 19 to Figure 24. For treatment pair dexh6 versus control, Figure 20 lists the 30 most significant DE genes which are overlapping for all three methods, and the bar charts with different colors represent the rank of p-values from the three methods. A shorter bar indicates a smaller p-value and therefore a more significantly differentially expressed genes under dex high and control comparison. The DEHOGT method selects genes *ROR1*, *FAT3*, *TLR4*, *CERNA2*, *ADPRHL1*, *NID2*, *CRISPLD2*, and *ABCA8* as the top 8 significant DE genes, and these genes are also among the top significant DE genes selected by EdgeR and DESeq. In general, our method provides a list of the top significant DE genes which is consistent with the DESeq and EdgeR in comparing dexh6 versus control. This cross-validation on DE genes via the three methods confirms the association



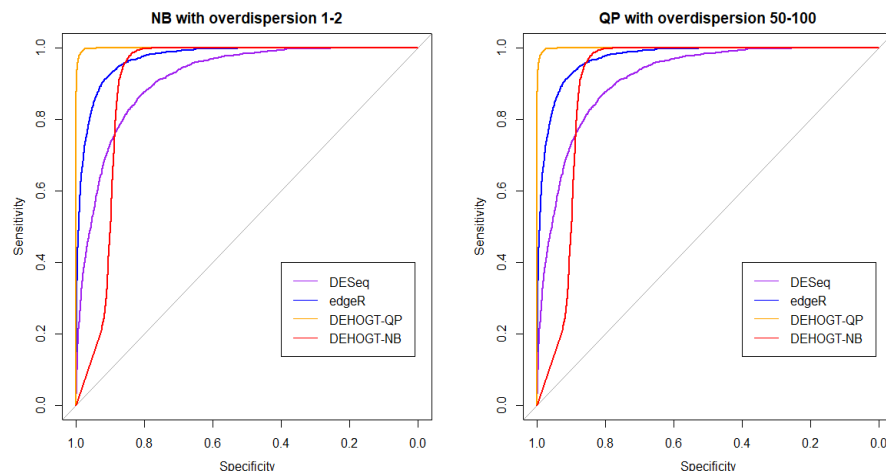


Figure 10: The ROC curve from different methods when the read counts follow the negative binomial distribution with  $\theta_g^{NB} \in (1, 2)$  and quasi-Poisson distribution with  $\theta_g^{QP} \in (50, 100)$ .

between PTSD and the top DE genes which are identified by the DEHOGT. In particular, the previously mentioned genes *TSC22D3* and *PSG1* are identified by all three methods for the vehicle treatment comparisons: dexvh 3 versus dexh 3, dexvl 3 versus dexl 3, and cortvh3 versus corth3. These results provide evidence of further need to explore their roles in the formulation of PTSD.

Methods	Treatment Pairs						
	dexh3 vs control	dexh6 vs control	corth3 vs control	corth6 vs control	dexvh3 vs dexh3	dexvl3 vs dexl3	cortvh3 vs corth3
DESeq	221	45	166	1	255	118	112
EdgeR	419	115	308	3	469	180	216
DEHOGT	981	237	355	582	383	148	256

Table 1: The number of selected DE genes from microglia RNA-seq data under different treatment comparisons.

## 5 Discussion

In this paper, we propose a revised differential expression analysis procedure DEHOGT for identifying differentially expressed genes based on overdispersed RNA-seq read count data. DEHOGT adopts a joint estimation of logfold changes that incorporates samples from all treatments simultaneously to utilize cross-treatment information. In addition, the proposed method takes advantage of within-treatment independence structures among genes to increase the effective sample size,



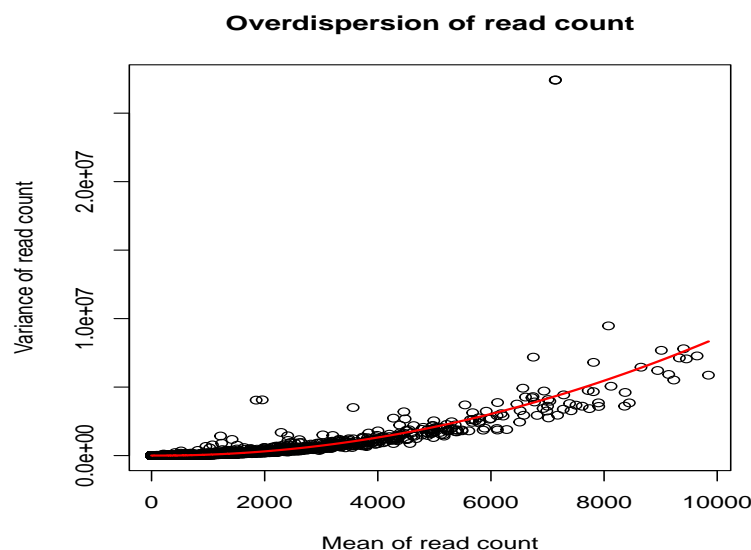


Figure 11: The dispersion of genewise RNA read counts. Each dot corresponds to a sample count from a specific gene.

which leads to stronger power in detecting DE genes. Furthermore, our method enjoys flexibility in utilizing different read count generating distributions instead of fixing only one negative binominal distribution as in the popular methods such as EdgeR and DESeq. This allows us to choose a generating distribution adopted to the empirical dispersion level. Therefore, DEHOGT has the potential to be applied for other genetic datasets with similar challenges of heterogeneous overdispersion levels.

In our simulation study, we demonstrate that the proposed method achieves better performance in detecting DE genes compared with the EdgeR and DESeq methods, especially when the per-treatment sample size is relatively small. The numerical experiments suggests that DEHOGT is less conservative in selecting DE genes due to adopting the individual fitting procedure. This property enables our method to have improved performance in controlling the false negative rate which is more critical for downstream analysis.

We further apply our method and compare it with EdgeR and DESeq on a real application in a microglia RNA-seq dataset collected by our team. Specifically, our method identifies more potential genes which may be potentially more relevant to PTSD than either EdgeR and DESeq.

Treatments	Methods		
	DEHOGT	EdgeR	DESeq
dexh3 vs control	BACE1, H19, LOC102724852, VSIR, SLC4A4 ITGB3, LTBP2, ELF1, EPS8, DUBR	PSG11, WNT7B, KLF15, SNORA74A, IGFN1 TNFRSF11B, EPB41L3, LCT-AS1, MYO5B, MMP1	
dexh6 vs control	SLC16A12, GCNT1, SLC7A2, FGD4, AOX1 DUBR, DGKI, NFKBIA, EMILIN3, KCNIP3 MARCHF1, TXNIP, IGFBP7, EPSTI1, PLXNA2	AMIG02, SRPX2, PAD1I, STAMBPL1, SHC3 KCNJ8, CYP24A1, KRT18P29, FAP, SOX3 LTBR, SLC28A3, TRHDE, LINC00402, MYO5B	
corth3 vs control	PRKACB, ABCB1, KLF9, SLC16A12, LTBP1 UGT2B7, MARCHF1, LAMA5, H19, LOC102724852 KCNIP3, COL5A1, NFKBIA, MEGF6, FGD4	RGMB, CTSC, FGF1, PSG4, ST6GAL2 FOSL1, SERPINE2, TRPC4, CREB5, NPPB ADAMTS1, ADAMTSL1, NLRP10, EVA1A, AGTR1	
corth6 vs control	ROR1, NID2, H3-2, ACTBL2, CHST2 EMILIN3, INAVA, RPEL1, PTGER2, MROH6 BBS12, CAMSAP3, TYSND1, TMEM116, MRPL38	NHS	
dexh3 vs dexvh3	GCNT1, ABCB1, SLC26A2, ITGB3, LTBP2 IGFBP7, COL5A1, DUBR, FKBP5, ADAMTS7 EFEMP1, PLXNA2, CPM, LPIN3, BIRC3	MAP3K7CL, QSOX1, SMURF2, SLC8A1, CLDN11 KHDRBS3, CREB5, GPRC5A, NR2F2, KRT17 ADAMTSS, NR2F2-AS1, HIF1A-AS3, SRPX2, KCNMA1	
dexl3 vs dexvl3	AGTR1, SCN9A, ABCB1, ELOVL6, SPSB1 DEPTOR, SPOCK1, NCAM2, CLDN1, KRT18P29 ITGA7, USP44, SRP14-DT, DPYD, STARD8	PLAUR, KRT17, AJAP1, MARCH, COL13A1 ITGB3, EPSTI1, COL4A4, DNER, TGFB3 HAS3, ANGPTL4, HCN3, ALPP, DNAH8	TGFB1, CPA4, CHST2, DGKI, KRT18P11 SNORA74A, LOC102724434, HNRNPA1P33, KRTAP4-8, RHEX SENCR, F8, PALMD, CCL26, KRTAP2-4
corth3 vs corth3	SLC26A2, MARCHF1, DOCK4, EIF1B-AS1, CREB5 EVA1A, PDZK1, COL13A1, FRMD6-AS1, FBXW4 CYP24A1, AJAP1, NUPR1, ATP6V1G2, MGST2	EFHD1, SLC1A3, SAA1, WFDC21P, ITGA1 CEBPD, GJD2, EGR2, KDR, GSX2 CCDC30, ASAH1-AS1, GGTLC3, LINC00886, LOC100506207	

Table 2: Top 15 unique DE genes unique selected by DEHOGT, EdgeR, and DESeq.

In addition, the cross-validation among EdgeR, DESeq and the proposed method provides a rich and robust candidate pool for genes relevant to PTSD. These results were obtained in the microglia dataset despite having issues of overdispersion and small sample size.

The popular existing methods DESeq and EdgeR identify differentially expressed genes by adopting an aggregate estimation strategy for read count overdispersion levels, which relies on the key assumption that genes with similar expression levels have similar overdispersion levels. The numerical results in this paper indicate that this assumption might be questionable under the scenario when heterogeneity of gene expression level is high. The violation of this assumption can undermine the detection power of methods based on aggregate estimators of overdispersion, especially when the overdispersion level is high. In contrast, estimating overdispersion levels for each gene separately can be more robust under high heterogeneity in gene expressions. On the other hand, the proposed independent estimation scheme integrates samples from different treatments instead from different genes, which might lose a certain amount of statistical testing power especially when the sample size is small. One direction worth of further exploration is to incorporate neighborhood similarity structures among genes such that the overdispersion estimation of a spe-

cific gene can borrow the information of samples from correlated genes, therefore we can increase the effective sample size for estimating overdispersion levels. A potential strategy could utilize gene-wise covariate variables or develop an adaptive fused-type penalty on gene overdispersion levels.

## 6 Appendix: Microglia cell experiment design

The microglial cell line HMC3 (ATCC CRL-3304, Manassas, Virginia) was used for in vitro experimentation following successful cell line authentication and Mycoplasma testing (Genetica, Burlington, NC). HMC3 cells (passage eight) were seeded in T-25 flasks with  $2 \times 10^5$  viable cells and incubated at 37°C and 5% CO<sub>2</sub>. After 24 hours, the growth medium in each T-25 was replaced with one of the following treatments: dexamethasone (1 or 0.01 µM), hydrocortisone (10 or 0.01 µM), vehicle (ethanol alcohol) or control (untreated media). Cells were incubated in treatment media for three days at 37°C and 5% CO<sub>2</sub> and imaged daily using the Axio Vert.A1 inverted microscope (Zeiss Oberkochen, Germany). At three days post-exposure (D3), cells were collected from each flask individually, quantified on the Countess II cell counter (Invitrogen Waltham, MA), seeded at  $2 \times 10^5$  viable cells/flask in new T-25 flasks with normal growth medium, and incubated for three additional days (i.e., washout period). The remaining D3 cell suspension for each flask was divided equally between two microcentrifuge tubes, pelleted and washed with PBS. One cell pellet per flask was placed in -80°C storage for future DNA extraction; the remaining cell pellet underwent RNA extraction using the RNeasy Mini Kit (QIAGEN, Hilden, Germany) protocol adapted for the QIAcube automated system (QIAGEN). On the final day of the washout period (D6), cells from each flask were imaged, collected in suspension and then quantified on the Countess II. Cell suspensions were split equally into two aliquots and then prepped for nucleic acid extraction as described for D3.

RNA samples from D3 and D6 were DNase treated (Dnase I kit; Sigma), quantified on the Qubit (RNA BR Assay Kit; Invitrogen) and scored for RNA integrity on the TapeStation (High

Sensitivity RNA ScreenTape; Agilent). Library preparation was performed following the Illumina TruSeq Stranded Total RNA Library Prep Kit protocol (Illumina, San Diego, CA) with TruSeq RNA Single Indexes (Set A and B; Illumina). Library quantity and quality were assessed using the Qubit 1X dsDNA HS Assay (Invitrogen), TapeStation High Sensitivity D1000 ScreenTape (Agilent), and using the KAPA Library Quantification Kit (Roche Basel, Switzerland) for the LightCycler 96 (Roche). RNA sequencing was conducted on the NextSeq 550 (Illumina) using the High Output Kit with 76 paired-end cycles (Illumina).

## 7 Acknowledgements

The authors would like to acknowledge the USF Genomics Core for their support of the microglia cell experiment.

## 8 Availability of data and materials

The datasets analysed during the current study are available in the NCBI's Gene Expression Omnibus (GEO) repository and are accessible through GEO Series accession number GSE219208 and link <https://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE219208>.

The experiment description and algorithm implementation are available via the following weblinks: <https://github.com/xiaobai0518/DEHOGT>. Operating systems: Windows, Linux, MacOS  
Programming language: R. Other requirements: RStudio. License: GPL-3.0.

## References

- [1] Anders, S. and Huber, W. (2010). Differential expression analysis for sequence count data. *Nature Precedings*, pages 1–1.
- [2] Appel, K., Schwahn, C., Mahler, J., Schulz, A., Spitzer, C., Fenske, K., Stender, J., Barnow, S., John, U., Teumer, A., et al. (2011). Moderation of adult depression by a polymorphism in

- the FKBP5 gene and childhood physical abuse in the general population. *Neuropsychopharmacology*, 36(10):1982–1991.
- [3] Benjamini, Y. and Hochberg, Y. (1995). Controlling the false discovery rate: a practical and powerful approach to multiple testing. *Journal of the Royal statistical society: series B (Methodological)*, 57(1):289–300.
- [4] Binder, E. B. (2009). The role of FKBP5, a co-chaperone of the glucocorticoid receptor in the pathogenesis and therapy of affective and anxiety disorders. *Psychoneuroendocrinology*, 34:S186–S195.
- [5] Blois, S. M., Sulkowski, G., Tirado-González, I., Warren, J., Freitag, N., Klapp, B. F., Rifkin, D., Fuss, I., Strober, W., and Dveksler, G. S. (2014). Pregnancy-specific glycoprotein 1 (PSG1) activates TGF- $\beta$  and prevents dextran sodium sulfate (DSS)-induced colitis in mice. *Mucosal Immunology*, 7(2):348–358.
- [6] Brewin, C. R., Andrews, B., and Valentine, J. D. (2000). Meta-analysis of risk factors for posttraumatic stress disorder in trauma-exposed adults. *Journal of Consulting and Clinical Psychology*, 68(5):748.
- [7] Cari, L., Ricci, E., Gentili, M., Petrillo, M. G., Ayroldi, E., Ronchetti, S., Nocentini, G., and Riccardi, C. (2015). A focused Real Time PCR strategy to determine GILZ expression in mouse tissues. *Results in Immunology*, 5:37–42.
- [8] Chambers, J. M. (2008). Software for data analysis: programming with R. *Springer*, 2(1).
- [9] Franco, L. M., Gadkari, M., Howe, K. N., Sun, J., Kardava, L., Kumar, P., Kumari, S., Hu, Z., Fraser, I. D., Moir, S., et al. (2019). Immune regulation by glucocorticoids can be linked to cell type-dependent transcriptional responses. *Journal of Experimental Medicine*, 216(2):384–406.
- [10] Himes, B. E., Jiang, X., Wagner, P., Hu, R., Wang, Q., Klanderman, B., Whitaker, R. M., Duan, Q., Lasky-Su, J., Nikolos, C., et al. (2014). RNA-Seq transcriptome profiling identifies *crispld2* as a glucocorticoid responsive gene that modulates cytokine function in airway smooth muscle cells. *PloS One*, 9(6):e99625.
- [11] Ising, M., Maccarrone, G., Brückl, T., Scheuer, S., Hennings, J., Holsboer, F., Turck, C. W.,

- Uhr, M., and Lucae, S. (2019). FKBP5 gene expression predicts antidepressant treatment outcome in depression. *International Journal of Molecular Sciences*, 20(3):485.
- [12] Kessler, R. C., Aguilar-Gaxiola, S., Alonso, J., Benjet, C., Bromet, E. J., Cardoso, G., Degenhardt, L., de Girolamo, G., Dinolova, R. V., Ferry, F., et al. (2017). Trauma and PTSD in the WHO world mental health surveys. *European Journal of Psychotraumatology*, 8(sup5):1353383.
- [13] Kim, G. S., Smith, A. K., Xue, F., Michopoulos, V., Lori, A., Armstrong, D. L., Aiello, A. E., Koenen, K. C., Galea, S., Wildman, D. E., et al. (2019). Methylomic profiles reveal sex-specific differences in leukocyte composition associated with post-traumatic stress disorder. *Brain, Behavior, and Immunity*, 81:280–291.
- [14] Klengel, T., Mehta, D., Anacker, C., Rex-Haffner, M., Pruessner, J. C., Pariante, C. M., Pace, T. W., Mercer, K. B., Mayberg, H. S., Bradley, B., et al. (2013). Allele-specific FKBP5 DNA demethylation mediates gene–childhood trauma interactions. *Nature Neuroscience*, 16(1):33–41.
- [15] Lowe, S. R., Galea, S., Uddin, M., and Koenen, K. C. (2014). Trajectories of post traumatic stress among urban residents. *American Journal of Community Psychology*, 53(1):159–172.
- [16] Mardis, E. R. (2008). The impact of next-generation sequencing technology on genetics. *Trends in Genetics*, 24(3):133–141.
- [17] Mehta, D., Gonik, M., Klengel, T., Rex-Haffner, M., Menke, A., Rubel, J., Mercer, K. B., Pütz, B., Bradley, B., Holsboer, F., et al. (2011). Using polymorphisms in FKBP5 to define biologically distinct subtypes of posttraumatic stress disorder: evidence from endocrine and gene expression studies. *Archives of General Psychiatry*, 68(9):901–910.
- [18] Mills, K. L., McFarlane, A. C., Slade, T., Creamer, M., Silove, D., Teesson, M., and Bryant, R. (2011). Assessing the prevalence of trauma exposure in epidemiological surveys. *Australian & New Zealand Journal of Psychiatry*, 45(5):407–415.
- [19] Mortazavi, A., Williams, B. A., McCue, K., Schaeffer, L., and Wold, B. (2008). Mapping and quantifying mammalian transcriptomes by RNA-Seq. *Nature Methods*, 5(7):621–628.

- [20] Patterson, T. A., Lobenhofer, E. K., Fulmer-Smentek, S. B., Collins, P. J., Chu, T.-M., Bao, W., Fang, H., Kawasaki, E. S., Hager, J., Tikhonova, I. R., et al. (2006). Performance comparison of one-color and two-color platforms within the MicroArray Quality Control (MAQC) project. *Nature Biotechnology*, 24(9):1140–1150.
- [21] Peart, M. J., Smyth, G. K., Van Laar, R. K., Bowtell, D. D., Richon, V. M., Marks, P. A., Holloway, A. J., and Johnstone, R. W. (2005). Identification and functional significance of genes regulated by structurally different histone deacetylase inhibitors. *Proceedings of the National Academy of Sciences*, 102(10):3697–3702.
- [22] Ranabir, S. and Reetu, K. (2011). Stress and hormones. *Indian Journal of Endocrinology and Metabolism*, 15(1):18.
- [23] Robinson, M. D., McCarthy, D. J., and Smyth, G. K. (2010). edgeR: A bioconductor package for differential expression analysis of digital gene expression data. *Bioinformatics*, 26(1):139–140.
- [24] Robinson, M. D. and Oshlack, A. (2010). A scaling normalization method for differential expression analysis of RNA-Seq data. *Genome biology*, 11(3):1–9.
- [25] Robinson, M. D. and Smyth, G. K. (2007). Moderated statistical tests for assessing differences in tag abundance. *Bioinformatics*, 23(21):2881–2887.
- [26] Ronchetti, S., Migliorati, G., and Riccardi, C. (2015). GILZ as a mediator of the anti-inflammatory effects of glucocorticoids. *Frontiers in Endocrinology*, 6:170.
- [27] Sarapas, C., Cai, G., Bierer, L. M., Golier, J. A., Galea, S., Ising, M., Rein, T., Schmeidler, J., Müller-Myhsok, B., Uhr, M., et al. (2011). Genetic markers for PTSD risk and resilience among survivors of the World Trade Center attacks. *Disease Markers*, 30(2-3):101–110.
- [28] SNYDER, S. K., WESSELLS, J. L., WATERHOUSE, R. M., DVEKSLER, G. S., WESSNER, D. H., WAHL, L. M., and ZIMMERMANN, W. (2001). Pregnancy-specific glycoproteins function as immunomodulators by inducing secretion of IL-10, IL-6 and TGF- $\beta$ 1 by human monocytes. *American Journal of Reproductive Immunology*, 45(4):205–216.
- [29] Sonesson, C. and Delorenzi, M. (2013). A comparison of methods for differential expression

- analysis of RNA-Seq data. *BMC Bioinformatics*, 14(1):1–18.
- [30] Trapnell, C., Hendrickson, D. G., Sauvageau, M., Goff, L., Rinn, J. L., and Pachter, L. (2013). Differential analysis of gene regulation at transcript resolution with RNA-Seq. *Nature Biotechnology*, 31(1):46–53.
- [31] Yehuda, R. (2002). Post-traumatic stress disorder. *New England Journal of Medicine*, 346(2):108–114.
- [32] Yehuda, R., Cai, G., Golier, J. A., Sarapas, C., Galea, S., Ising, M., Rein, T., Schmeidler, J., Müller-Myhsok, B., Holsboer, F., et al. (2009). Gene expression patterns associated with post traumatic stress disorder following exposure to the World Trade Center attacks. *Biological Psychiatry*, 66(7):708–711.
- [33] Yehuda, R. and LeDoux, J. (2007). Response variation following trauma: a translational neuroscience approach to understanding PTSD. *Neuron*, 56(1):19–32.
- [34] Zhang, H., Pounds, S. B., and Tang, L. (2013). Statistical methods for overdispersion in mRNA-Seq count data. *The Open Bioinformatics Journal*, 7(1).
- [35] Zohar, J., Fostick, L., Cohen, A., Bleich, A., Dolfin, D., Weissman, Z., Doron, M., Kaplan, Z., Klein, E., and Shalev, A. Y. (2009). Risk factors for the development of posttraumatic stress disorder following combat trauma: A semiprospective study. *The Journal of Clinical Psychiatry*, 70(12):18399.



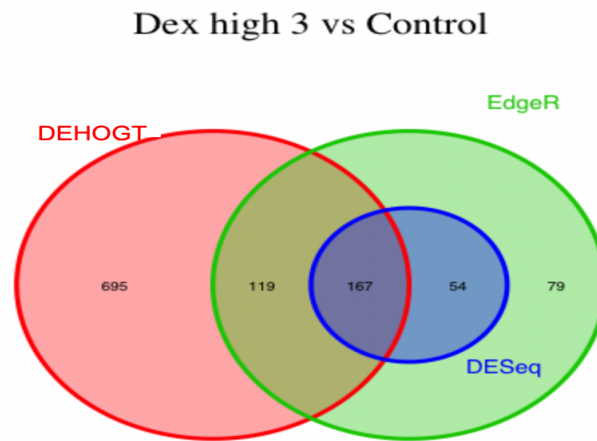


Figure 12: The selected DE genes from DEHOGT, DESeq, EdgeR under treatment comparison dexh3 versus control.

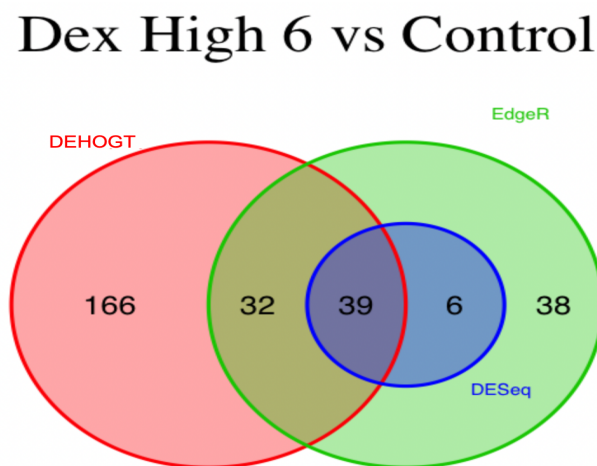


Figure 13: The selected DE genes from DEHOGT, DESeq, EdgeR under treatment comparison dexh6 versus control.

## Cort high 3 vs Control

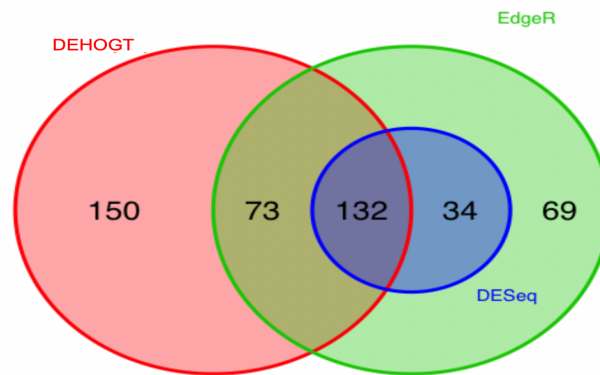


Figure 14: The selected DE genes from DEHOGT, DESeq, EdgeR under treatment comparison corth3 versus control.

## Cort high 6 vs Control

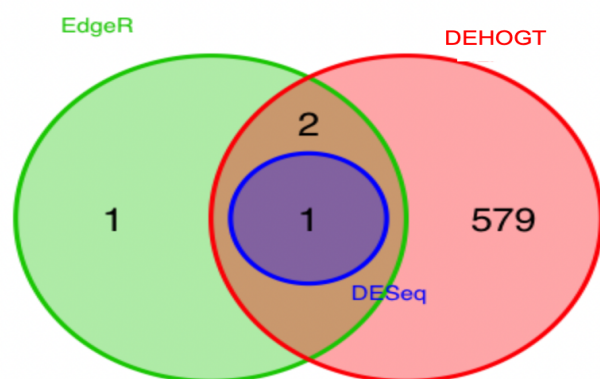


Figure 15: The selected DE genes from DEHOGT, DESeq, EdgeR under treatment comparison corth6 versus control.

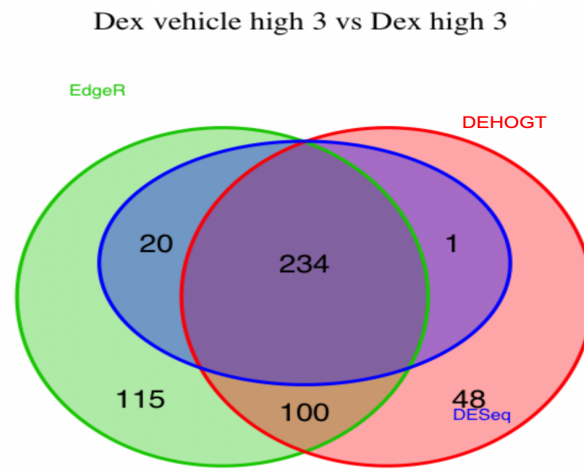


Figure 16: The selected DE genes from DEHOGT, DESeq, EdgeR under treatment comparison dexvh3 versus dexh3.

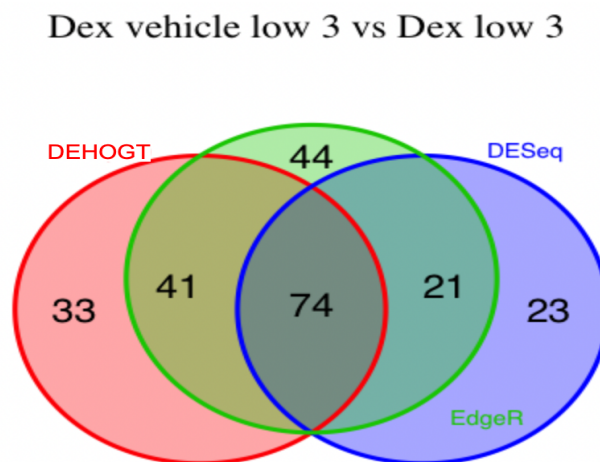


Figure 17: The selected DE genes from DEHOGT, DESeq, EdgeR under treatment comparison dexvl3 vs dexl3.

### Cort vehicle high 3 vs Cort high 3

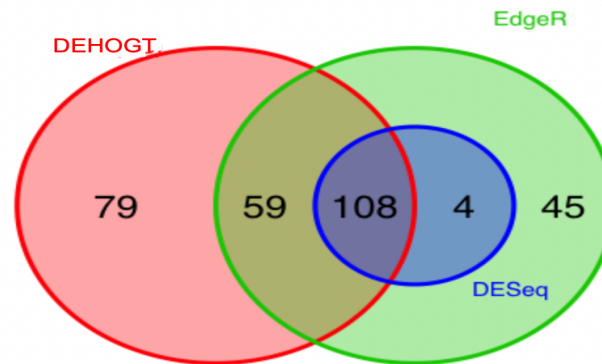


Figure 18: The selected DE genes from DEHOGT, DESeq, EdgeR under treatment comparison cortvh3 vs corth3.

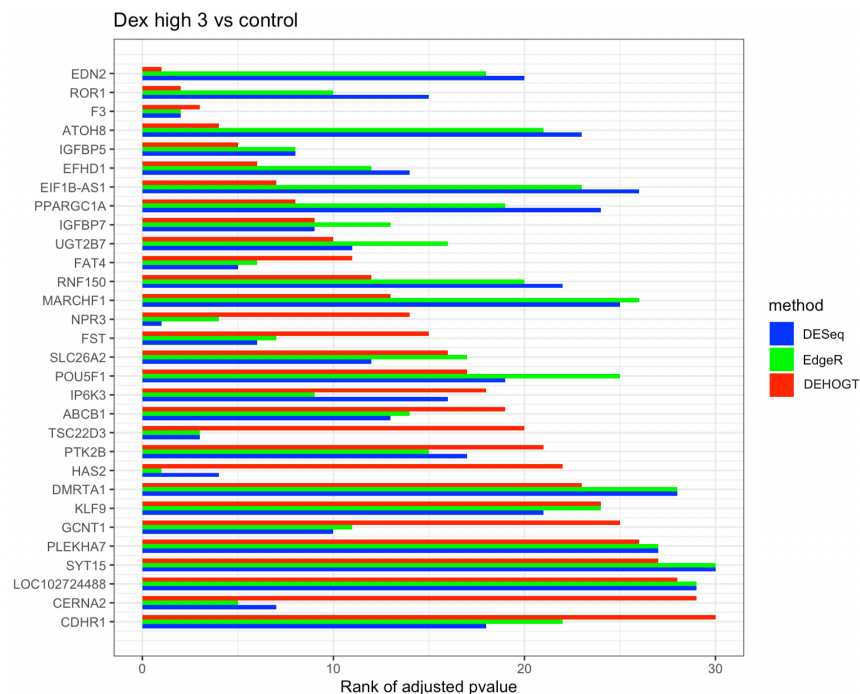


Figure 19: The rank of p-value of selected genes under treatment comparison dexh3 versus control, a shorter bar indicates a smaller p-value (more significantly differently expressed).

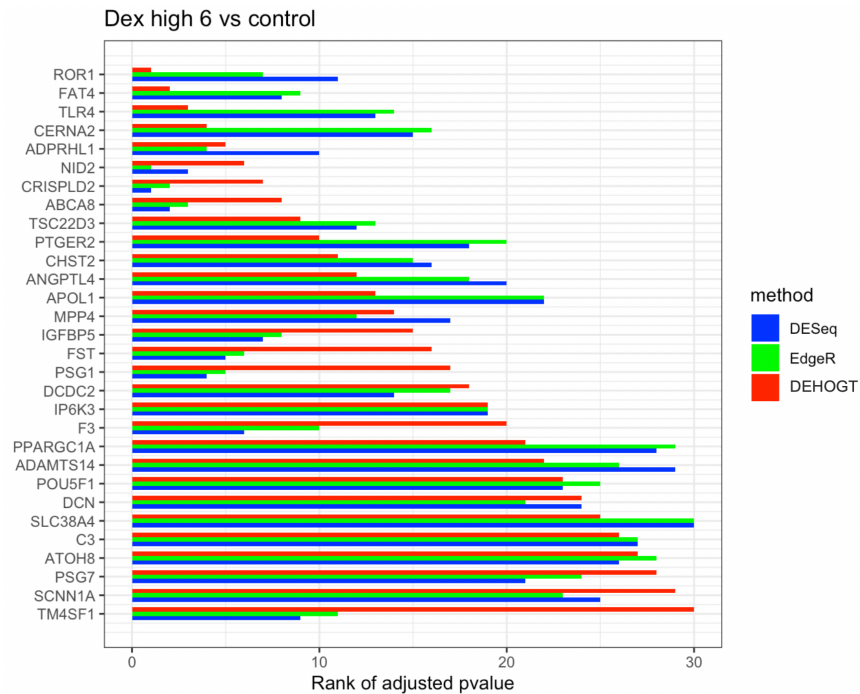


Figure 20: The rank of p-value of selected genes under treatment comparison dexh6 versus control, and shorter bar indicates a smaller p-value.

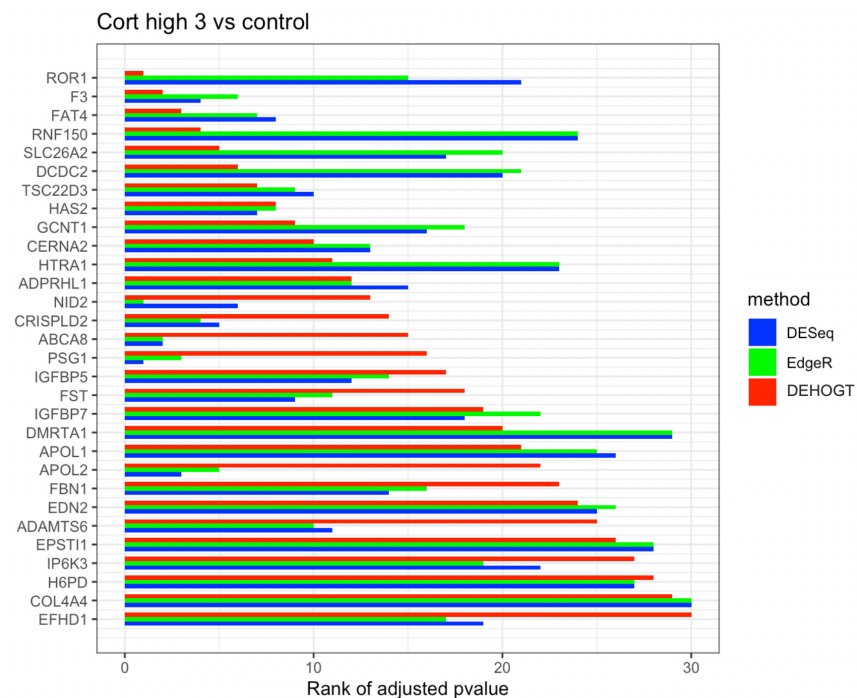


Figure 21: The rank of p-value of selected genes under treatment comparison corth3 versus control, and shorter bar indicates a smaller p-value.

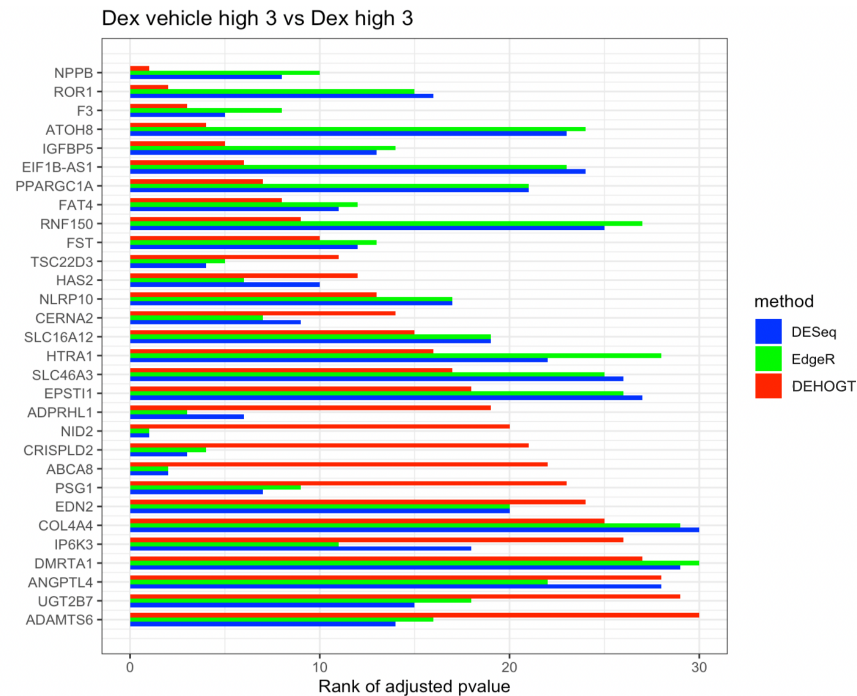


Figure 22: The rank of p-value of selected genes under treatment comparison dexvh versus dexh, and shorter bar indicates a smaller p-value.

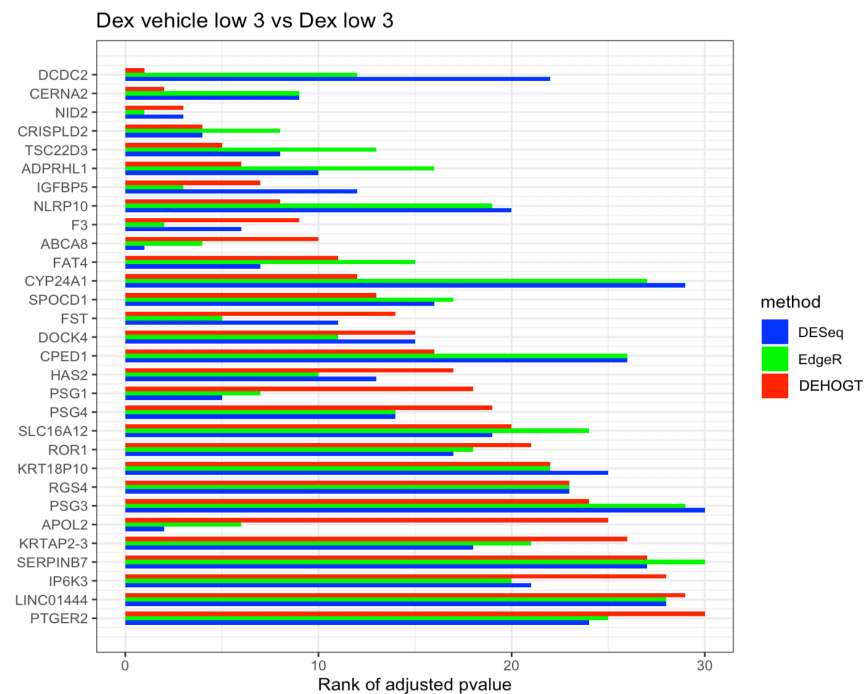


Figure 23: The rank of p-value of selected genes under treatment comparison dexvl versus dexl, and shorter bar indicates a smaller p-value.

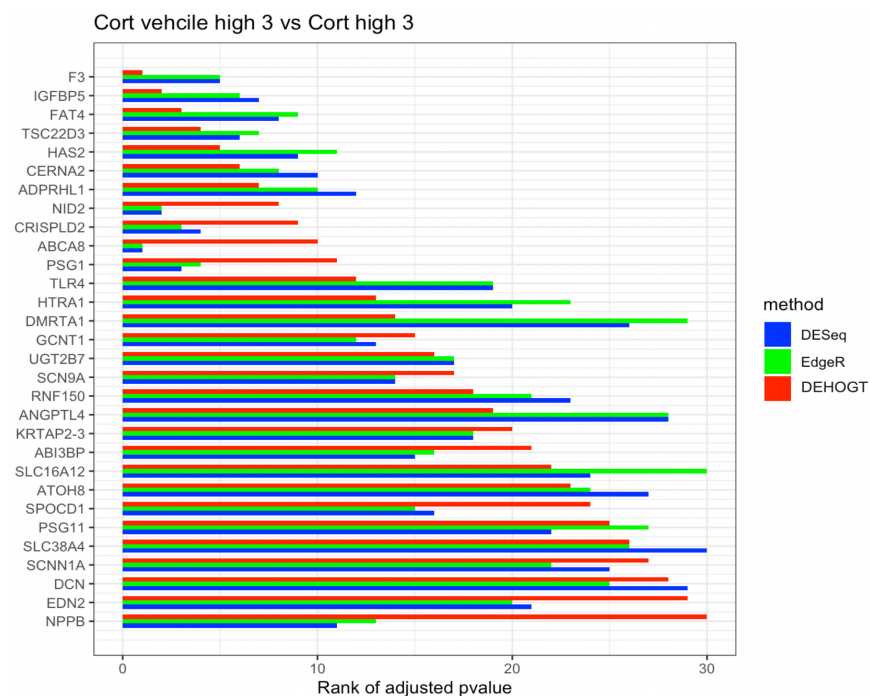


Figure 24: The rank of p-value of selected genes under treatment comparison cortvh versus cortth, and shorter bar indicates a smaller p-value.