



miRNA Temporal Analyzer (mirnaTA): a bioinformatics tool for identifying differentially expressed microRNAs in temporal studies using normal quantile transformation

Cer *et al.*

TECHNICAL NOTE

Open Access

miRNA Temporal Analyzer (mirnaTA): a bioinformatics tool for identifying differentially expressed microRNAs in temporal studies using normal quantile transformation

Regina Z Cer^{1,2*}, J Enrique Herrera-Galeano^{1,2}, Joseph J Anderson³, Kimberly A Bishop-Lilly^{1,2} and Vishwesh P Mokashi¹

Abstract

Background: Understanding the biological roles of microRNAs (miRNAs) is an active area of research that has produced a surge of publications in PubMed, particularly in cancer research. Along with this increasing interest, many open-source bioinformatics tools to identify existing and/or discover novel miRNAs in next-generation sequencing (NGS) reads become available. While miRNA identification and discovery tools are significantly improved, the development of miRNA differential expression analysis tools, especially in temporal studies, remains substantially challenging. Further, the installation of currently available software is non-trivial and steps of testing with example datasets, trying with one's own dataset, and interpreting the results require notable expertise and time. Subsequently, there is a strong need for a tool that allows scientists to normalize raw data, perform statistical analyses, and provide intuitive results without having to invest significant efforts.

Findings: We have developed miRNA Temporal Analyzer (mirnaTA), a bioinformatics package to identify differentially expressed miRNAs in temporal studies. mirnaTA is written in Perl and R (Version 2.13.0 or later) and can be run across multiple platforms, such as Linux, Mac and Windows. In the current version, mirnaTA requires users to provide a simple, tab-delimited, matrix file containing miRNA name and count data from a minimum of two to a maximum of 20 time points and three replicates. To recalibrate data and remove technical variability, raw data is normalized using Normal Quantile Transformation (NQT), and linear regression model is used to locate any miRNAs which are differentially expressed in a linear pattern. Subsequently, remaining miRNAs which do not fit a linear model are further analyzed in two different non-linear methods 1) cumulative distribution function (CDF) or 2) analysis of variances (ANOVA). After both linear and non-linear analyses are completed, statistically significant miRNAs ($P < 0.05$) are plotted as heat maps using hierarchical cluster analysis and Euclidean distance matrix computation methods.

Conclusions: mirnaTA is an open-source, bioinformatics tool to aid scientists in identifying differentially expressed miRNAs which could be further mined for biological significance. It is expected to provide researchers with a means of interpreting raw data to statistical summaries in a fast and intuitive manner.

Keywords: microRNA, miRNA Temporal Analyzer, mirnaTA, Time series, Differential expression, DE, Quantile normalization, Linear model, Normal quantile transformation

* Correspondence: regina.cer@med.navy.mil

¹Biological Defense Research Directorate, Naval Medical Research Center-Frederick, 8400 Research Plaza, Fort Detrick, MD 21702, USA

²Henry M. Jackson Foundation for the Advancement of Military Medicine, 6720-A Rockledge Drive, Suite 100, Bethesda, MD 20817, USA

Full list of author information is available at the end of the article

Findings

MicroRNAs (miRNAs) are short single-stranded non-coding RNAs approximately 19–22 nucleotide long, which are critical regulators of gene expression and have been implicated in a wide range of physiological processes, such as apoptosis and growth, as well as pathological processes, including inflammatory responses, cancer, neurodegenerative and cardiovascular diseases [1-7]. This rapid growth is evident by the exponentially increasing number of miRNAs reported in the recent Release 21 (June 2014) of miR-Base [8,9] which contains 35,828 mature miRNA products in 223 different species. Accompanying this growth is the development of many miRNA discovery bioinformatics tools including, but not limited to miRscan [10], miR-Finder [11], miRDeep [12] and miRanalyzer [13], to help researchers identify miRNAs from existing miRNA databases and/or predict novel miRNAs from NGS data.

In recent years, the number of miRNA experiments involving multiple time-series has largely increased [14]. This is not surprising since biological processes are often dynamic, and therefore, a cross-sectional study based on a single time point would not provide important information that time-course studies can provide. Particularly in miRNA studies, researchers often implement multiple-time points to capture how certain miRNAs may display transient expression changes over the course of treatment or infection over time. For instance, Jayaswal *et al.* carried out a drug study involving a multiple myeloma cell line, U266, and consisting of six time points—0, 2, 4, 8, 24, and 48 hours with two biological replicates per time point for both miRNA and mRNA [15]. In another study by Li Z *et al.*, miRNA expression profiles were assessed in the mouse livers in a time-course experiment at days 1, 3, 7, 15, 30 and 120 post treatment [16].

Despite the fact that there are several publications dealing with the statistical analysis of time-series expression data, there are significant computational challenges in making sense of the massive and complex datasets produced by time-series experiments. There are a number of differential expression (DE) tool packages including STEM [17], MaTSE [18], Linear Models for Microarray Data (LIMMA) [19], Significance Analysis of Microarrays (SAM) [20], Extraction of Differential Gene Expression (EDGE) [21], and Bayesian Estimation of Temporal Regulation (BETR) [22]. In miRNA research in particular, the two most commonly used DE tools are edgeR [23] and DEseq [24]. However, learning to use these tools requires the expertise of bioinformaticians as they often come with several package dependencies and not every research laboratory employs these specialists nor has advanced computing infrastructure. There is a crucial need for a simple tool that would allow any bench scientist to analyze the differential expression of miRNAs or other subjects in temporal studies.

Development of mirnaTA

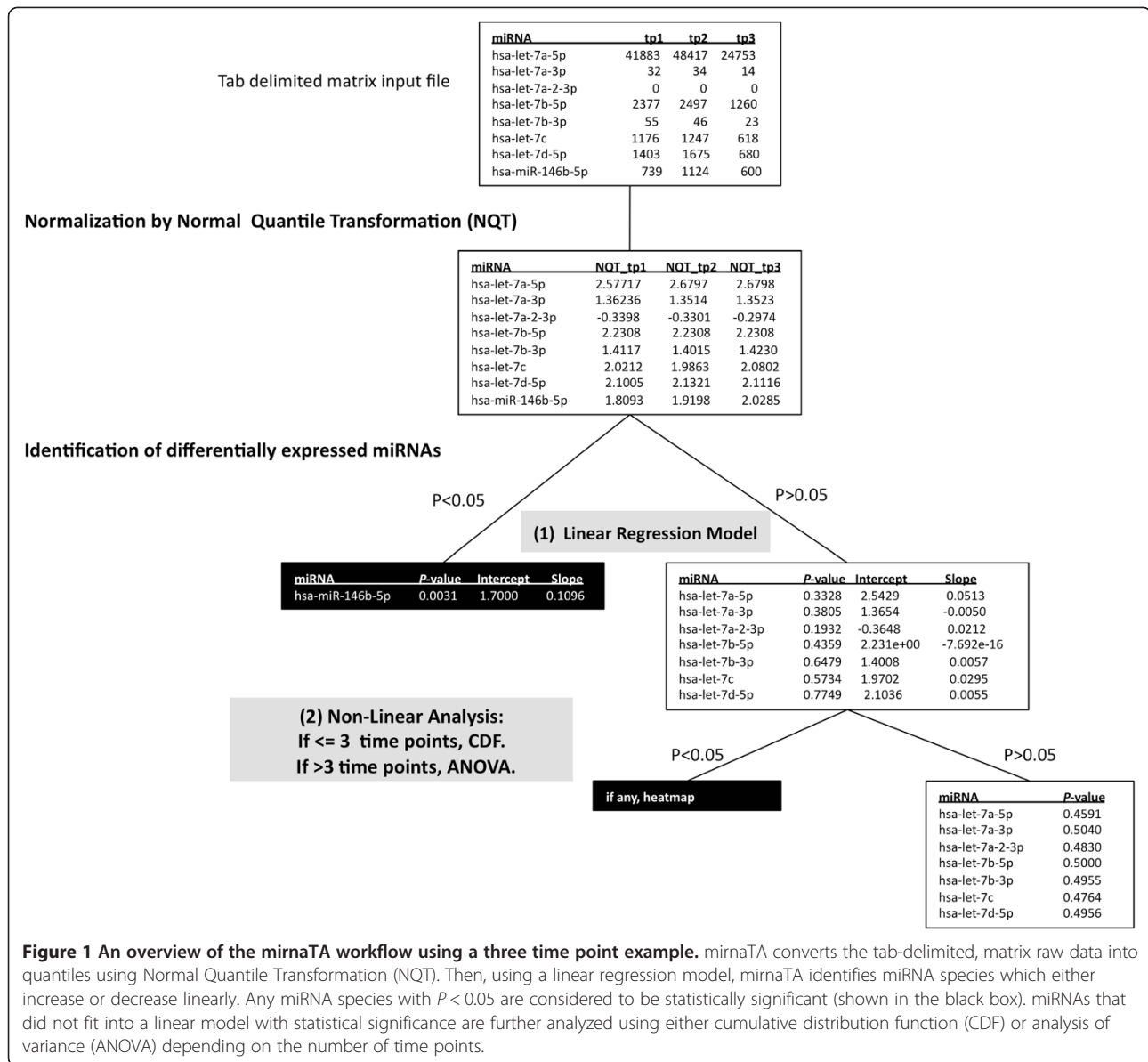
The miRNA Temporal Analyzer (mirnaTA) is a by-product of a recent study where human peripheral blood mononuclear cells (hPBMCs) were exposed to bacterial pathogens and assayed at three different time points. Briefly, total RNA was extracted, miRNA sequencing libraries were constructed, sequencing was performed on the MiSeq sequencer, and miRNA sequences were identified from sequencing reads using a standard procedure (see in Additional file 1: Figure S1). Then, mirnaTA was used to identify differentially expressed miRNAs for further validation in the laboratory [Cer *et al.*, unpublished data]. mirnaTA performs a number of distinct steps (Figure 1): (i) normalization of the raw count data into quantiles using Normal Quantile Transformation (NQT), (ii) analysis of NQT data to locate any miRNA species which are either increasing or decreasing linearly using linear regression model, (iii) further analysis of miRNA species that did not fit in linear model by two different methods: (a) normal distribution function known as cumulative distribution function (CDF) if the number of time points is ≤ 3 or (b) analysis of variances (ANOVA) if the number of time points is >3 , (iv) generation of heat maps for any miRNAs that are differentially expressed with statistical significance ($P < 0.05$), and (v) providing intuitive HTML output formats. Some of these steps are explored in more detail below.

Input file requirement

Any datasets resulting from RNA-seq experiments, array-based miRNA profiling or RT-PCR experiments, can be used as input data for the mirnaTA as long as the tab-delimited input files contain two columns, miRNA name and sequence read count of the miRNA. The number of time points that can be submitted to mirnaTA is a minimum of 2 to a maximum of 20. mirnaTA is able to handle to three replicates as follows: if the correlation coefficient (r) value between two replicates is more than or equal to 0.70, the count values from replicates are averaged and used in the study. Incremental correlation is applied for each additional replicate (see User Guide available at [25] for more details). If the r value is less than 0.7, the replicate data is not used in the calculations. The correlation coefficient, r value 0.7 was chosen arbitrarily as it was deemed reasonable for replicates to be considered sufficiently correlated. However, users are able to change this to any value although we do not recommend to use any value below 0.7.

Data normalization by NQT

Data normalization, a critical step of DE analysis, attempts to remove technical variability while preserving biological significance. Although there is no consensus on the best method for the normalization of microRNA sequencing



data, when seven commonly used normalization methods, namely, global normalization, Lowess normalization, trimmed mean method (TMM), quantile normalization, scaling normalization, variance stabilization (VSN) and invariant method (IN), were evaluated, the Lowess normalization and the quantile normalization were found to be the best approaches [26]. It should be noted that edgeR uses TMM and DESeq uses a negative binomial (NB) model. mirnaTA is different in that it implements quantile normalization (also referred here as (NQT)) chosen for its robustness and reduced bias. Quantile normalization is a rank-based procedure which scales data within each quantile separately and has been shown to have bias near zero relative to qRT-PCR [27]. Another factor considered in choosing quantile-based scaling is its

medium computational complexity [28]. Briefly, raw data is normalized using NQT as follows: random numbers are drawn from the normal distribution based on the number of observations (here the number of miRNA species) and these numbers are sorted (Equation 1). Then quantiles are assigned according to each rank (Equations 2–3). All the equations included here are R language commands unless or otherwise specified:

$$rn = \text{sort}(rnorm(\text{no_of_miRNAs})) \quad (1)$$

$$qt = \text{rank}(x) / \text{length}(x) \quad (2)$$

$$nqt = rn[\text{rank}(qt)] \quad (3)$$

Identification of differentially expressed miRNAs by a linear regression model

Using a linear regression, quantile transformed data is analyzed to identify any miRNA species which display either increasing or decreasing expression. Using NQT data over different time points, the corresponding *P*-values, intercept, and slope values are calculated (Equations 4–8), and any miRNAs with *P* < 0.05 are considered to be statistically significant:

$$lmod = lm(y \sim x) \text{ where } x = \text{time}, y = \text{NQT data} \quad (4)$$

$$s = \text{summary}(lmod) \quad (5)$$

$$p = \text{as.character}(pf(s\$fstatistic[1], s\$fstatistic[2], s\$fstatistic[3], \text{lower.tail} = \text{FALSE})) \quad (6)$$

$$\text{intercept} = s\$coefficients[1] \quad (7)$$

$$\text{slope} = s\$coefficients[2] \quad (8)$$

Identification of other differentially expressed miRNA species using a non-linear analysis

The remaining miRNA species that did not fit in a linear model are further analyzed depending on the number of time points in the experiment: (a) the cumulative distribution function (CDF) is applied when the number of time points is ≤ 3 (Equation 9), and (b) without the assumption of equal variances for all groups, the analysis of variance (ANOVA) [29] is applied when the number of time points is > 3 . Similar to linear analysis criteria, any miRNAs with *P* < 0.05 are considered to be statistically significant.

$$\begin{aligned} cdf_{pval} &= 1 - pnorm((as.numeric(tpn) \\ &\quad - as.numeric(tp1)), \text{mean} = 0, \text{sd} \\ &= 1, \text{lower.tail} = \text{TRUE}, \text{log.p} \\ &= \text{FALSE})) \text{ where } tp1 \\ &= \text{nqt data for time point 1 and } tpn \\ &= \text{nqt data for time point 2 or time point 3} \end{aligned} \quad (9)$$

$$\text{anova}_{pval} = \text{oneway.test}(y \sim x) \text{ where } y = \text{nqt data and } x = \text{time points divided into two groups} \quad (10)$$

Syntax for running mirnaTA

The syntax to run mirnaTA is simple, as follows: *perl nmrc_mirnata.pl -i <input_file> -n <x> -t <y>*, where *nmrc_mirnata.pl* is the wrapper Perl script, *<input_file>* is a tab-delimited, matrix file, *x* = the number of replicates,

y = the number of time points. For example, to run mirnaTA for a dataset with three time points, two replicates and correlation coefficient value of 0.85 (default is 0.7), the syntax would be “*perl nmrc_mirnata.pl -i example_datasets/3tp2rep_dataset -n 2 -t 3 -c 0.85*”. To demonstrate the functionality of mirnaTA, a comprehensive User Guide with examples using different datasets and expected results are provided at our SourceForge [25] page and also archived in GigaDB [30].

Visualizations of mirnaTA outputs

mirnaTA provides a rich source of HTML format data, including raw data, quantile transformed data, a heat map of differentially expressed miRNAs (*P* < 0.05) which were identified to be linearly increasing or decreasing using a linear regression model, as well as a heat map of differentially expressed miRNAs (*P* < 0.05) which were identified to be increasing or decreasing using either a cumulative distribution function (CDF) or an analysis of variance (ANOVA), if present. Also included are several intermediate files including a list of significant miRNA species with their *P*-values, slope, and intercept values (Figure 2). All these information is integrated into an HTML interface which can be viewed in any web browser. It should be noted that the heat map text files may also be visualized in other software packages such as PermutMatrix [31].

Conclusions

mirnaTA is an open-source bioinformatics tool that can be run in Linux, Mac or Windows with Perl and R package dependencies. While there are many other time-series analysis tools available, the main advantages of mirnaTA are:

- (1). **Simplicity:** Even though it is a command line program and users need to be aware of corrections for multiple testing, they need to provide only one input file and enter one line of Perl syntax, and do not need to have a strong statistical or computational background.
- (2). **Reliability:** The statistically significant miRNA species identified by linear regression model have met very stringent criteria, and therefore could be further examined for confirmation or validation in wet laboratories.
- (3). **Practicality:** The mirnaTA only requires a simple input file which can be prepared by anyone, takes up to three replicates and 20 time points, and generates publication quality graphics in PNG format.

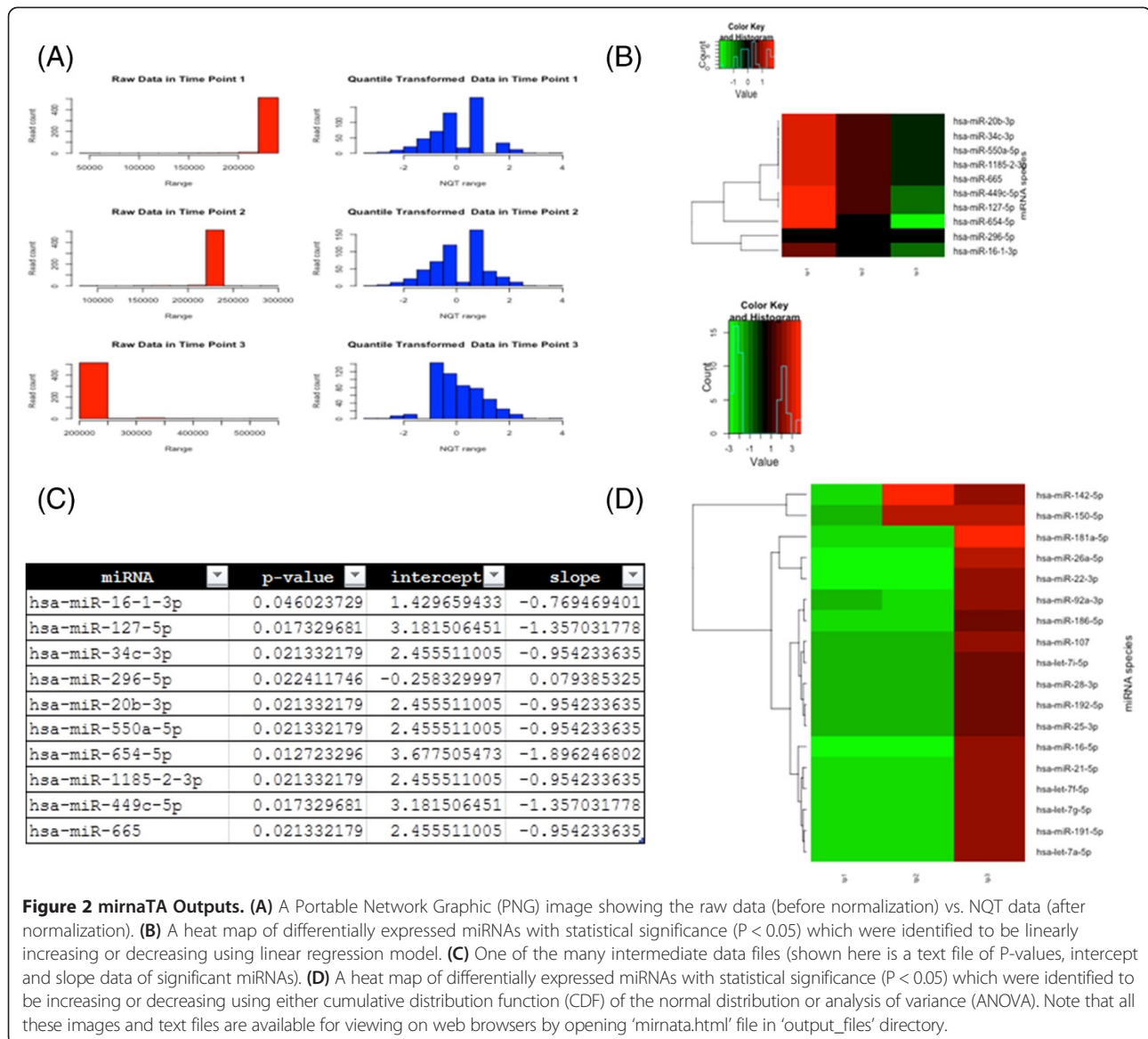


Figure 2 mirnaTA Outputs. (A) A Portable Network Graphic (PNG) image showing the raw data (before normalization) vs. NQT data (after normalization). (B) A heat map of differentially expressed miRNAs with statistical significance ($P < 0.05$) which were identified to be linearly increasing or decreasing using linear regression model. (C) One of the many intermediate data files (shown here is a text file of P-values, intercept and slope data of significant miRNAs). (D) A heat map of differentially expressed miRNAs with statistical significance ($P < 0.05$) which were identified to be increasing or decreasing using either cumulative distribution function (CDF) of the normal distribution or analysis of variance (ANOVA). Note that all these images and text files are available for viewing on web browsers by opening 'mirnata.html' file in 'output_files' directory.

(4). Wide applicability: The mirnaTA can be used to analyze any kind of data (gene, weather data, weight, etc.), and does not have to be restricted to miRNA data. All a user has to do is to replace the miRNA name in the first column with another subject of study.

In summary, although mirnaTA does not yet provide researchers with “a push-button” solution for differential expression analysis, it gives them the power to use a simple Perl syntax to not only identify differentially expressed miRNAs, but also automatically generate intuitive graphical outputs. This simple and automatic implementation of previously detached concepts in mirnaTA could be instrumental in saving a significant amount of time for many basic research scientists.

Availability and requirements

Project name: miRNA Temporal Analyzer (mirnaTA)

Project home page: <http://sourceforge.net/projects/mirnata>

Operating system (s): Linux, Mac and Windows

Programming language: Perl and R

Other Requirements: None

License: The GNU Lesser General Public License, version 3.0 (LGPL-3.0)

Any restrictions to use by non-academics: None

Availability of supporting data

Archival version of the supporting files, user guide and additional datasets used in the paper are hosted in the *GigaScience* GigaDB database [30], and for the most up to date versions please see the source forge page: <http://sourceforge.net/projects/mirnata>.

Additional file

Additional file 1: Figure S1. Detailed steps for generating input files for mirnaTA. FASTQ files generated from any NGS sequencing platform are converted into FASTA files. Artificially introduced 3' adapter sequences are trimmed, and post-trimmed reads that are a minimum of 15 base pairs are filtered against contaminants. Reads that do not match to contaminants are screened for mature miRNA species (black box) which are further analyzed for statistical significance using mirnaTA.

Abbreviations

ANOVA: Analysis of variance; CDF: Cumulative distribution function; DE: Differential expression; miRNA: microRNA; NGS: Next-generation sequencing; NQT: Normal quantile transformation; PNG: Portable network graphics; TMM: Trimmed mean method.

Competing interests

The authors declare that they have no competing interests.

Authors' contributions

RZC wrote Perl and R scripts, packaged the workflow, released code and prepared the manuscript. JEH wrote custom R functions and oversaw R statistical analyses. JJA tested the package and provided patches. KAB tested the package and edited the manuscript. VPM oversaw the project and gave scientific advice. All authors read, contributed and approved the final manuscript.

Acknowledgements

mirnaTA was developed as a part of a study supported by the Defense Threat Reduction Agency (DTRA) project CBM.DIAGB.03.10.NM.028. JJA and VPM are military service members or employees of the U.S. Government and this work was prepared as part of their official duties. Title 17 U.S.C. §105 provides that 'Copyright protection under this title is not available for any work of the United States Government'. Title 17 U.S.C. §101 defines a U.S. Government work as a work prepared by a military service member or employee of the U.S. Government as part of that person's official duties. The opinions or assertions contained herein are the private ones of the author (s) and are not to be construed as official or reflecting the views of either the Department of the Navy or the Department of Defense, nor the U.S. Government.

Author details

¹Biological Defense Research Directorate, Naval Medical Research Center-Frederick, 8400 Research Plaza, Fort Detrick, MD 21702, USA. ²Henry M. Jackson Foundation for the Advancement of Military Medicine, 6720-A Rockledge Drive, Suite 100, Bethesda, MD 20817, USA. ³Chem Bio Research Center of Excellence, Defense Threat Reduction Agency, 2800 Bush River Road E2800-b198, Aberdeen Proving Ground, MD 21010, USA.

Received: 8 April 2014 Accepted: 30 September 2014

Published: 13 October 2014

References

- Kim VN, Han J, Siomi MC: Biogenesis of small RNAs in animals. *Nat Rev Mol Cell Biol* 2009, **10**(2):126–139.
- Xiao C, Rajewsky K: MicroRNA control in the immune system: basic principles. *Cell* 2009, **136**(1):26–36.
- Tili E, Michaille JJ, Calin GA: Expression and function of micro-RNAs in immune cells during normal or disease state. *Int J Med Sci* 2008, **5**(2):73–79.
- Matsushima K, Isomoto H, Inoue N, Nakayama T, Hayashi T, Nakayama M, Nakao K, Hirayama T, Kohno S: MicroRNA signatures in *Helicobacter pylori*-infected gastric mucosa. *Int J Cancer* 2011, **128**(2):361–370.
- Wang WX, Wilfred BR, Madathil SK, Tang G, Hu Y, Dimayuga J, Stromberg AJ, Huang Q, Saatman KE, Nelson PT: miR-107 regulates granulin/progranulin with implications for traumatic brain injury and neurodegenerative disease. *Am J Pathol* 2010, **177**(1):334–345.
- Garza-Manero S, Pichardo-Casas I, Arias C, Vaca L, Zepeda A: Selective distribution and dynamic modulation of miRNAs in the synapse and its possible role in Alzheimer's disease. *Brain Res* 2013, **1584**:80–93.
- Ono K, Kuwabara Y, Han J: MicroRNAs and cardiovascular diseases. *FEBS J* 2011, **278**(10):1619–1633.
- Ambros V, Bartel B, Bartel DP, Burge CB, Carrington JC, Chen X, Dreyfuss G, Eddy SR, Griffiths-Jones S, Marshall M, Matzke M, Ruvkun G, Tuschl T: A uniform system for microRNA annotation. *RNA* 2003, **9**(3):277–279.
- Kozomara A, Griffiths-Jones S: miRBase: integrating microRNA annotation and deep-sequencing data. *Nucleic Acids Res* 2011, **39**(Database issue):D152–D157.
- Lim LP, Lau NC, Weinstein EG, Abdelhakim A, Yekta S, Rhoades MW, Burge CB, Bartel DP: The microRNAs of *Caenorhabditis elegans*. *Genes Dev* 2003, **17**(8):991–1008.
- Huang TH, Fan B, Rothschild MF, Hu ZL, Li K, Zhao SH: MiRFinder: an improved approach and software implementation for genome-wide fast microRNA precursor scans. *BMC Bioinformatics* 2007, **8**:341.
- An J, Lai J, Lehman ML, Nelson CC: miRDeep*: an integrated application tool for miRNA identification from RNA sequencing data. *Nucleic Acids Res* 2013, **41**(2):727–737.
- Hackenbergh M, Rodriguez-Ezpeleta N, Aransay AM: miRanalyzer: an update on the detection and analysis of microRNAs in high-throughput sequencing experiments. *Nucleic Acids Res* 2011, **39**(Web Server issue):W132–W138.
- Barrett T, Wilhite SE, Ledoux P, Evangelista C, Kim IF, Tomashevsky M, Marshall KA, Phillippy KH, Sherman PM, Holko M, Yefanov A, Lee H, Zhang N, Robertson CL, Serova N, Davis S, Soboleva A: NCBI GEO: archive for functional genomics data sets—update. *Nucleic Acids Res* 2013, **41**(Database issue):D991–D995.
- Jayaswal V, Lutherborrow M, Ma DD, Hwa Yang Y: Identification of microRNAs with regulatory potential using a matched microRNA-mRNA time-course data. *Nucleic Acids Res* 2009, **37**(8):e60.
- Li Z, Branham WS, Dial SL, Wang Y, Guo L, Shi L, Chen T: Genomic analysis of microRNA time-course expression in liver of mice treated with genotoxic carcinogen N-ethyl-N-nitrosourea. *BMC Genomics* 2010, **11**:609.
- Ernst J, Bar-Joseph Z: STEM: a tool for the analysis of short time series gene expression data. *BMC Bioinformatics* 2006, **7**:191.
- Craig P, Cannon A, Kukla R, Kennedy J: MaTSE: the gene expression time-series explorer. *BMC Bioinformatics* 2013, **14**:S1.
- Smyth GK: Linear models and empirical bayes methods for assessing differential expression in microarray experiments. *Stat Appl Genet Mol Biol* 2004, **3**:Article 3. Epub 2004 Feb 12.
- Grace C, Nacheva EP: Significance Analysis of Microarrays (SAM) offers clues to differences between the genomes of adult philadelphia positive ALL and the lymphoid blast transformation of CML. *Cancer Inform* 2012, **11**:173–183.
- Leek JT, Monsen E, Dabney AR, Storey JD: EDGE: extraction and analysis of differential gene expression. *Bioinformatics* 2006, **22**(4):507–508.
- Aryee MJ, Gutierrez-Pabello JA, Kramnik I, Maiti T, Quackenbush J: An improved empirical bayes approach to estimating differential gene expression in microarray time-course data: BETR (Bayesian Estimation of Temporal Regulation). *BMC Bioinformatics* 2009, **10**:409.
- Robinson MD, McCarthy DJ, Smyth GK: edgeR: a Bioconductor package for differential expression analysis of digital gene expression data. *Bioinformatics* 2010, **26**(1):139–140.
- Anders S, Huber W: Differential expression analysis for sequence count data. *Genome Biol* 2010, **11**(10):R106.
- SourceForge: <http://sourceforge.net/projects/mirna>.
- Garmire LX, Subramaniam S: Evaluation of normalization methods in mammalian microRNA-Seq data. *RNA* 2012, **18**(6):1279–1288.
- Bullard JH, Purdom E, Hansen KD, Dudoit S: Evaluation of statistical methods for normalization and differential expression in mRNA-Seq experiments. *BMC Bioinformatics* 2010, **11**:94.
- McCormick KP, Willmann MR, Meyers BC: Experimental design, preprocessing, normalization and differential expression analysis of small RNA sequencing experiments. *Silence* 2011, **2**(1):2.
- Dalgaard P: *Statistics and Computing: Introductory Statistics with R*. New York: Springer Science; 2002.
- Cer RZ, Herrera-Galeano JE, Anderson JJ, Bishop-Lilly KA, Mokashi VP: Software and Supporting Material for: "mirnaTA: a Bioinformatics Tool for

Identifying Differentially Expressed microRNAs in Temporal Studies using Normal Quantile Transformation (NQT). GigaScience Database; 2014.
<http://dx.doi.org/10.5524/100107>.

31. Caraux G, Pinloche S: **PermutMatrix: a graphical environment to arrange gene expression profiles in optimal linear order.** *Bioinformatics* 2005, **21**(7):1280–1281.

doi:10.1186/2047-217X-3-20

Cite this article as: Cer et al.: miRNA Temporal Analyzer (mirnaTA): a bioinformatics tool for identifying differentially expressed microRNAs in temporal studies using normal quantile transformation. *GigaScience* 2014 **3**:20.

**Submit your next manuscript to BioMed Central
and take full advantage of:**

- Convenient online submission
- Thorough peer review
- No space constraints or color figure charges
- Immediate publication on acceptance
- Inclusion in PubMed, CAS, Scopus and Google Scholar
- Research which is freely available for redistribution

Submit your manuscript at
www.biomedcentral.com/submit

