


SCIENTIFIC REPORTS



OPEN

GREIN: An Interactive Web Platform for Re-analyzing GEO RNA-seq Data

Naim Al Mahi¹, Mehdi Fazel Najafabadi¹, Marcin Pilarczyk¹, Michal Kouril² & Mario Medvedovic¹ 

The vast amount of RNA-seq data deposited in Gene Expression Omnibus (GEO) and Sequence Read Archive (SRA) is still a grossly underutilized resource for biomedical research. To remove technical roadblocks for reusing these data, we have developed a web-application GREIN (GEO RNA-seq Experiments Interactive Navigator) which provides user-friendly interfaces to manipulate and analyze GEO RNA-seq data. GREIN is powered by the back-end computational pipeline for uniform processing of RNA-seq data and the large number (>6,500) of already processed datasets. The front-end user interfaces provide a wealth of user-analytics options including sub-setting and downloading processed data, interactive visualization, statistical power analyses, construction of differential gene expression signatures and their comprehensive functional characterization, and connectivity analysis with LINCS L1000 data. The combination of the massive amount of back-end data and front-end analytics options driven by user-friendly interfaces makes GREIN a unique open-source resource for re-using GEO RNA-seq data. GREIN is accessible at: <https://shiny.ilincs.org/grein>, the source code at: <https://github.com/uc-bd2k/grein>, and the Docker container at: <https://hub.docker.com/r/ucbd2k/grein>.

Depositing RNA-seq datasets in Gene Expression Omnibus (GEO)¹ and Sequence Read Archive (SRA)² repositories ensures the reproducibility of published studies and facilitates its reuse. Re-analysis of these data can lead to novel scientific insights³ and it has been routinely used to inform the design of new studies⁴. However, reuse of GEO RNA-seq data is made difficult by the complexity of the processing protocols⁵ and analytical tools⁶ which are often inaccessible to biomedical scientists not specializing in bioinformatics.

Recent efforts at re-processing GEO/SRA RNA-seq data alleviate this problem by providing access to large number of processed and per-transcript summarized RNA-seq datasets which significantly simplifies its use⁷⁻¹⁰. Other resources provide access and analysis tools for specific datasets^{11,12}. While these platforms are extremely useful, they do not support additional functionalities for downstream analyses. For example, exploratory data analysis, differential expression analysis with batch effect adjustment, or statistical power analysis (Table 1). Therefore, open-source user-friendly tools with comprehensive analytical toolbox for re-analysis of public RNA-seq data are still lacking. We address this problem by developing and deploying GEO RNA-seq Experiments Interactive Navigator (GREIN) web tool for analysis of GEO RNA-seq data. In addition to the rich repertoire of analysis tools, GREIN provides access to more than 6,500 uniformly processed human, mouse, and rat GEO RNA-seq datasets with >400,000 samples that are ready for analysis. These datasets were retrieved from GEO and uniformly reprocessed by the back-end GEO RNA-seq experiments processing pipeline (GREP2). The pipeline also curates metadata for each of the datasets and annotates each sample with biomedical ontologies provided by MetaSRA¹³. As the number of new studies are included in GEO, more datasets are processed and added to GREIN on a regular basis. Apart from the preprocessed datasets, GREIN also facilitates user requested processing of GEO RNA-seq datasets on-the-fly (Table 1). We also release GREP2 as an R¹⁴ package and GREIN as a Docker¹⁵ container for easy local deployment of the complete infrastructure which can be used to reproduce GREIN results off-line.

¹Division of Biostatistics and Bioinformatics, Department of Environmental Health, University of Cincinnati, 3223 Eden Avenue, Cincinnati, OH, 45220, USA. ²Division of Biomedical Informatics, Cincinnati Children's Hospital Medical Center, Cincinnati, OH, USA. Correspondence and requests for materials should be addressed to M.M. (email: medvedm@ucmail.uc.edu)

Comparison	RNA-seq processed data resources					
	GREIN	Recount2 ⁷	ARCHS4 ⁸	Toil ⁹	Skymap ¹⁰	Expression Atlas ¹¹
Data availability	Gene and transcript read counts	Gene and transcript read counts	Gene and transcript read counts	Transcript read counts	Allelic and transcript read counts	Gene read counts
Quantification type	Read mapping (Salmon)	Read alignment (Rail-RNA)	Read mapping (Kallisto)	Read mapping (Kallisto)	Read alignment (Bowtie2)	Read alignment (Tophat)
Number of samples	>400,000*	>80,000	>300,000	>20,000	>400,000	>120,000
Interactive exploratory analysis	Yes	No	No	No	No	No
QC report	Yes	No	No	Yes	No	Yes
Power analysis	Yes	No	No	No	No	No
Differential expression analysis on-the-fly	Yes	No	No	No	No	No
Enrichment analysis for differential expression signature	Yes	No	No	No	No	Yes
Pipeline availability	R package and Docker container	R package and Docker container	Docker container	Docker container	Python scripts	NA
Process data on-the-fly	Yes	No	No	No	No	No
Ontological annotations of samples	Yes	No	No	No	No	No

Table 1. Comparison of different RNA-seq data resources. All these resources provide access to processed RNA-seq data, however most of them do not provide interactive interfaces for further manipulation and analyses of the datasets; *>250,000 are already processed and the rest are accessible using user-initiated processing.

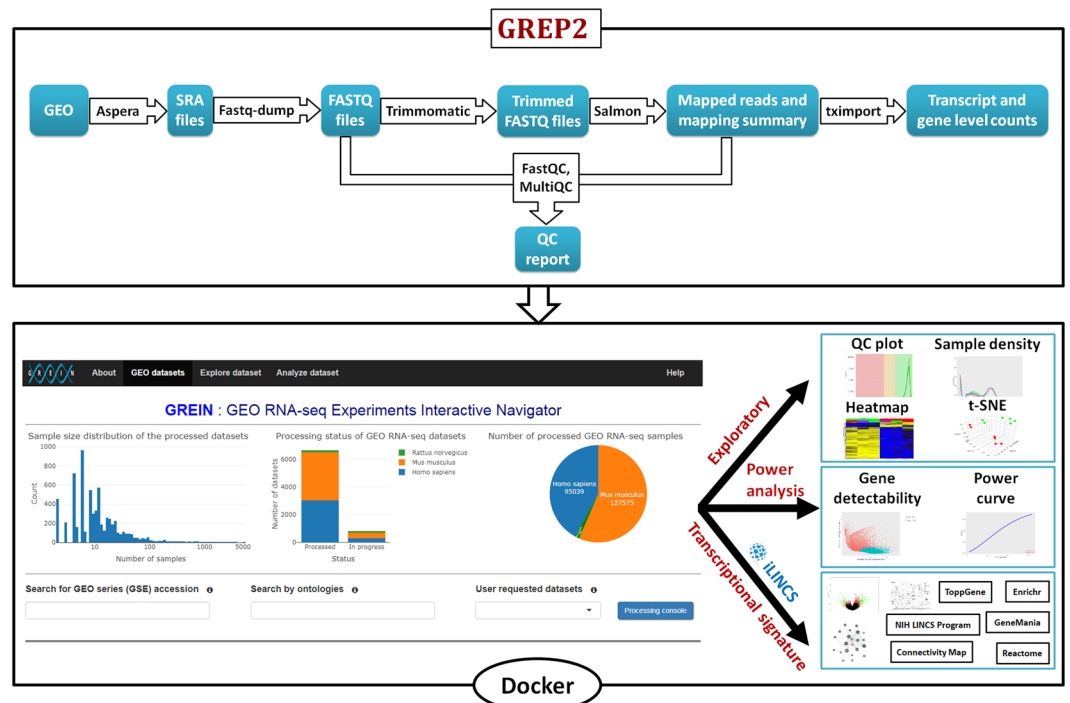


Figure 1. Schematic workflow of GREP2, web interface and outputs of GREIN. GEO datasets are systematically processed using GREP2 pipeline and stored within the back-end dataset library. GUI-driven GREIN workflows facilitate comprehensive analysis and visualization of processed datasets.

Results

The conceptual outline of GREIN is showed in Fig. 1. Individual RNA-seq datasets are processed by the GREP2 pipeline and stored locally as R Expression Sets. User can access and analyze preprocessed datasets via GREIN graphical user interface (GUI) or submit for processing datasets that have not yet been processed. GUI-driven

workflows facilitate examination and visualization of data, statistical analysis, transcriptional signature construction, and systems biology interpretation of differentially expressed (DE) genes. Both GREIN and the back-end pipeline (GREP2) are written in R and released as Docker container and R package respectively. Graphical user interfaces for GREIN are implemented in Shiny¹⁶, a web framework for building dynamic web applications in R. The web instance at <https://shiny.ilincs.org/grein> is deployed via robust Docker swarm of load-balanced Shiny servers. The complete GREIN infrastructure, including processing pipeline is deployed via Docker containers.

User friendly GUI driven workflows in GREIN facilitate typical reuse scenarios for RNA-seq data such as examination of quality control measures and visualization of expression patterns in the whole dataset, sample size and power analysis for the purpose of informing experimental design of future studies, statistical differential gene expression, gene list enrichment, and network analysis. Besides standard two-group comparison, the differential gene expression analysis module also supports fitting of a generalized linear model that accounts for covariates or batch-effects. The interactive visualization and exploration tools implemented include cluster analysis, interactive heatmaps, principal component analysis (PCA), t-distributed stochastic neighbor embedding (t-SNE), etc. (Supplementary Table S1). User can also search for ontological annotations of human RNA-seq samples and datasets provided by the MetaSRA project¹³. Each processed human RNA-seq sample is labelled with MetaSRA mapping of biomedical ontologies including Disease Ontology, Cell Ontology, Experimental Factor Ontology, Cellosaurus, and Uberon. Biological interpretation of differential gene expressions is aided by direct links to other online tools for performing typical post-hoc analyses such as the gene list and pathway enrichment analysis and the network analysis of differentially expressed (DE) genes. The connection to these analytical web services is implemented by submitting the differential gene expression signature (i.e., the list of average changes in gene expression and associated p-values for all/up/down regulated genes analyzed) to iLINC¹⁷ (Integrative LINC). iLINC also provides the signatures connectivity analysis for recently released Connectivity Map L1000 signatures¹⁸. Detailed step-by-step instructions about GREIN analysis workflows are provided in the Supplementary Material and 'Help' section in GREIN.

Key functionalities. *Search or submit for processing.* User can either search for an already processed GEO data set in the 'Search for GEO series (GSE accession)' box or submit a dataset for processing if the dataset is not already processed (Supplementary Fig. S2). At this point in time, the vast majority of GEO human, mouse, and rat RNA-seq datasets have been preprocessed and the user-submission of GEO datasets for processing will be required only occasionally. User can check the processing status of the requested dataset in the 'Processing console' tab (Supplementary Fig. S3). Other search options include keyword search through metadata of the datasets and search samples through biomedical ontologies via MetaSRA ontological annotations.

Explore dataset. GREIN allows access to both raw and normalized (counts per million and transcript per million) gene and transcript level data. GREIN comes with several interactive and customizable tools to visualize expression patterns such as interactive heatmaps of clustered genes and samples, density plots for all or a subset of samples, between and within group variability analysis through 2D and 3D dimensionality reduction analyses and visualizations such as PCA and t-SNE (Fig. 2). User can also visualize expression profile of each gene separately (Supplementary Fig. S6).

Quality control. The quality of the RNA-seq data in public repositories continues to be a major problem. In a recent study by Deelen *et al.*¹⁹, more than half of 65,000 processed public RNA-seq samples had to be removed due to QC issues. Rather than removing samples, GREIN provides a comprehensive quality control (QC) report of raw sequence data and sequence mapping for each sample (Supplementary Fig. S7), and allows user to make a decision about which samples should be excluded from downstream analyses.

Statistical power analysis. The power analysis module in GREIN facilitates calculation and visualization of statistical power of detecting differentially expressed genes in future studies utilizing similar biological samples. Estimating appropriate sample size for future studies with similar biological samples is often the key motivating factor in re-analysis of RNA-seq data. Power analysis also facilitates the post-hoc analysis of false negative rates in the current dataset. The lack of statistical power and differences in statistical power between genes can produce false negative results leading to wrong conclusions²⁰. The 'Power curve' segment provides power estimates for different number of samples based on a single gene (Fig. 3A). User can modify the default values of the parameters. The 'Detectability of genes' plot visualizes power estimate of each of the genes based on the selected groups and gene-wise dispersion (Fig. 3B). Mean coverage of the genes are plotted against their biological variability and are displayed in two sets based on their detectability status (power ≥ 0.8 and power < 0.8).

Differential gene expression. Creating and interpreting differential gene expression signature is a typical analysis scenario in RNA-seq experiments. With GREIN, user can create a signature by comparing gene expression between two groups of samples with or without adjustments for experimental covariates or batch effects. GREIN can handle complex experimental designs by providing the flexibility of rearranging groups and sub-groups or selecting specific samples. Differential expression signature can be visualized via interactive graphics that include heatmap of top differentially deregulated genes (Supplementary Fig. S15) ranked by false discovery rate (FDR), log fold change vs. log average expression (MA) plot (Supplementary Fig. S16), and gene detectability plot (Supplementary Fig. S17). Differential expression signature, with or without accounting for potentially false negative results, can be directly exported to iLINC for enrichment and connectivity analysis.

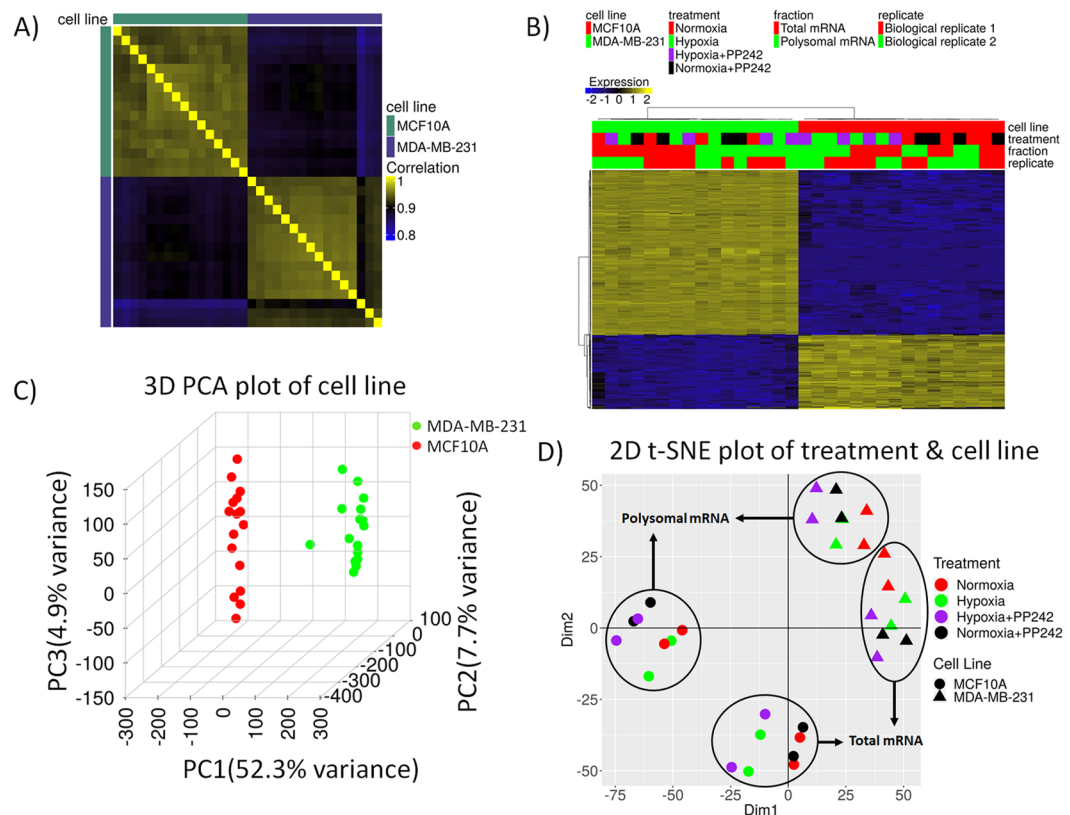


Figure 2. Exploratory analysis plots in GREIN. (A) Correlation heatmap shows a higher correlation within cell lines and low correlation between cell lines. Generally high correlations within each cell line indicate high quality of transcriptional profiles. (B) Hierarchical clustering based on Pearson correlation of top 500 most variable genes based on median absolute deviation as the variability measure. Data is normalized and centered to the mean. (C) Three-dimensional principal component analysis plot of the cell lines. (D) Two-dimensional t-SNE plot of treatment condition and cell line shows clear separation of the cell lines, and then the RNA fractions indicating two dominant sources of the variability between RNA-seq profiles.

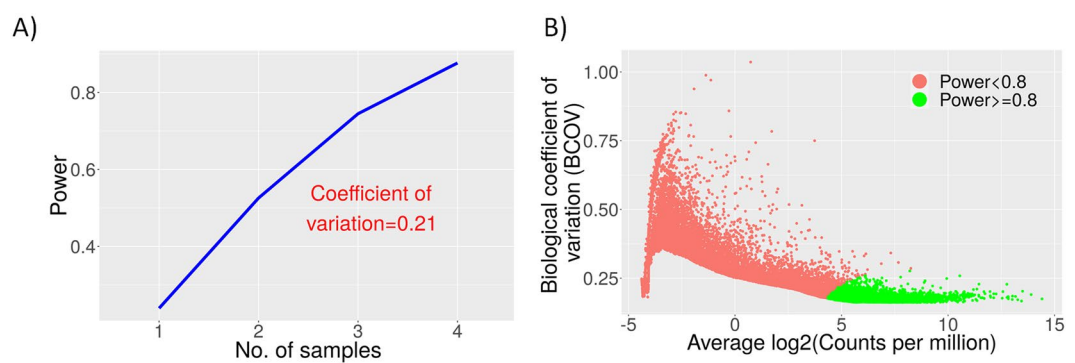


Figure 3. Power analysis for assessing transcriptional changes in non-malignant MCF10A cell line. (A) Single-gene based power estimates for different number of samples in each group with a minimal fold change of 2 and statistical significance $\alpha = 0.01$. (B) Gene-wise detectability on the log₂CPM-BCOV plane with FDR ≤ 0.1 and two samples in each group.

Use case: Analysis of transcriptional and translational regulation of hypoxia in non-malignant breast epithelial and triple-negative breast cancer cell lines. We demonstrate the usage of GREIN by re-analyzing a recently published GEO RNA-seq data (GSE104193). Sesé *et al.*²¹ examined the transcriptional and translational regulation of hormone-refractory triple-negative breast cancer (TNBC) subtype under a combination of hypoxia and mTOR (mechanistic target of rapamycin) inhibitor treatment. In particular, the authors analyzed the expression profiles of TNBC (MDA-MB-231) and non-malignant breast epithelial (MCF10A) cells exposed to normoxic (21% O₂) and hypoxic (0.5% O₂) conditions and/or treated with an mTORC1 and -2

Cell line	Comparison groups	Average sequencing depth (in million)	Common BCOV	Power (2 samples)	Power (3 samples)	Power (4 samples)
MCF10A	Hypoxia and normoxia	41.72	0.21	0.52	0.74	0.87
MCF10A	Hypoxia + PP242 and normoxia	38.18	0.31	0.28	0.44	0.59
MDA-MB-231	Hypoxia and normoxia	27.86	0.24	0.38	0.73	0.86
MDA-MB-231	Hypoxia + PP242 and normoxia	29.52	0.19	0.51	0.73	0.86

Table 2. Statistical power analysis to assess transcriptional changes in malignant and non-malignant cell lines. With a minimal fold change of 2 between the groups and statistical significance at $\alpha = 0.01$, all the comparisons are under-powered with two samples in each group. However, power increases as we increase the sample size.

inhibitor PP242. Each of the samples were sequenced for total (T) and polysome-bound (P) mRNA. The dataset contains 32 samples, representing two biological replicates for each combination of cell line, oxygen level, treatment status, and mRNA fraction.

Exploratory analysis of the processed dataset in GREIN (Fig. 2) shows that the strongest source of variation in between samples comes from differences between the two cell lines. This is re-enforced by the correlation analysis of full expression profiles (Fig. 2A), the hierarchical clustering of top 500 highly variable genes based on median absolute deviation (Fig. 2B), 3D PCA plot of the samples (Fig. 2C), and the 2D t-SNE plot (Fig. 2D). Furthermore, high correlations between expression profiles for the same cell line (Fig. 2A) indicates good signal-to-noise in the gene expression measurements. The additional substructure of data indicated by the 2D t-SNE plot has been examined by painting samples according to different attributes (Supplementary Fig. S1). This analysis revealed that separations within each cell line are induced by different mRNA fractions and then differences between experimental conditions.

Next, we used GREIN to perform statistical power analysis based on the pattern of biological variability observed in this dataset. We considered transcriptional profiles of each cell line exposed to hypoxia and treated with or without PP242 which leads to four comparisons. Assuming an expression difference of at least two-fold between the groups, at the statistical significance of $\alpha = 0.01$, and with only two replicates in each group, statistical power of a gene to be detected as differentially expressed is below 0.55 in all the comparisons (Table 2). Our analysis indicates that one would need four replicates per group to achieve 80% power detecting two-fold change in expression (Table 2 and Fig. 3A). In a typical RNA-seq experiment, a sequencing depth of 20–30 million is sufficient to quantify gene expression for almost all genes^{4,22} which is also evident in this dataset. We also evaluated statistical power of each gene to be detected as differentially expressed from the ‘Detectability of genes’ plot. Average log of counts per million (CPM) values of the genes were plotted against gene-wise biological coefficient of variation (BCOV) and power was calculated for the corresponding genes (Fig. 3B). A controlled false discovery rate of 0.05 and expected percentage of true positives of 10% was used to estimate statistical significance. We define a gene to be detectable as differentially expressed in hypoxic condition if its power is 0.8 or above. As expected, there exists an inverse relationship between BCOV and power (Fig. 3B). Also, power to detect differential expression of a gene increases with a higher log CPM or effect size.

One of the goals of the study was to analyze transcriptional changes in hypoxic and normoxic conditions with and without PP242 treatment in both MCF10A and MDA-MB-231 cell lines. We created transcriptional signatures of hypoxia and hypoxia + PP242 in total mRNA by differential expression analysis between hypoxia and hypoxia + PP242 samples respectively against the control samples while adjusting for batch effect by treating ‘replicate’ as a covariate, for each cell line separately. We found a higher number of genes differentially expressed (DE) in MCF10A cell lines compared to MDA-MB-231 in both hypoxia and hypoxia + PP242 (Fig. 4A) indicating that perhaps the tumor cell line is better equipped to deal with hypoxia. This analysis also showed that most non-differentially expressed genes are also not detectable, indicating that they may represent false negative results. This is in accordance with the power analysis showing that 4 samples per group would be needed to consistently identify differentially expressed genes with average BCOV. To identify lower expressed genes an even higher sample size would be required.

To interpret differentially expressed genes in terms of affected biological pathways, we submitted the differential gene expression signatures of hypoxia to online enrichment tools (DAVID²³, ToppGene²⁴, Enrichr²⁵, and Reactome²⁶) via iLINCS. The submitted signatures included a combined list of DE and NDE&DT genes representing likely true positive and true negatives. Genes were selected based on a cutoff of 0.7 and 0.01 for statistical power and FDR respectively. Figure 5 illustrates the enrichment results obtained from ToppGene for the MCF10 hypoxia signature. Significantly enriched (FDR < 0.05) top 10 gene ontology (GO) categories from ToppGene and DAVID functional annotation tool include response to hypoxia, response to decreased oxygen levels, angiogenesis, regulation of cell proliferation, oxidation-reduction process, and response to abiotic stimulus that are common in both cell lines (Supplementary Table S2 and Supplementary Table S3). Most of these categories are consistent with the original study. In addition, ToppGene suite identified hypoxia induced factor (HIF-1-alpha) transcription factor network that was activated in both cell lines (Supplementary Table S4 and Supplementary Table S5).

Finally, we utilized GREIN connection with iLINCS to “connect” the uploaded signature with LINCS²⁷ consensus (CGS) gene knockdown signatures¹⁸. We found 3,727 LINCS consensus gene knockdown signatures that were significantly (pValue < 0.05) connected with our uploaded signature. The target genes of top 100

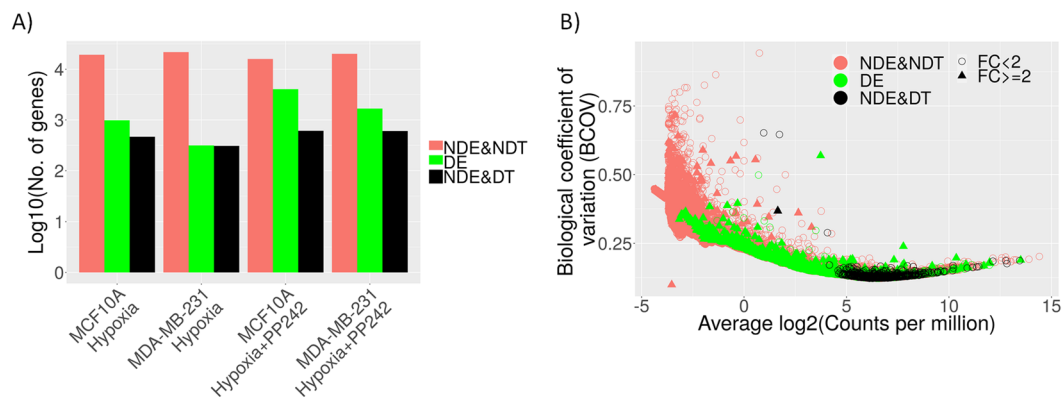


Figure 4. Differential expression and detectability of the genes. **(A)** The number of genes (log10 scale) not differentially expressed and not detectable (NDE&NDT), differentially expressed (DE), and not differentially expressed but detectable (NDE&DT) in the comparisons with normoxia for total mRNA fraction. We call a gene detectable (DT) if its power ≥ 0.8 and differentiable if FDR < 0.05 . **(B)** The gene detectability plot for the first comparison (MCF10A and hypoxia) which visualizes the above-mentioned list of genes along with their respective fold changes (FC).

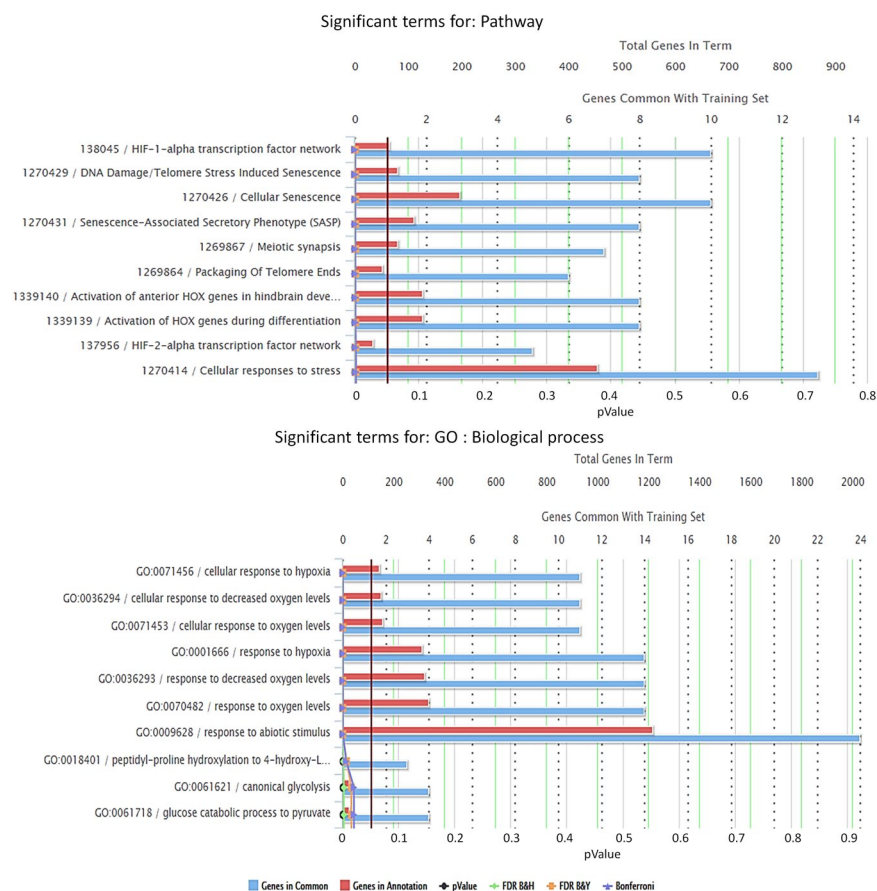


Figure 5. Snapshot of some of the significant pathway and gene ontology (GO) categories from ToppGene via iLINC. These categories are found in the comparison between hypoxia and normoxia in MCF10A cell line using a combined list of DE and NDE&DT genes. The red vertical line is the selected cutoff of 0.05.

connected signatures were selected for further enrichment analysis. We found cellular response to hypoxia and regulation of Hypoxia-inducible Factor (HIF) by oxygen in the list of top 10 activated pathways in both cell lines (Supplementary Table S6 and Supplementary Table S7). While this analysis yields similar enriched functional categories as the initial enrichment analysis, it complements the original analysis by implicating several target genes

that are not differentially expressed although they are sufficiently highly expressed to be detectable according to our power analyses. Tying these two results together implicates these genes as potential higher level regulators of the response to hypoxia.

Discussion

The combination of the access to a vast number of preprocessed mRNA expression data and the diversity of the analytical tools implemented makes GREIN a unique new resource for reuse of public domain RNA-seq data. GREIN not only removes technical barriers in reusing GEO RNA-seq data for biomedical scientists, but also facilitates easy assessment of the validity and reproducibility of the analysis results while uncovering new insights of the experiments. Our case-study analysis indicated that GREIN can be used for quickly reproducing results of the published studies as well as for gleaning additional insights that go beyond the original analysis. The complete analysis can be reproduced in less than 10 minutes. The web instance deployed on our server provides no-overhead use and requires no technical expertise. Open source R packages and the Docker container provide a flexible and transparent way for computationally savvy users to deploy the complete infrastructure locally with very little effort.

The GUI-driven analysis workflows implemented by GREIN covers a large portion of use cases for RNA-seq data analysis, making it the only tool that a scientist may need to meaningfully re-analyze GEO RNA-seq data. In Table 1 we provide a comprehensive comparison of GREIN with existing resources in terms of data processed and implemented analysis tools (Recount2⁷, ARCHS4⁸, Toil⁹, Skymap¹⁰, and Expression Atlas¹¹). Two key features distinguish GREIN from other currently available resources, it provides access to an exhaustive collection of RNA-seq samples in GEO, and it provides by far the most comprehensive set of interactive tools for analyzing RNA-seq datasets. Of existing tools, ARCHS4 and SkyMap both aim to provide access to the similar set of processed samples, but lack the interactive analytical toolbox. Other resources provide access to fewer processed samples and interactive analysis tools. In addition to standard analysis and visualization tools, the GREIN's power analysis workflow provides means to assess the statistical reasons for detecting or not-detecting specific differentially expressed genes. This kind of analyses is not common in the standard RNA-seq pipelines, but the results can be extremely useful when assessing false negative results. For example, when performing gene list enrichment analysis of differentially expressed genes, genes that do not meet a statistical significance cut-off are considered not differentially expressed. However, our power analysis indicates that the vast majority of these genes may simply be below detection limits for the available number of samples. GREIN allows user to distinguish between genes not detectable due to sample size limitations and/or their low expression levels, vs genes that are expressed at high enough expression levels but were not differentially expressed. This kind of resolution allows for more nuanced interpretation of analysis results. Off-line use of GREP2 and GREIN packages, as well as flexible export options enable additional analyses of GREIN-processed and pre-analyzed datasets.

Methods

GREIN back-end pipeline. To consistently process GEO RNA-seq datasets through a robust and uniform system, we have developed GEO RNA-seq experiments processing pipeline (GREP2) which is available as an R package in CRAN. Both GREP2 and GREIN are simultaneously running on different Docker containers. The whole processing workflow can be summarized in the following steps:

1. Obtain GEO series accession ID (series type: Expression profiling by high throughput sequencing) from GEO that contain at least two human, mouse, or rat RNA-seq samples. We then retrieve metadata for each GEO series accession using Bioconductor package GEOquery²⁸. We also obtain metadata files from the Sequence Read Archive (SRA) to get corresponding run information and merge both the GEO and SRA metadata.
2. Download corresponding experiment run files from SRA using 'ascp' utility of Aspera Connect²⁹ and convert them into FASTQ file format using NCBI SRA toolkit³⁰. All the downloaded files are stored in the local repository until processed.
3. Run FastQC³¹ on each FASTQ file to generate QC reports and remove adapter sequences if necessary using Trimmomatic³².
4. Quantify transcript abundances by mapping reads to reference transcriptome using Salmon³³.
5. Transcript level abundances are then then summarized to gene level using Bioconductor package tximport³⁴. We use *lengthScaledTPM* option in the summarization step which gives estimated counts scaled up to library size while considering for transcript length. Gene annotation for Homo sapiens (GRCh38), Mus musculus (GRCm38), and Rattus norvegicus (Rnor_6.0) are obtained from Ensemble³⁵ (release-91).
6. Compile FastQC reports and Salmon log files into a single interactive HTML report using MultiQC³⁶.

Power analysis. The power analysis in GREIN is performed using the Bioconductor package RNASeqPower⁴ which uses the following formula:

$$\left(z_{1-\frac{\alpha}{2}} + z_{\beta}\right)^2 = \frac{n \log_e \Delta^2}{2\left(\frac{1}{\mu} + \sigma^2\right)\#} \quad (1)$$

where, α is the target false positive rate, β is the target false negative rate or $1-\beta$ is power, n is the sample size, Δ is the effect size, μ is the average sequencing depth, and σ is the biological coefficient of variation (BCOV)

calculated as the square root of the dispersion. We use common dispersion and tagwise dispersion estimates from Bioconductor package *edgeR*³⁷ for computing power of a single gene and multiple genes respectively.

Typically, thousands of genes are tested simultaneously for differential expression in RNA-seq experiments. Therefore, the above method for estimating power needs further adjustment to correct for multiple testing. Jung *et al.*³⁸ derived an FDR correction formula and consequently the significance level (α^*) to calculate sample size for microarray data in the following form:

$$\alpha^* = \frac{r_1 f}{m_0(1 - f)} \# \quad (2)$$

where, r_1 is the expected number of true positives in m_1 rejected null hypotheses, m_0 denote the number of genes for which null hypotheses are true, and f implies desired FDR level. Hence, to calculate power for each of the genes, we replace α with α^* in eq. (1).

Differential expression analysis. GREIN uses negative binomial generalized linear model as implemented in *edgeR* to find differentially expressed genes between sample groups. Data is normalized using trimmed mean of M-values (TMM) as implemented in *edgeR*. All the analyses are based on CPM values and genes are filtered at the onset with a cutoff of CPM > 0 in m samples, where m is the minimum sample size in any of the groups. Besides two-group comparison, GREIN also supports adjustment for experimental covariates or batch effects. A design matrix is constructed with the selected variable and groups. We use gene-wise negative binomial generalized linear models with quasi-likelihood tests and gene-wise exact tests to calculate differential expression between groups with and without covariates respectively. P-values are adjusted for multiple testing correction using Benjamini-Hochberg method³⁹. Interactive visualization of the differentially expressed genes is also available via heatmap of the top ranked genes, MA plot, and gene detectability plot.

References

- Edgar, R., Domrachev, M. & Lash, A. E. Gene Expression Omnibus: NCBI gene expression and hybridization array data repository. *Nucleic Acids Res.* **30**, 207–210, <https://doi.org/10.1093/nar/30.1.207> (2002).
- Leinonen, R., Sugawara, H. & Shumway, M. & on behalf of the International Nucleotide Sequence Database, C. The Sequence Read Archive. *Nucleic Acids Res.* **39**, D19–D21, <https://doi.org/10.1093/nar/gkq1019> (2011).
- Rung, J. & Brazma, A. Reuse of public genome-wide gene expression data. *Nat. Rev. Genet.* **14**, 89, <https://doi.org/10.1038/nrg3394> (2012).
- Hart, S. N., Therneau, T. M., Zhang, Y., Poland, G. A. & Kocher, J.-P. Calculating Sample Size Estimates for RNA Sequencing Data. *J. Comput. Biol.* **20**, 970–978, <https://doi.org/10.1089/cmb.2012.0283> (2013).
- Conesa, A. *et al.* A survey of best practices for RNA-seq data analysis. *Genome Biol.* **17**, 13, <https://doi.org/10.1186/s13059-016-0881-8> (2016).
- Huber, W. *et al.* Orchestrating high-throughput genomic analysis with Bioconductor. *Nat Meth* **12**, 115–121, <https://doi.org/10.1038/nmeth.3252> (2015).
- Collado-Torres, L. *et al.* Reproducible RNA-seq analysis using recount2. *Nat. Biotechnol.* **35**, 319, <https://doi.org/10.1038/nbt.3838> (2017).
- Lachmann, A. *et al.* Massive mining of publicly available RNA-seq data from human and mouse. *Nat. Commun.* **9**, 1366, <https://doi.org/10.1038/s41467-018-03751-6> (2018).
- Vivian, J. *et al.* Toil enables reproducible, open source, big biomedical data analyses. *Nat. Biotechnol.* **35**, 314, <https://doi.org/10.1038/nbt.3772> (2017).
- Tsui, B. Y., Dow, M., Skola, D. & Carter, H. Extracting allelic read counts from 250,000 human sequencing runs in Sequence Read Archive. *bioRxiv*, 386441, <https://doi.org/10.1101/386441> (2018).
- Papatheodorou, I. *et al.* Expression Atlas: gene and protein expression across multiple studies and organisms. *Nucleic Acids Res.* **46**, D246–D251, <https://doi.org/10.1093/nar/gkx1158> (2018).
- Pimentel, H., Sturmfels, P., Bray, N., Melsted, P. & Pachter, L. The Lair: a resource for exploratory analysis of published RNA-Seq data. *BMC Bioinformatics* **17**, 490, <https://doi.org/10.1186/s12859-016-1357-2> (2016).
- Bernstein, M. N., Doan, A. & Dewey, C. N. MetaSR: normalized human sample-specific metadata for the Sequence Read Archive. *Bioinformatics* **33**, 2914–2923, <https://doi.org/10.1093/bioinformatics/btx334> (2017).
- Team, R. C. R language definition. *Vienna, Austria: R foundation for statistical computing* (2013).
- Merkel, D. Docker: lightweight linux containers for consistent development and deployment. *Linux Journal* **2014**, 2 (2014).
- Chang, W., Cheng, J., Allaire, J. J., Xie, Y. & McPherson, J. Shiny: web application framework for R. *R package version 0.11.1*, 106 (2015).
- iLINCS <http://www.ilincs.org/> (accessed, 5 October 2018).
- Subramanian, A. *et al.* A Next Generation Connectivity Map: L1000 Platform and the First 1,000,000 Profiles. *Cell* **171**, 1437–1452, e1417, <https://doi.org/10.1016/j.cell.2017.10.049> (2017).
- Deelen, P. *et al.* Improving the diagnostic yield of exome-sequencing, by predicting gene-phenotype associations using large-scale gene expression analysis. *bioRxiv*, 375766, <https://doi.org/10.1101/375766> (2018).
- Norris, A. W. & Kahn, C. R. Analysis of gene expression in pathophysiological states: Balancing false discovery and false negative rates. *Proc. Natl. Acad. Sci. USA* **103**, 649 (2006).
- Sesé, M. *et al.* Hypoxia-mediated translational activation of ITGB3 in breast cancer cells enhances TGF- β signaling and malignant features *in vitro* and *in vivo*. *Oncotarget* **8**, 114856–114876, <https://doi.org/10.18632/oncotarget.23145> (2017).
- Wang, Y. *et al.* Evaluation of the coverage and depth of transcriptome by RNA-Seq in chickens. *BMC Bioinformatics* **12**, S5, <https://doi.org/10.1186/1471-2105-12-S10-S5> (2011).
- Huang, D. W., Sherman, B. T. & Lempicki, R. A. Systematic and integrative analysis of large gene lists using DAVID bioinformatics resources. *Nat. Protoc.* **4**, 44, <https://doi.org/10.1038/nprot.2008.211> (2008).
- Chen, J., Bardes, E. E., Aronow, B. J. & Jegga, A. G. ToppGene Suite for gene list enrichment analysis and candidate gene prioritization. *Nucleic Acids Res.* **37**, W305–W311, <https://doi.org/10.1093/nar/gkp427> (2009).
- Chen, E. Y. *et al.* Enrichr: interactive and collaborative HTML5 gene list enrichment analysis tool. *BMC Bioinformatics* **14**, 128, <https://doi.org/10.1186/1471-2105-14-128> (2013).
- Joshi-Tope, G. *et al.* Reactome: a knowledgebase of biological pathways. *Nucleic Acids Res.* **33**, D428–D432, <https://doi.org/10.1093/nar/gki072> (2005).

27. Keenan, A. B. *et al.* The Library of Integrated Network-Based Cellular Signatures NIH Program: System-Level Cataloging of Human Cells Response to Perturbations. *Cell Syst.* **6**, 13–24, <https://doi.org/10.1016/j.cels.2017.11.001> (2018).
28. Davis, S. & Meltzer, P. S. GEOquery: a bridge between the Gene Expression Omnibus (GEO) and BioConductor. *Bioinformatics* **23**, 1846–1847, <https://doi.org/10.1093/bioinformatics/btm254> (2007).
29. Aspera Connect <https://www.asperasoft.com> (accessed, 5 October 2018).
30. NCBI SRA toolkit <http://trace.ncbi.nlm.nih.gov/Traces/sra/sra.cgi?view=software> (accessed, 5 October 2018).
31. Andrews, S. FastQC: a quality control tool for high throughput sequence data <http://www.bioinformatics.babraham.ac.uk/projects/fastqc> (2010).
32. Bolger, A. M., Lohse, M. & Usadel, B. Trimmomatic: a flexible trimmer for Illumina sequence data. *Bioinformatics* **30**, 2114–2120, <https://doi.org/10.1093/bioinformatics/btu170> (2014).
33. Patro, R., Duggal, G., Love, M. I., Irizarry, R. A. & Kingsford, C. Salmon provides fast and bias-aware quantification of transcript expression. *Nat. Methods* **14**, 417, <https://doi.org/10.1038/nmeth.4197> (2017).
34. Sonesson, C., Love, M. I. & Robinson, M. D. Differential analyses for RNA-seq: transcript-level estimates improve gene-level inferences. *F1000Res.* **4**, 1521, <https://doi.org/10.12688/f1000research.7563.2> (2015).
35. Zerbino, D. R. *et al.* Ensembl 2018. *Nucleic Acids Res.* **46**, D754–D761, <https://doi.org/10.1093/nar/gkx1098> (2018).
36. Ewels, P., Magnusson, M., Lundin, S. & Käller, M. MultiQC: summarize analysis results for multiple tools and samples in a single report. *Bioinformatics* **32**, 3047–3048, <https://doi.org/10.1093/bioinformatics/btw354> (2016).
37. Robinson, M. D., McCarthy, D. J. & Smyth, G. K. edgeR: a Bioconductor package for differential expression analysis of digital gene expression data. *Bioinformatics* **26**, 139–140, <https://doi.org/10.1093/bioinformatics/btp616> (2010).
38. Jung, S.-H. Sample size for FDR-control in microarray data analysis. *Bioinformatics* **21**, 3097–3104, <https://doi.org/10.1093/bioinformatics/bti456> (2005).
39. Benjamini, Y. & Hochberg, Y. Controlling the False Discovery Rate: A Practical and Powerful Approach to Multiple Testing. *J. Roy. Stat. Soc. Ser. B. (Stat. Method.)* **57**, 289–300 (1995).

Acknowledgements

This work was supported by the grants from National Institutes of Health: LINCS DCIC (U54HL127624) and Center for Environmental Genetics (P30ES006096).

Author Contributions

N.A.M. developed the pipeline and web application, M.M. conceived the project, supervised software development and data processing, M.M. and N.A.M. wrote the manuscript, M.F.N. developed and maintain the Docker containers, M.P. and M.K. maintain the web server and implemented APIs for connecting with iLINCS. All authors reviewed the manuscript.

Additional Information

Supplementary information accompanies this paper at <https://doi.org/10.1038/s41598-019-43935-8>.

Competing Interests: The authors declare no competing interests.

Publisher's note: Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.



Open Access This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons license, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons license, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons license and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this license, visit <http://creativecommons.org/licenses/by/4.0/>.

© The Author(s) 2019