

# Technical Support Center: Normalizing RNA-seq Data from NCI-H38 Experiments

**Author:** BenchChem Technical Support Team. **Date:** December 2025

## Compound of Interest

Compound Name: NCD38

Cat. No.: B609494

[Get Quote](#)

This guide provides troubleshooting advice and frequently asked questions for researchers, scientists, and drug development professionals working with RNA-seq data, particularly from experiments involving the NCI-H38 cancer cell line.

## Frequently Asked Questions (FAQs)

### Q1: What is RNA-seq data normalization, and why is it crucial?

A: RNA-seq data normalization is a critical step in the analysis pipeline that aims to correct for technical variations in the data, allowing for more accurate comparisons of gene expression levels between different samples.<sup>[1][2]</sup> Without normalization, variations in factors like sequencing depth (the total number of reads per sample) and RNA composition can lead to erroneous conclusions about gene expression differences.<sup>[1][3]</sup> The primary goal is to minimize these technical variations while preserving the true biological differences.<sup>[1]</sup>

### Q2: What are the common sources of technical variation in RNA-seq data?

A: Several factors can introduce technical variability into RNA-seq experiments:

- **Sequencing Depth:** The total number of reads can vary significantly from sample to sample. A sample with more reads will naturally show higher gene counts, which is not necessarily due to a biological difference.<sup>[3][4]</sup>

- **Gene Length:** Longer genes will have more reads mapped to them than shorter genes, even if they are expressed at the same level.[3]
- **Library Composition:** Differences in the overall composition of the RNA population between samples can affect the quantification of individual genes.[5] For instance, if a few genes are very highly expressed in one sample, they will consume a large proportion of the sequencing reads, making other genes appear to have lower expression.
- **GC-Content:** Biases related to the GC-content of genes can affect how efficiently their corresponding transcripts are amplified and sequenced.[6]
- **Batch Effects:** When samples are processed in different batches or on different dates, systematic technical variations can be introduced that may obscure true biological differences.[7]

### Q3: What are the most common RNA-seq normalization methods?

A: There are several normalization methods, each with its own assumptions and ideal use cases. The most common methods can be broadly categorized as those for within-sample comparisons and those for between-sample comparisons.

- For within-sample comparison (comparing gene expression within a single sample):
  - **TPM (Transcripts Per Million):** This method normalizes for both gene length and sequencing depth.[3][8] It is generally considered a better method for comparing the proportion of reads mapped to a gene across different samples.[8]
- For between-sample comparison (required for differential expression analysis):
  - **RPKM/FPKM (Reads/Fragments Per Kilobase of transcript per Million mapped reads):** These were among the first methods developed to normalize for sequencing depth and gene length.[3] However, they have been shown to be less effective for comparing gene expression between samples due to issues with library composition.[1][9]
  - **CPM (Counts Per Million):** This method only accounts for sequencing depth and not gene length.[7] It is not recommended for comparing expression levels of different genes within

the same sample.

- DESeq2's Median of Ratios: This method calculates a size factor for each sample based on the median of the ratios of gene counts to a pseudo-reference sample (the geometric mean of each gene across all samples). This approach is robust to the presence of highly expressed genes and differences in library composition.<sup>[5][10]</sup>
- TMM (Trimmed Mean of M-values): Used by the edgeR package, this method is similar to DESeq2's in that it calculates a normalization factor for each sample. It is also designed to be robust to differences in library composition.

## Q4: Which normalization method should I choose for my NCI-H38 differential expression experiment?

A: For differential expression analysis, methods that account for library composition, such as DESeq2's median of ratios or edgeR's TMM, are generally recommended.<sup>[9]</sup> Studies have shown that normalized count data from methods like DESeq2 or TMM perform better than TPM and FPKM in clustering replicate samples and are more reproducible.<sup>[9]</sup> While TPM is useful for visualizing the expression of a gene across samples, the raw counts or DESeq2/edgeR normalized counts should be used as input for differential expression analysis software.<sup>[11]</sup>

## Q5: Can I directly compare TPM values across different experiments or studies?

A: Caution should be exercised when comparing TPM values across different experiments, especially if they were generated using different library preparation protocols or sequencing platforms.<sup>[1]</sup> Variations in these protocols can significantly impact the measurement of different RNA classes, making direct comparisons of TPM values unreliable.<sup>[1]</sup> It is generally safer to re-process the raw data from different studies through the same analysis pipeline to ensure consistency.

## Troubleshooting Guide

**Problem: My PCA plot does not show clear separation between my experimental groups after normalization.**

- Possible Cause 1: Insufficient biological signal. The experimental conditions may not have induced strong enough changes in gene expression to be the dominant source of variation in your data.
- Troubleshooting Steps:
  - Check for batch effects: Color your PCA plot by other metadata variables, such as the date of library preparation or the sequencing lane. If samples cluster by these technical variables instead of your biological condition, you may have a batch effect.
  - Examine a larger number of principal components: The separation may be apparent in later principal components (e.g., PC3 vs. PC4).
  - Use a different dimensionality reduction technique: Try t-SNE or UMAP to see if they reveal a better separation of your groups.

## Problem: I suspect a batch effect in my data.

- Possible Cause: Samples were processed at different times or by different technicians. This can introduce systematic, non-biological variation into the data.[\[7\]](#)
- Troubleshooting Steps:
  - Visualize the batch effect: Use PCA plots or hierarchical clustering to see if samples group by batch.
  - Correct for the batch effect:
    - Include batch in your statistical model: For differential expression analysis, you can often include the batch as a covariate in the design formula of tools like DESeq2 or limma. This will account for the variation due to the batch effect when testing for differences between your biological groups.
    - Use batch correction tools: For visualization or clustering, you can use tools like ComBat-seq or the removeBatchEffect function in limma to adjust the normalized expression data.

## Problem: My differential expression analysis yields a very large number of significant genes, or results that don't align with biological expectations.

- Possible Cause 1: Inappropriate normalization method. Using a method like RPKM/FPKM for differential expression can sometimes lead to an inflated number of false positives.[\[1\]](#)[\[9\]](#)
- Possible Cause 2: Confounding variables. A hidden technical or biological variable that is correlated with your experimental condition of interest might be driving the observed differences.
- Troubleshooting Steps:
  - Verify your normalization method: Ensure you are using a method appropriate for between-sample comparisons, such as DESeq2's median of ratios or TMM.[\[9\]](#)
  - Re-examine your experimental design and metadata: Look for any potential confounding factors. For example, were all the control samples processed in one batch and all the treated samples in another?
  - Use independent filtering: Tools like DESeq2 automatically perform independent filtering to remove genes with very low read counts, which can increase the power to detect differentially expressed genes.[\[12\]](#) Ensure this step is being performed.
  - Check for outliers: Use PCA plots and sample distance heatmaps to identify any outlier samples that may be skewing the results. Consider removing them if they are clear technical failures.

## Data Presentation

### Comparison of Common RNA-seq Normalization Methods

| Normalization Method          | Normalizes for Gene Length? | Normalizes for Sequencing Depth? | Accounts for Library Composition? | Best For   |
|-------------------------------|-----------------------------|----------------------------------|-----------------------------------|--|
| CPM (Counts Per Million)      | No                          | Yes                              | No                                | Gross estimation of gene expression; not for DE analysis |
| RPKM/FPKM                     | Yes                         | Yes                              | No                                | Within-sample gene comparison (with caution)             |
| TPM (Transcripts Per Million) | Yes                         | Yes                              | Partially                         | Within-sample comparison and visualization               |
| DESeq2 Median of Ratios       | No (uses raw counts)        | Yes                              | Yes                               | Differential Expression Analysis                         |
| edgeR TMM                     | No (uses raw counts)        | Yes                              | Yes                               | Differential Expression Analysis                         |

## Experimental Protocols

### Protocol: Standard RNA-seq Normalization and Differential Expression Workflow using DESeq2

This protocol outlines the key computational steps for normalizing RNA-seq count data and performing differential expression analysis.

- Input Data:
  - A "count matrix": a table where rows represent genes and columns represent samples, with the values being the raw read counts.

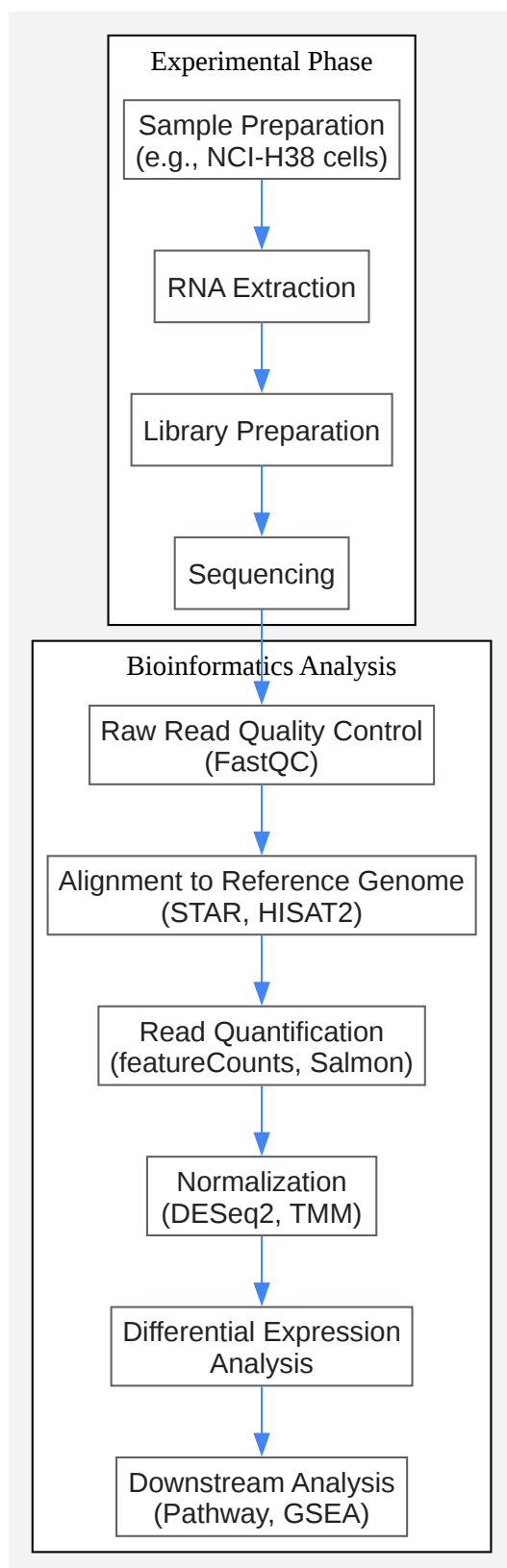
- A "metadata table": a table describing the experimental design, where rows correspond to the samples in the count matrix and columns contain variables of interest (e.g., treatment group, batch number).
- Create a DESeqDataSet Object:
  - In R, use the DESeqDataSetFromMatrix function to combine the count matrix and metadata into a single object.
  - Specify the "design formula" at this stage (e.g., ~ batch + condition), which tells DESeq2 which variables to test for and which to correct for.
- Pre-filtering (Optional but Recommended):
  - Remove genes with very low counts across all samples. A common approach is to keep only genes that have a count of at least 10 in a minimum number of samples.[\[12\]](#)
- Normalization:
  - Run the DESeq function. This single function will perform the following steps:
    - Estimate size factors: This is the normalization step, where DESeq2 calculates a size factor for each sample using the "median of ratios" method.[\[5\]](#)
    - Estimate gene-wise dispersions: This step models the variance of each gene's counts.
    - Fit a negative binomial model and perform statistical testing: This step identifies genes that are differentially expressed between your conditions of interest.
- Extracting Results:
  - Use the results function to generate a table of differentially expressed genes, which will include log2 fold changes, p-values, and adjusted p-values (to correct for multiple testing).
- Normalized Data for Visualization:
  - To obtain normalized counts for downstream applications like heatmaps or PCA plots, use functions like vst (variance stabilizing transformation) or rlog (regularized log

transformation) on the DESeqDataSet object. These transformations create homoscedastic data (i.e., the variance is similar across the range of mean values), which is important for many visualization and clustering methods.

## Visualizations

### Experimental and Analytical Workflow

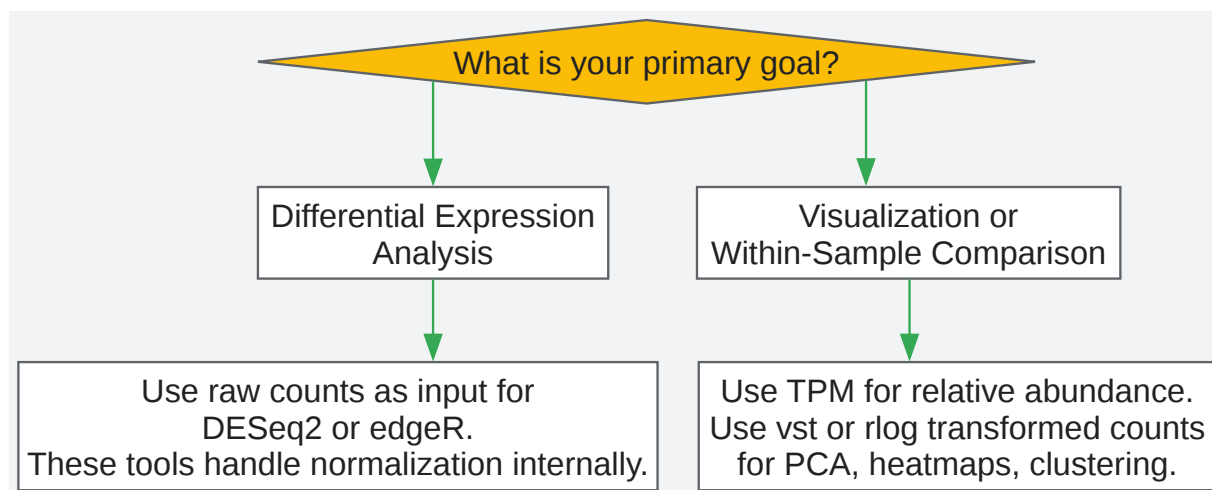




[Click to download full resolution via product page](#)

Caption: A typical workflow for an RNA-seq experiment.

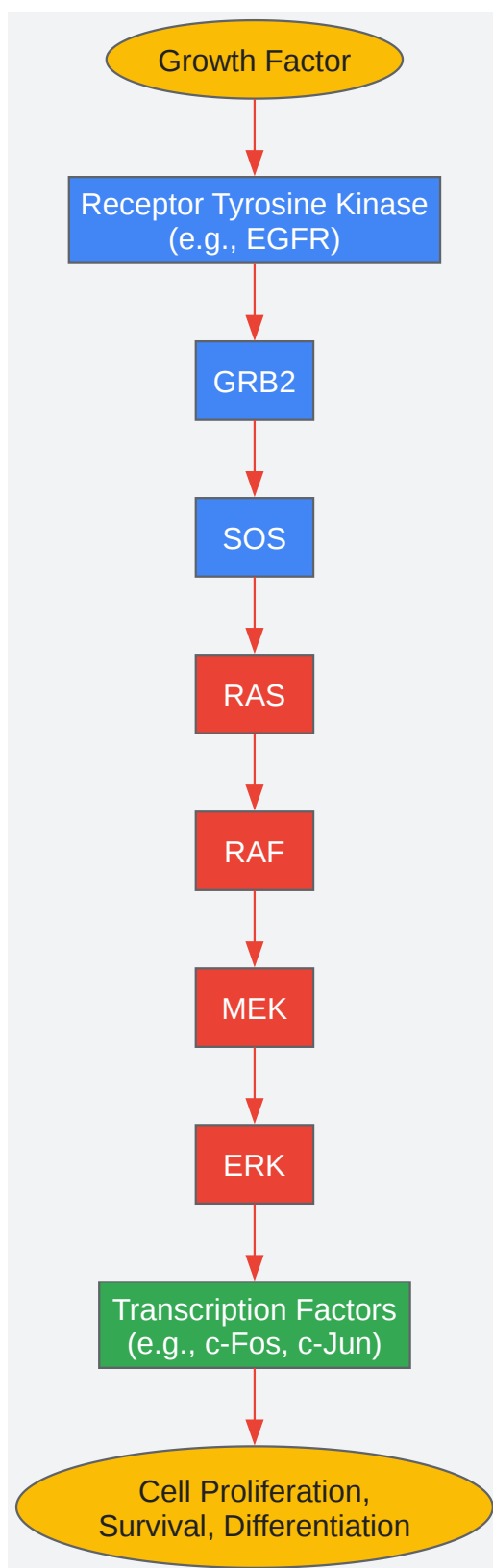
## Decision Tree for Choosing a Normalization Method



[Click to download full resolution via product page](#)

Caption: A guide to selecting an appropriate normalization strategy.

## MAPK/ERK Signaling Pathway



[Click to download full resolution via product page](#)

Caption: The MAPK/ERK pathway, often dysregulated in cancer.

**Need Custom Synthesis?**

BenchChem offers custom synthesis for rare earth carbides and specific isotopic labeling.

Email: [info@benchchem.com](mailto:info@benchchem.com) or [Request Quote Online](#).

## References

- 1. Normalizing Behaviors and RNASeq: A Tale of Good and Bad Practices | by Decode Box | Medium [medium.com]
- 2. rna-seqblog.com [rna-seqblog.com]
- 3. pluto.bio [pluto.bio]
- 4. youtube.com [youtube.com]
- 5. m.youtube.com [m.youtube.com]
- 6. discovery.researcher.life [discovery.researcher.life]
- 7. bigomics.ch [bigomics.ch]
- 8. youtube.com [youtube.com]
- 9. tpm-fpk-m-or-normalized-counts-a-comparative-study-of-quantification-measures-for-the-analysis-of-rna-seq-data-from-the-nci-patient-derived-models-repository - Ask this paper | Bohrium [bohrium.com]
- 10. youtube.com [youtube.com]
- 11. How important is normalization for RNA\_Seq data ? [biostars.org]
- 12. m.youtube.com [m.youtube.com]
- To cite this document: BenchChem. [Technical Support Center: Normalizing RNA-seq Data from NCI-H38 Experiments]. BenchChem, [2025]. [Online PDF]. Available at: [https://www.benchchem.com/product/b609494#normalizing-rna-seq-data-from-ncd38-experiments]

---

**Disclaimer & Data Validity:**

The information provided in this document is for Research Use Only (RUO) and is strictly not intended for diagnostic or therapeutic procedures. While BenchChem strives to provide

accurate protocols, we make no warranties, express or implied, regarding the fitness of this product for every specific experimental setup.

**Technical Support:** The protocols provided are for reference purposes. Unsure if this reagent suits your experiment? [[Contact our Ph.D. Support Team for a compatibility check](#)]

**Need Industrial/Bulk Grade?** [Request Custom Synthesis Quote](#)

## BenchChem

Our mission is to be the trusted global source of essential and advanced chemicals, empowering scientists and researchers to drive progress in science and industry.

### Contact

Address: 3281 E Guasti Rd  
Ontario, CA 91761, United States  
Phone: (601) 213-4426  
Email: [info@benchchem.com](mailto:info@benchchem.com)