# Testing Multiple Hypotheses through IMP Weighted FDR Based on a Genetic Functional Network with Application to a New Zebrafish Transcriptome Study

Jiang Gui*, Walter Taylor, Jason H. Moore
Institute for Quantitative Biomedical Sciences
Dartmouth College
Hanover, NH, US

Con Sullivan, Carol H. Kim
Department of Molecular and
Biomedical Sciences  University
of Maine
Orono, ME, US

Casey S. Greene*
Department of Genetics
Geisel School of Medicine at
Dartmouth
Hanover, NH, US

*Abstract—* **In genome-wide studies, thousands of hypothesis tests are carried out at the same time. Bonferroni correction and False Discovery Rate (FDR) can effectively control type I error but often yield a high false negative rate. We aim to develop a more powerful method to detect differential expressed genes. We present an Weighted False Discovery Rate (WFDR) method that incorporate biological knowledge from genetic networks. We first identify weights using Integrative Multi-species Prediction (IMP) and then apply the weights in WFDR to identify differentially expressed genes through a IMP-WFDR algorithm. We conducted a simulation study to characterize the performance of this method. We performed genomic characterization to identify potential synergistic and antagonist interactions between the highly-conserved zebrafish cftr gene and the environmental toxicant arsenic, particularly in the context of a systemic infection with Pseudomonas aeruginosa. Zebrafish were exposed to arsenic at 10 parts per billion and/or infected with P. aeruginosa. Appropriate controls were included. We then applied IMP-WFDR during the analysis of differentially expressed genes. We compared the mRNA expression for each group and found over 200 differentially expressed genes and several enriched pathways including defense response pathways, arsenic response pathways, and the Notch signaling pathway.**

*Keywords— false discovery rate; family-wise error rate; genomic studies; data integration*

## INTRODUCTION

With the rapid development of novel high-throughput deep sequencing technology, the study of functional transcriptomes has changed dramatically. Compared to microarray based measurements, RNA-sequencing (RNA-seq) technology can profile RNA transcript abundance within greater depth and accuracy. It can effectively detect alternative splicing variants and novel transcripts and does not require an assembled genome sequence. Study [1] compared high-throughput sequencing data in Illumina and Affymetrix platform and showed that sequencing data are highly replicable, with much smaller technical variation when compared to microarray data. In many cases, it is sufficient to sequence each mRNA sample only once.

RNA-seq analysis raises significant challenges when tens of thousands of transcripts are tested at the same time in order to find those that are differentially expressed under certain conditions. Multiple testing correction is required to adequately control type I error. Bonferroni correction [2] effectively controls the family-wise error rate (FWER), but it is too strict for most uses and has been utilized primarily in studies where only a few null hypotheses are expected to be rejected. In the context of high-dimensional data analysis, using a procedure that guards against any single false positive occurring can lead to many missed findings i.e. increased Type II error rates. Benjamini and Hochberg [3] proposed a procedure to control the false discovery rate (FDR), which is defined as the proportion of null hypotheses that are rejected erroneously. This criterion is less stringent than equivalent FWER-based procedures and provides a useful compromise between the loss of power attributable to the Bonferroni correction and the lack of control of Type I errors associated with comparisons unadjusted for multiplicity. Much additional research has been done on this approach, including the q-value method by Storey [4] as a generalization of the p-value to the FDR setting, and the local FDR introduced by Efron et al. [5-7]. Because of its strengths, the FDR method has been widely applied to microarray analyses to detect differentially expressed genes, and is accordingly incorporated into popular software packages, e.g. SAM (Significance Analysis of Microarrays) and LIMMA (Linear Models for Microarray Data) in R.

Although it improves Type II error rates relative to FWER-based methods, the FDR method lacks power when there are a large number of tests. To further improve power, Genovese et al. [8] proposed extending the FDR method to include weighted p-values and proved that as long as the sum of weights equals the total number of tests, the weighted method still effectively controls the FDR at the nominal level (i.e. with WFDR at a nominal level of 0.05, five percent of the

discoveries would be false positives just as with the unweighted FDR). In the field of biology, researchers have access to substantial information about biology from both published experiments and existing biological data. Leveraging this knowledge to improve the power of genomic tests would provide significant advantages; therefore, methods for p-value weighting have become an active and important research area.

The Integrative Multi-species Prediction [9, 10] webserver has collected over 2,000 microarray datasets from the NCBI Gene Expression Omnibus (GEO), large collection of biochemical experiments from the Database of Protein and Genetic Interactions (BioGRID), tissue-specificity and phenotype characterizations from The Zebrafish Model Organism Database (ZFIN) and other sources for seven organisms. IMP integrates these data into functional relationship networks using naïve Bayesian classifiers. Functional networks represent an efficient and comprehensive summary of existing publicly available data. Specifically, in each functional network, each node represents a gene and each edge between two genes is the posterior probability that the two genes are involved in the same process or pathway [11]. IMP is implemented on a web server (http://imp.princeton.edu/) that enables investigators to analyze their experimental results in the functional context of gene-gene networks from multiple organisms. IMP has previously been used to direct functional experiments by identifying novel gene participants in a pathway or process [10]. While these networks are valuable for guiding targeted experiments, methods that allow us to use these networks to effectively analyze new large-scale experiments would address a currently unmet need.

In this paper, we perform a simulation to compare power among weighted FDR, Benjamin and Hochberg's (BH) procedure for FDR control and Bonferroni correction. We develop and evaluate a novel IMP-WFDR method that uses state-of-the-art genetic network information to infer appropriate weights for WFDR and adjusts the p-values from simultaneous tests to correct for multiple testing issue.

Cystic fibrosis (CF) is a multi-system genetic disorder caused by a mutation in the cystic fibrosis transmembrane conductance regulator (CFTR) gene. Eighty percent of CF patients develop chronic infections with P. aeruginosa [12]. CF is a complex disease undoubtedly affected by the environment. In particular, there is emerging evidence linking CFTR function with the environmental toxicant arsenic [13, 14]. Using our zebrafish model for CF, we have shown that cftr and arsenic each influence aspects of innate immunity [15]. In the current study, we are establishing linkages that show how arsenic affects the innate immune response through Cftr

function. To accomplish this objective, we deep sequenced the transcriptomes of zebrafish that lack a functional Cftr, were exposed to arsenic at environmentally relevant doses (2 ppb and 10 ppb), and/or were infected by P. aeruginosa, along with the appropriate controls. We apply our method to a portion of this RNA-seq dataset to detect genes that were differentially expressed upon exposure to arsenic and/or infection with P. aeruginosa. Our proposed IMP-WFDR method integrates large data compendia through functional networks and focuses these data on the analysis of a new genome-scale experiment to better identify relevant candidates.

METHOD

A. *Simulation Study*

To demonstrate the efficacy of the weighted FDR approach, we simulated 1,000,000 test statistics from a mixture of two normal distributions:

$$\mathbf{T} \sim (\mathbf{1} - \boldsymbol{\pi}) * \mathbf{N}(\mathbf{0}, \mathbf{1}) + \boldsymbol{\pi} * \mathbf{N}(\mathbf{2.5}, \mathbf{1}) \qquad (1)$$
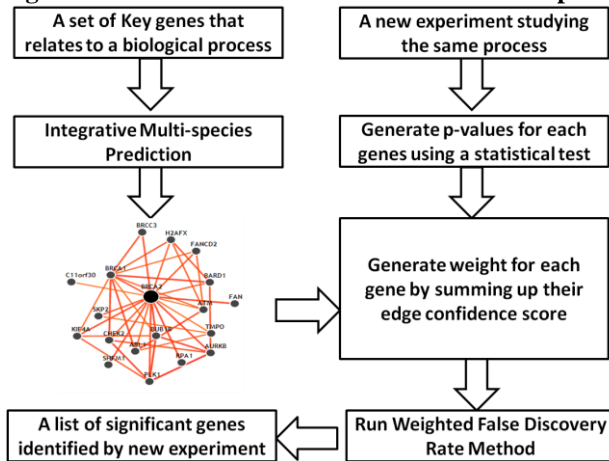
We assume that null statistics are generated from the standard normal distribution and the alternative statistics are from a normal distribution with a mean of 2.5. Here we set $\boldsymbol{\pi} = \mathbf{0.001}$ to indicate that most of the test statistics are generated from null distribution. We calculate the p-values using an inverse normal calculator and run Bonferroni, BH and weighted FDR adjustment at the 0.05 level. We assigned a weight of 10 to p-values from the alternative and 0.99 to those from the null. This demonstrates the impact of appropriate weighting on the test.

B. *IMP-WFDR Method*

Genovese et al. [8] proposed a simple weighted FDR procedure in which non-negative weights $W_i$ are assigned to each p-value such that $\sum_{i=1}^{m} W_i = m$. The BH procedure is applied directly to $Q_i = P_i/W_i$, where $P_i$ denotes the unweighted p-value. They have proven that weighted FDR controls FDR at the nominal level.

Given the freedom of assigning weights to p-values, weighted FDR provides an excellent platform to incorporate expert biological knowledge of each gene's biological function. Because we already have substantial knowledge about biology, we can use this procedure to most strongly focus on genes that represent likely candidates. If we assign weights only to the strongest candidates (i.e. those with known literature involvement), we limit our ability to make novel discoveries. If we assign weights evenly to all genes, we do not provide adequate weight to genes with some support. Consequently functional networks, which integrate publicly

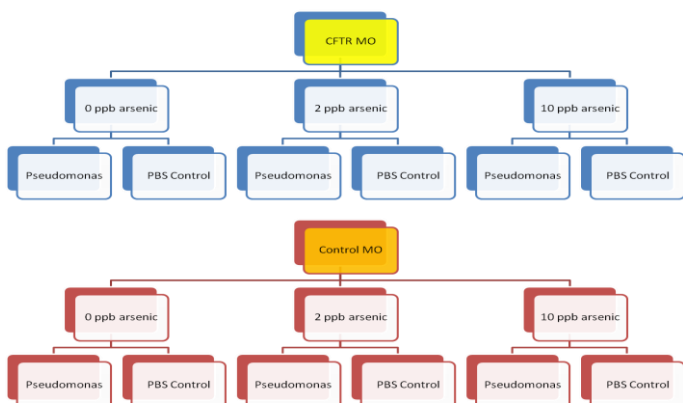**Figure 1 Demonstration of WFDR-IMP procedure**



available data using biological knowledge, provide both the opportunity to make new discoveries while also allowing us to focus on candidates that are more promising than randomly selected genes [16, 17]. We use functional networks from IMP to assign weights and apply weighted FDR to adjust for multiple testing. The algorithm is depicted in figure 1 with the following steps:

- Identify a small set of key genes that relate to a biological process of interest.

- Input the set of gene names into IMP (https://imp.princeton.edu) and identify networks that linked to those key genes.

- For each gene in the network, calculate the total relationship confidence $RC_i$, which is the likelihood of connecting with other genes in the network.

- Set $W_i \propto RC_i$ and apply weighted FDR to adjust for multiple testing.

*Zebrafish Experiment and RNA-seq dataset*

There is a growing interest in establishing cystic fibrosis (CF)-environment linkages. Using our zebrafish model, we

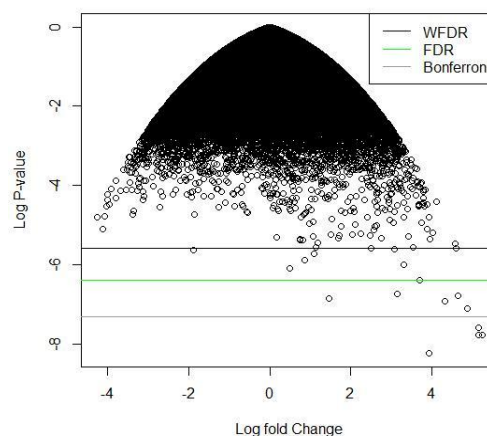**Figure 2 Zebrafish experiment**



have recently shown that the gene responsible for CF, the cystic fibrosis transmembrane conductance regulator (Cftr), and arsenic each separately mediate aspects of innate immunity [15][18-20]. We hypothesize that arsenic disrupts the innate immune response at least in part through its effects on Cftr signaling pathways. In our current studies, we are establishing linkages that show how arsenic affects the innate immune response through its effects on Cftr function. Patients with CF typically succumb to infection by the Gram-negative bacterium P. aeruginosa. We recently established a zebrafish model for P. aeruginosa infection and have applied it to our Cftr morphant fish [20]. We have shown that Cftr morphant fish are indeed more susceptible to P. aeruginosa infection. In this study, we generated twelve groups of zebrafish, each with three replicates, based upon (1) expression of Cftr, (2) exposure to different levels of arsenic, and (3) infection with P. aeruginosa (Figure 2) for each of the three time points tested (6, 12, and 18 h post-P. aeruginosa infection). At 6, 12, and 18 h post-infection, larvae were collected and homogenized in trizol reagent, using a motorized mortar and pestle. Total RNA was extracted, as per manufacturer's instructions (Life Technologies, Carlsbad, CA). RNA integrity was ensured using Nano 6000 RNA chips with the Agilent 2100 Bioanalyzer (Santa Clara, CA). RNA concentrations were determined with the Nanodrop UV/Vis spectrophotometer (Wilmington, DE). One microgram of total RNA (RIN>7) was used as input material to the Illumina TruSeq v2 high throughput library construction procedure (Illumina Inc.) Library validation was performed using the Agilent Bioanalyzer 2100 followed by quantitation using the Qubit HS DNA assay and qPCR kit for Illumina (Kapa Biosystems Inc). Libraries were diluted using the Qubit or qPCR information and sequenced on the HiSeq 2000 (Agilent). The RPKM (reads per kilobase per million) value is calculated for each gene using CLC genome bench 4.5.

RESULTS

We developed the novel IMP-weighted FDR (IMP-WFDR) procedure to assign FDR weights. Our evaluations using simulations demonstrated that the weighted FDR approach

**Figure 3 Compare power of WFDR, FDR and Bonferroni correction.**

appropriately controlled error and can provide additional power. We applied this method to an RNA-seq experiment that tests an organism's response to exposure to pathogens and arsenic, and found that IMP-WFDR effectively identified important players.

### A. Simulation Result

Figure 3 plots the test statistics versus the log-p-values of the 1,000,000 simulated test statistics. The three colored horizontal lines denote the cut-off defined by Bonferroni correction, FDR and weighted FDR. WFDR rejects more p-values than the other two procedures. It shows that when the weight is assigned appropriately, e.g. via the IMP weighting procedure, it provides additional power over unweighted FDR and Bonferroni correction.

### B. Zebrafish Experiment Result

We calculated p-values for each treatment by comparing mean RPKM (reads per kilobase per million) between PBS control samples and P. aeruginos-treated samples while varying the arsenic exposure at 0 and 10 ppb concentrations at 6 or 12 hours after treatment. To calculate weights for each gene using the IMP-WFDR procedure, we first identified 14 important genes that are associated with arsenic exposure (fn1, notch1a, notch1b, pik3r1, akt2, ass1, nfkb2, foxo5) and immune response (il1b, tnfa, il8, mmp9, irak3, ifnphi1), based on prior biological knowledge.. We queried IMP for these genes and obtained the subnetwork of all genes connected to these 14 starting genes and applied the IMP-WFDR procedure to generate weights. We compared FDR using Benjamin and Hochberg's (BH) procedure and IMP-WFDR to adjust all p-values and used 0.2 as a cut-off. There were 4 comparisons where at least one gene's adjusted p-value was under 0.2. Table 1 shows the number of findings identified by IMP-WFDR and FDR. IMP-WFDR outperforms FDR in 3 comparisons. This indicated that IMP-WFDR can greatly improve power in detecting differentially expressed genes. Indeed, IMP-WFDR analysis allowed us to identify numerous genes in each of the comparison groups that had previously been associated with arsenic exposure and/or immune challenge. Most of these same genes were not identifd by FDR.

TABLE 1. COMPARE IMP-WFDR AND FDR METHOD IN ZEBRAFISH EXPERIMENT

| Arsenic 0 vs 10 ppb for PBS samples at 6H | | Arsenic 0 vs 10 ppb for Pseudomonas at 6H | | PBS vs Pseudomonas for Arsenic=0 at 12H | | PBS vs Pseudomonas for Arsen=10 at 6H | |
|---|---|---|---|---|---|---|---|
| IMP-WFDR | FDR | IMP-WFDR | FDR | IMP-WFDR | FDR | IMP-WFDR | FDR |
| 116 | 25 | 121 | 10 | 157 | 1 | 12 | 35 |

TABLE 2. GSEA RESULTS FOR OVERLAPPING GENES.

| Biological process | Network Freq. | Genome Freq. | Adjusted p-values | Genes |
|---|---|---|---|---|
| Notch signaling pathway | 19.4% (7/36) | 0.8% (49/6131) | 8.38E-06 | notch1b, jag2 gro2, dlb , dla , dld , jag1a |
| regulation of defense response | 5.6% (2/36) | 0.0% (3/6131) | 2.92E-02 | ifnphi1, mmp9 |
| regulation of immune effector process | 5.6% (2/36) | 0.1% (4/6131) | 3.87E-02 | ifnphi1, mmp9 |
| response to arsenic containing substance | 5.6% (2/36) | 0.1% (5/6131) | 4.82E-02 | il1b , ifnphi1 |

From the genes identified by IMP-WFDR, 36 genes were identified in more than one comparison. We anticipate that these genes are likely to play an important role in regulating responses to pathogens such as P. aeruginosa in the context of environmental exposures such as arsenic. We calculated functional enrichment p-values using this set of genes [9] to identify biological processes that were enriched within this group. Table 2 shows these processes. It is interesting to note that we identified 2 out of 6 genes (ifnphi1, mmp9) that are known to regulate the immune response in zebrafish. The ifnphi1 gene encodes a type I interferon most often associated with the antiviral ressponse [21]. The mmp9 gene encodes a collagenase that degrades extracellular matrix proteins and facilitates the migration of immune cells like neutrophils [22].

The differential expression of the ifnphi1 gene recapitulates findings that that mRNA expression of the antiviral cytokines interferon (IFN) was differentially expressed in arsenic-exposed zebrafish before and after viral infection in a previous study [14]. It is also important to note that the unweighted FDR procedure fails to identify these differentially expressed genes in any of the 4 comparisons. However, IMP-WFDR identifies the pair in 3 out of the 4 comparisons. This indicates that our proposed method can identify more biologically relevant genes than an unweighted FDR method, which ignores existing biological data and knowledge.

DISCUSSION

In this paper we present a novel IMP-WFDR method that derives weights from a state-of-the-art data integration algorithm and incorporates them in WFDR to more effectively account for multiple test given the context of available biological data and knowledge. IMP has previously been used to guide morpholino assays in zebrafish [9] and our IMP-WFDR procedure greatly extends IMP's application by improving the analysis of genomic experiments through such functional networks. We demonstrate through both a simulation study and analysis of RNA-seq data that our proposed method identifies more differentially expressed genes than unweighted FDR or Bonferroni correction, while maintaining appropriate control of the false discovery rate. We used RPKM to obtain p-values as our primary goal was to

study IMP Weighted FDR, but in practice edgeR[23] or DEseq[24] can be used to calculate p-values for RNA-Seq reads.

The advantage of IMP-WFDR is that it can effectively incorporate expert biological knowledge from published experiments and existing biological data as weights and control false discovery rate to adjust for multiple testing. Roeder and Wasserman proposed a two-valued weighting scheme [25] to up weight all p-values in a pre-specified priority list and down-weight the rest. The drawback of this approach is that it demotes novel findings since all the unknown genes are down-weighted and have less of a chance to be detected than in the un-weighted case. Rather than up- weight only genes in the priority list, IMP-WFDR obtains its weight based on the functional subgenetic network defined by the priority genes. Thus, the choice of weight is much less dependent on the given priority list. From the Table 1, we see that there were over 200 genes identified by the IMP-WFDR approach while the input list contains only 14 genes. Furthermore, neither of the two key genes, ifnohi1 and mmp9 identified as differentially expressed are included in the input list and would be down-weighted by a priority list alone.

In conclusion IMP-WFDR is a powerful analytic tool that embraces state-of-the-art genetic network information in identifying differentially expressed genes in high-dimensional settings. We believe that it will play an important role as part of a research strategy to understand genetic influences on disease outcomes that embraces the complexity of the genotype-phenotype mapping relationship.

### REFERENCES

[1] J. C. Marioni, C. E. Mason, S. M. Mane, M. Stephens and Y. Gilad. RNA-seq: An assessment of technical reproducibility and comparison with gene expression arrays. Genome Res. 18(9), pp. 1509-1517. 2008. . DOI: 10.1101/gr.079558.108 [doi].

[2] C. E. Bonferroni. Il Calcolo Delle Assicurazioni Su Gruppi Di Teste 1935.

[3] Y. Benjamini and Y. Hochberg. Controlling the false discovery rate: A practical and powerful approach to multiple testing. Journal of the Royal Statistical Society.Series B (Methodological) pp. 289-300. 1995.

[4] J. D. Storey. The positive false discovery rate: A bayesian interpretation and the q-value. Annals of Statistics pp. 2013-2035. 2003.

[5] B. Efron. Size, power and false discovery rates. The Annals of Statistics pp. 1351-1377. 2007.

[6] B. Efron. Large-scale simultaneous hypothesis testing. Journal of the American Statistical Association 99(465), 2004.

[7] B. Efron, R. Tibshirani, J. D. Storey and V. Tusher. Empirical bayes analysis of a microarray experiment. Journal of the American Statistical Association 96(456), pp. 1151-1160. 2001.

[8] C. R. Genovese, K. Roeder and L. Wasserman. False discovery control with p-value weighting. Biometrika 93(3), pp. 509-524. 2006.

[9] A. K. Wong, C. Y. Park, C. S. Greene, L. A. Bongo, Y. Guan and O. G. Troyanskaya. IMP: A multi-species functional genomics portal for integration, visualization and prediction of protein functions and networks. Nucleic Acids Res. 40(Web Server issue), pp. W484-90. 2012. . DOI: 10.1093/nar/gks458 [doi].

[10] C. Y. Park, A. K. Wong, C. S. Greene, J. Rowland, Y. Guan, L. A. Bongo, R. D. Burdine and O. G. Troyanskaya. Functional knowledge transfer for high-accuracy prediction of under-studied biological processes. PLoS Computational Biology 9(3), pp. e1002957. 2013.

[11] C. Huttenhower, E. M. Haley, M. A. Hibbs, V. Dumeaux, D. R. Barrett, H. A. Coller and O. G. Troyanskaya. Exploring the human genome with functional maps. Genome Res. 19(6), pp. 1093-1106. 2009. . DOI: 10.1101/gr.082214.108 [doi].

[12] R. L. Young, K. C. Malcolm, J. E. Kret, S. M. Caceres, K. R. Poch, D. P. Nichols, J. L. Taylor-Cousar, M. T. Saavedra, S. H. Randell and M. L. Vasil. Neutrophil extracellular trap (NET)-mediated killing of pseudomonas aeruginosa: Evidence of acquired resistance within the CF airway, independent of CFTR. PLoS One 6(9), pp. e23637. 2011.

[13] J. M. Bomberger, B. A. Coutermarsh, R. L. Barnaby and B. A. Stanton. Arsenic promotes ubiquitinylation and lysosomal degradation of cystic fibrosis transmembrane conductance regulator (CFTR) chloride channels in human airway epithelial cells. J. Biol. Chem. 287(21), pp. 17130-17139. 2012. . DOI: 10.1074/jbc.M111.338855 [doi].

[14] J. R. Shaw, J. M. Bomberger, J. VanderHeide, T. LaCasse, S. Stanton, B. Coutermarsh, R. Barnaby and B. A. Stanton. Arsenic inhibits SGK1 activation of CFTR cl$^-$ channels in the gill of killifish,$ $ fundulus heteroclitus. Aquatic Toxicology 98(2), pp. 157-164. 2010.

[15] A. S. Nayak, C. R. Lage and C. H. Kim. Effects of low concentrations of arsenic on the innate immune system of the zebrafish (danio rerio). Toxicol. Sci. 98(1), pp. 118-124. 2007. . DOI: kfm072 [pii].

[16] M. A. Hibbs, C. L. Myers, C. Huttenhower, D. C. Hess, K. Li, A. A. Caudy and O. G. Troyanskaya. Directing experimental biology: A case study in mitochondrial biogenesis. PLoS Computational Biology 5(3), pp. e1000322. 2009.

[17] Y. Guan, C. L. Myers, R. Lu, I. R. Lemischka, C. J. Bult and O. G. Troyanskaya. A genomewide functional network for the laboratory mouse. PLoS Computational Biology 4(9), pp. e1000165. 2008.

[18] C. R. Lage, A. Nayak and C. H. Kim. Arsenic ecotoxicology and innate immunity. Integr. Comp. Biol. 46(6), pp. 1040-1054. 2006. . DOI: 10.1093/icb/icl048 [doi].

[19] A. C. Hermann and C. H. Kim. Effects of arsenic on zebrafish innate immune system. Marine Biotechnology 7(5), pp. 494-505. 2005.

[20] R. T. Phennicie, M. J. Sullivan, J. T. Singer, J. A. Yoder and C. H. Kim. Specific resistance to pseudomonas aeruginosa infection in zebrafish is mediated by the cystic fibrosis transmembrane conductance regulator. Infect. Immun. 78(11), pp. 4542-4550. 2010. . DOI: 10.1128/IAI.00302-10 [doi].

[21] O. J. Hamming, G. Lutfalla, J. P. Levraud and R. Hartmann. Crystal structure of zebrafish interferons I and II reveals conservation of type I interferon structure in vertebrates. J. Virol. 85(16), pp. 8181-8187. 2011. . DOI: 10.1128/JVI.00521-11 [doi].

[22] L. M. Bradley, M. F. Douglass, D. Chatterjee, S. Akira and B. J. Baaten. Matrix metalloprotease 9 mediates neutrophil migration into the airways in response to influenza virus-induced toll-like receptor signaling. PLoS Pathogens 8(4), pp. e1002641. 2012.

[23] S. Anders, D.J. McCarthy, Y. Chen, M. Okoniewski, G.K. Smyth, W. Huber, and M.D. Robinson. Count-based differential expression analysis of RNA sequencing data using R and Bioconductor. Nature Protocols 8, 1765–1786. 2013.

[24] S. Anders and W. Huber. Differential expression analysis for sequence count data. Genome Biology, 11, pp. R106. 2010.

[25] K. Roeder and L. Wasserman. Genome-wide significance levels and weighted hypothesis testing. Stat. Sci. 24(4), pp. 398-413. 2009. . DOI: 10.1214/09-STS289 [doi].