

# sSNAPPY: a R/Bioconductor package for single-sample directional pathway perturbation analysis

Wenjun Liu<sup>1</sup>, Ville-Petteri Mäkinen<sup>2,3,4,5</sup>, Wayne D. Tilley<sup>1</sup>, and Stephen M. Pederson<sup>1,6,7</sup>

<sup>1</sup>Dame Roma Mitchell Cancer Research Laboratories, Adelaide Medical School, Faculty of Health and Medical Sciences, University of Adelaide, Adelaide, Australia

<sup>2</sup>Australian Centre for Precision Health, Cancer Research Institute, University of South Australia, Adelaide, Australia

<sup>3</sup>Computational and Systems Biology Program, Precision Medicine Theme, South Australian Health and Medical Research Institute, Adelaide, Australia

<sup>4</sup>Computational Medicine, Faculty of Medicine, University of Oulu, Oulu, Finland

<sup>5</sup>Center for Life Course Health Research, Faculty of Medicine, University of Oulu, Oulu, Finland

<sup>6</sup>Black Ochre Data Laboratories, Telethon Kids Institute, Adelaide, Australia

<sup>7</sup>John Curtin School of Medical Research, Australian National University, Canberra, Australia

**Abstract** When analysing RNA-Seq data, a common outcome is to detect biological pathways with significantly altered activity between the conditions under investigation. The most common strategies test for over-representation within pre-defined gene-sets for genes showing changed expression[1, 2], but without accounting for gene-gene interactions encoded by pathway topologies, and not being able to directly predict the directional change of pathway activity. To address these issues, we have developed a single-sample pathway perturbation analysis method *sSNAPPY*, now available as an R/Bioconductor package, which leverages pathway topology information to compute pathway perturbation scores, and predicts the potential direction of change across a set of pathways. Here, we demonstrate the use of *sSNAPPY* by applying the method to public scRNA-seq data, derived from ovarian cancer patient tissues collected before and after chemotherapy. Not only were we able to replicate results reported in the original study, but *sSNAPPY* was also able to detect significant perturbation of other biological processes, yielding far greater insight into the response to treatment. *sSNAPPY* represents a novel pathway analysis strategy that takes into consideration of pathway topology to predict impacted biology pathways, both within related samples and across treatment groups. In addition to not relying on the detection of differentially expressed genes, the method and associated R package offer important flexibility and provide powerful visualisation tools.

## Keywords

RNA-seq, pathway enrichment, R package, topology, KEGG, scRNA-seq

**R version:** R version 4.2.3 (2023-03-15)

**Bioconductor version:** 3.16

## Introduction

Using pathway enrichment analysis to gain biological insights from gene expression data is a pivotal step in the analysis and interpretation of RNA-seq data, for which numerous methods have been developed (reviewed in [3, 4]). Many existing methods tend to view pathways simply as a collection of gene names, as seen in those relying on the detection of differentially expressed genes and applying over-representation analysis (ORA) strategies, and those scoring all genes using functional class scoring (FCS), such as in Gene Set Enrichment Analysis (GSEA) [1], arguably the most widely-used approach. However, databases such as the Kyoto Encyclopaedia of Genes and Genomes (KEGG)[5] and WikiPathways[6] capture not only which genes are implicated in a certain biological process but also their interactions, activating or inhibitory roles, and their relative importance within the pathway, all of which are overlooked in ORA- and FCS-based approaches.

To fully utilise that additional information, the latest generation of pathway analysis approaches include many which are topology-based such as SPIA[7], DEGraph[8], NetGSA[9] and PRS[10], as well as others which explicitly model inter-gene correlations[11]. Despite differences in the null hypotheses tested across these approaches, overall, they have demonstrated enhanced sensitivity and specificity due to their abilities to take gene-gene interconnections into account[12, 13]. Nevertheless, most topology-based methods focus only on comparing activities of pathways between two treatment groups and cannot be used to score individual samples. However, in heterogenous data where more than one factor may be influencing our observations[14], incorporating scoring within paired samples may be desirable and may be able to reveal more nuanced insights. To address this gap, we present a Single-Sample Network and Pathway Perturbation analysis methodology called *sSNAPPY*, available as an R/Bioconductor package. This article defines how *sSNAPPY* computes changes in gene expression within paired samples, and propagates this through gene-set topologies to predict the perturbation in pathway activities within paired samples, before providing summarised results across an entire dataset which would include more robust levels of biological replication. The practical usage of the *sSNAPPY* R/Bioconductor package will be illustrated through the analysis of a public scRNA-seq dataset using pseudo-bulk strategies.

## Methods

### Implementation

*sSNAPPY* is an R package that has been reviewed and published on the open-source bioinformatics software platform Bioconductor with all source code available via GitHub. The methodology itself is topology-based, designed to compute directional, single-sample, pathway perturbation scores in gene expression datasets with a matched-pair, or nested design (eg. samples collected before and after treatment). This allows for the detection of pathway perturbations within all samples from a treatment group, but also within individual samples.

To run *sSNAPPY*, the only required data is a log-transformed expression matrix (e.g. logCPM) with matching sample metadata describing treatment groups and the nested structure. It is assumed that all pre-processing has been performed beforehand, such as the exclusion of low-signal genes or normalisation to minimise technical artefacts like GC-bias. The first step, performed internally by *sSNAPPY*, is to estimate sample-specific log fold-change ( $\delta_{ghi} = \mu_{ghi} - \mu_{g0i}$ ) for a treatment  $h$  across all genes  $g$  within each sample  $i$ , by subtracting expression estimates for the baseline samples  $\mu_{g0i}$  from those in the treatment group  $h$ . Since it has been shown that in RNA-seq data, genes with lower expression tend to have larger variance and larger estimates of change[15], we utilise a gene-level weighting strategy to de-emphasise logFC estimates for low-abundance genes. Gene-level weights  $w_g$  are obtained in a treatment-agnostic manner by fitting a loess curve through the relationship between observed gene-level variance ( $\hat{\sigma}_g$ ) and average signal ( $\bar{\mu}_{g..}$ ), and taking the inverse of the loess-predicted variance as the weight  $w_g = a/f(\bar{\mu}_{g..})$ , where  $f(\bar{\mu}_{g..})$  is the predicted value from the loess curve and the constant  $a$  ensures  $\sum w_g = 1$ . We then use these weighted estimates of logFC for calculation of all pathway perturbation scores.

*sSNAPPY* was built upon the group-level topology-based scoring algorithm initially proposed in R package SPIA[16] to propagate genes' changes in expression through pathway topologies to compute a perturbation score for each pathway. By modifying the algorithm to incorporate single-sample, weighted estimates of changes in expression we are able to quantify changes in a pathway within a given sample, and then model these across all samples within a treatment group. Thus, we define the single-sample perturbation score ( $S_{hip}$ ) for a given pathway  $p$ , sample  $i$  and treatment  $h$ :

$$S_{hip} = \sum_{g \in G_p} [S_{ghip} - \delta_{ghi}^*], \text{ where}$$

$$S_{ghip} = \delta_{ghi}^* + \sum_{g' \in U_{gp}} \beta_{gg'p} \frac{S_{g'hip}}{N_{g'p}}$$

where:

- $G_p$  represents the set of genes in pathway  $p$ , such that  $g \in G_p$
- $S_{ghip}$  is the gene-, treatment- and sample-specific perturbation score for pathway  $p$
- $\delta_{ghi}^* = w_g \delta_{ghi}$  is the weighted logFC of gene  $g$  as described above
- $U_{gp}$  is the subset of  $G_p$  containing only the genes directly upstream of gene  $g$
- $\beta_{gg'p}$  is the pair-wise gene-gene interactions[16] encoded by the topology matrix for genes  $g$  and  $g'$
- $N_{gp}$  is the number of downstream genes from any gene  $g$
- $S_{hip}$  is the accumulated pathway perturbation score for pathway  $p$  in treatment  $h$  for sample  $i$

The Bioconductor package `graphite`[17] provides functions that can be used to retrieve pathway topologies from a database and convert topology information to adjacency matrices. In order to streamline this process we have implemented a convenience function, where users only need to provide the name of the desired database to retrieve all topology information in the format required by the scoring algorithm with the correct type of gene identifiers.

To scale the single-sample pathway perturbation scores ( $S_{hip}$ ) so they are comparable across pathways and to test for significance of individual scores, null distributions of perturbation scores for each pathway are generated through a sample permutation strategy, retaining the correct correlation structure between genes within a pathway. With each permutation, column names (i.e. sample labels) for the logCPM matrix are randomly shuffled while the rest of the scoring algorithm remains unchanged. We recommend users to perform a minimum of 1000 permutations, requiring at least 8 unique samples. Subsequently, the median and median absolute deviation (MAD) of the permuted perturbation scores will be calculated and used to normalise the raw perturbation scores to robust Z-scores and obtain associated two-sided  $p$ -values. Since the method is single sample-based, the permutation strategy remains applicable regardless of experimental design.

Apart from assessing whether a pathway's activity changed significantly within an individual sample, users may also be interested in detecting changes at the group-level, which can be performed by modelling scores with regression models, incorporating Smyth's moderated  $t$ -statistics[18] as implemented in `limma`[19]. The single-sample nature of `sSNAPPY`'s pathway perturbation scores is particularly helpful for datasets with complex experimental designs or known confounding factors as these can also be incorporated into the final regression models.

## Operation

The package has been tested on all operating systems, requiring R > 4.2.0, and can be installed using BiocManager as follows.

```
if (!requireNamespace("BiocManager", quietly = TRUE))
  install.packages("BiocManager")
BiocManager::install("sSNAPPY")
```

## Use Cases

### Data

We used A publicly available scRNA-seq data of patient-derived ovarian tissue collected prior to and after 11 homogeneously treated high-grade serous ovarian cancer (HGSOC) patients were subjected to chemotherapy[20] to demonstrate the use of `sSNAPPY` here. Pre-processed count data were retrieved from Gene Expression Omnibus (GEO) with accession code GSE165897.

In the original study, cells were classified into epithelial, stromal, and immune cells but we have chosen to only focus on epithelial cells as they were what the original study primarily focused on. Since `sSNAPPY` was designed for bulk RNA-seq data, counts of epithelial cells from the same samples were first summed into pseudo-bulk profiles, giving rise to a total of 22 samples. We considered a gene detectable if we observed

>1.5 counts per million in >11 samples out of 22, representing all samples from a complete treatment group. A total of 11,101 (33.8%) out of 32,847 annotated genes passed the selection criteria and were included in downstream analyses. Conditional quantile normalisation[21] was applied to mitigate potential biases introduced by gene length and GC content. The logCPM matrix of the processed dataset and sample metadata can be downloaded from here.

To start, firstly load all the packages that will be used in this workflow:

```
library(sSNAPPY)
library(tidyverse)
library(magrittr)
library(ggplot2)
library(cowplot)
library(kableExtra)
library(AnnotationHub)
library(edgeR)
```

To read in the data:

```
logCPM <- readRDS(here::here("data/logCPM.rds"))
sample_meta <- readRDS(here::here("data/sample_meta.rds"))
head(sample_meta)
```

```
## # A tibble: 6 x 8
##   sample          treatment    patie~1 anato~2    Age Stage    PFI    CRS
##   <chr>          <chr>        <chr>   <chr>    <dbl> <chr> <dbl> <dbl>
## 1 EOC372_treatment-naive treatment-naive EOC372 Perito~    68 IIIC    460    1
## 2 EOC372_post-NACT      post-NACT      EOC372 Perito~    68 IIIC    460    1
## 3 EOC443_post-NACT      post-NACT      EOC443 Omentum    54 IVA     177    3
## 4 EOC443_treatment-naive treatment-naive EOC443 Omentum    54 IVA     177    3
## 5 EOC540_treatment-naive treatment-naive EOC540 Omentum    62 IIIC    126    2
## 6 EOC540_post-NACT      post-NACT      EOC540 Omentum    62 IIIC    126    2
## # ... with abbreviated variable names 1: patient_id, 2: anatomical_location
```

### Data preparation and retrieval pathway topology

To apply sSNAPPY, the rownames of the logCPM matrix must be converted to Entrez IDs. Genes without an Entrez IDs were removed.

```
ah <- AnnotationHub() %>%
  AnnotationHub::subset(rdataclass == "EnsDb") %>%
  AnnotationHub::subset(str_detect(description, "101")) %>%
  AnnotationHub::subset(genome == "GRCh38")
stopifnot(length(ah) == 1)
ensDb <- ah[[1]]
rownames(logCPM) <- mapIds(ensDb, rownames(logCPM),
  "ENTREZID", keytype = "GENENAME")
```

```
## Warning: Unable to map 210 of 10311 requested IDs.
```

```
# Remove genes that couldn't be matched to entrez IDs
logCPM <- logCPM[!is.na(rownames(logCPM)),]
```

Next, pathway topology information needs to be retrieved from a chosen database. Using KEGG as an example, the retrieved topology information will be stored as a list where each element corresponds to a pathway and the numbers in the matrices encode gene-gene interaction.

```
gsTopology <- retrieve_topology(database = "kegg")
```

Instead of downloading the topology matrices of all pathways, it is also possible to specify specific pathways' names to focus on. Customised weights could be assigned to different types of gene-gene interaction type by providing a named numeric vector.

```
# Only retrieve the topology matrices of 3 specific pathways
gsTopology_sub <- retrieve_topology(
  database = "kegg",
  pathwayName = c(
    "Glycolysis / Gluconeogenesis",
    "Citrate cycle (TCA cycle)",
    "Pentose phosphate pathway"
  )
)
```

### Score single-sample pathway perturbation

To compute the single-sample logFCs needed for perturbation scores, samples must be in matched pairs and the factor defining those pairs must be specified in the `weight_ss_fc()` function. In our example dataset, pre- and post-treatment samples are matched by patient IDs. Additionally, the sample metadata must include the treatment of all samples, one level of which must be the control level indicated as a parameter of the `weight_ss_fc()` function.

```
weightedFC <- weight_ss_fc(
  logCPM, sample_meta, factor = "patient_id",
  control = "treatment-naive"
)
names(weightedFC)
```

```
## [1] "weight" "logFC"
```

The output of the `weight_ss_fc` function is a list where one element is a matrix of single-sample logFCs with rows corresponding to genes and columns to treated samples and the other element is a vector of gene-wise weights that will be used to alleviate the influence of lowly expressed but highly variable genes.

Single-sample logFCs will then be piped through pathway topologies to compute the gene-wise perturbation scores for all genes included in a pathway. Apart from being summed into pathway-level perturbation scores, gene-wise perturbation scores can also be ranked to identify genes playing the most significant roles in each pathway, which will be further elaborated in the visualisation section below. The pathway-level perturbation scores will be returned as a data.frame containing sample and gene-set names.

```
# Compute perturbation scores at the gene-level within
# each pathway each treated sample
genePertScore <- raw_gene_pert(weightedFC$logFC, gsTopology)
# Sum gene-level scores to derive pathway-level scores
ssPertScore <- pathway_pert(genePertScore)
head(ssPertScore)
```

```
##   sample      tA                                     gs_name
## 1 EOC372 0.005332217 EGFR tyrosine kinase inhibitor resistance
## 2 EOC443 -0.001506460 EGFR tyrosine kinase inhibitor resistance
## 3 EOC540 -0.006199438 EGFR tyrosine kinase inhibitor resistance
## 4 EOC3   -0.004205656 EGFR tyrosine kinase inhibitor resistance
## 5 EOC87  -0.003840786 EGFR tyrosine kinase inhibitor resistance
## 6 EOC136 -0.008354838 EGFR tyrosine kinase inhibitor resistance
```

### Sample permutation for normalisation and significance testing

To estimate the significance of individual scores and transform the scores so they are comparable across pathways, sSNAPPY utilises a sample-permutation strategy to simulate the null distributions of perturbation scores. Since sample labels will be permuted randomly, sample metadata is not required by the `generate_permuted_scores` function, instead, users only need to specify the number of treatment groups in the study, including the control level. Since permutation requires a large amount of computational time and memory, sSNAPPY parallelises this step through functions provided by BiocParallel. Users can choose to customize the parallel back-end or the default one returned by `BiocParallel::bpparam()` will be used. If the number of samples or the size of the chosen pathway database is large, it is recommended to perform the permutation step on a high-performance computer.

```
permutedScore <- generate_permuted_scores(
  logCPM, numOffTreat = 2, NB = 1000,
  gsTopology = gsTopology, weight = weightedFC$weight
)
```

Next, using the function `normalise_by_permu`, the raw perturbation scores will be normalised to robust z-scores using the median and MAD of permuted scores and further converted to two-sided p-values. Defaulted to using the false-discovery rate, the p-values will be corrected for multiple testings using a user-defined approach. In this case, using an FDR of 0.05 as a cut-off, none of the pathways was considered to be significantly perturbed at the individual sample level.

```
normalisedScores <- normalise_by_permu(permutedScore, ssPertScore)
head(normalisedScores)
```

```
##                                gs_name      MAD      MEDIAN  sample
## 1 EGFR tyrosine kinase inhibitor resistance 0.0073873075 -1.837123e-04 EOC1005
## 2                               Endocrine resistance 0.0109646845  1.323193e-04 EOC1005
## 3                               Antifolate resistance 0.0017190861  2.189169e-05 EOC1005
## 4                               Platinum drug resistance 0.0085760397  9.387240e-05 EOC1005
## 5                               mRNA surveillance pathway 0.0001001364 -1.095895e-06 EOC1005
## 6                               RNA degradation 0.0015561033  1.181997e-06 EOC1005
##                                tA      robustZ      pvalue adjPvalue
## 1 -0.0044946212 -0.5835562  0.5595190  0.9996697
## 2  0.0051022014  0.4532627  0.6503596  0.9996697
## 3  0.0011010089  0.6277272  0.5301826  0.9996697
## 4 -0.0065039960 -0.7693374  0.4416930  0.9996697
## 5 -0.0001183863 -1.1713060  0.2414758  0.9996697
## 6 -0.0001092463 -0.0709646  0.9434259  0.9996697
```

A key biological question to answer in this study was what biological processes were impacted by chemotherapy across all patients, the answer to which could be obtained by applying t-tests to normalised scores for each pathway, with the null hypothesis being the mean of normalised perturbation scores equals to 0 for a given pathway. To avoid under-estimation of sample variance, the moderated t-statistics approach proposed by Smyth 2004<sup>14</sup> was adopted here. Being an empirical Bayes strategy, in moderated t-test, estimated sample variances will be adjusted towards the expected variances estimated by pooling information across all genes/pathways. Performances of the moderated t-statistics were evaluated through simulation studies, and the moderated approach demonstrated more advantageous performance in controlling the false discovery rate.

```
# Normalised perturbation scores were converted to a matrix,
# with rows corresponding to pathways and columns to samples
pert_matrix <- normalisedScores %>%
  dplyr::select(robustZ, gs_name, sample) %>%
  pivot_wider(
    names_from = sample,
    values_from = robustZ
  ) %>%
  column_to_rownames("gs_name") %>%
  as.matrix()
# Linear models were fitted for each pathway. No design matrix was
# specified as all samples were replicates of the same treatment group
fit_kg <- lmFit(pert_matrix)
# Moderated t-statistics are calculated using the eBayes function
fit_kg <- eBayes(fit_kg, trend = abs(rowMeans(pert_matrix)))
```

Pathways with an FDR < 0.05 in the moderated t-test were considered to be significantly perturbed at the group level. 22 out of the 315 tested KEGG pathways passed the selection threshold (Table 1).

```
# Use topTable to summarise the linear model fit and correct the p-values with FDR
table1 <- topTable(fit_kg,
  number = Inf) %>%
  rownames_to_column("gs_name") %>%
  mutate(
    logFC = round(logFC, 4),
```

```

    gs_name = as.factor(gs_name),
    Direction = ifelse(logFC < 0, "Inhibited", "Activated"),
    Significant = ifelse(adj.P.Val < 0.05, TRUE, FALSE)) %>%
dplyr::select(
  `Gene-set Name` = gs_name, `Change in Perturbation Score` = logFC,
  P.Value, FDR = adj.P.Val, Direction, Significant
)

```

[Table 1 about here.]

In the original study[20], to study the effect of chemotherapy, unsupervised clustering was performed on all cells labelled to be cancer cells. Identified cancer clusters were labelled by performing pathway enrichment testing on cluster marker genes. Two clusters, associated with proliferative DNA repair signatures and stress-related markers, were found to contain significantly higher numbers of post-chemotherapy cells than pre-treatment ones (Ta in Zhang et al. [20]). Marker genes reported for those two clusters were observed in pathways that were detected to be significantly perturbed by sSNAPPY (Table 2).

[Table 2 about here.]

Apart from treating all treated samples as biological replicates, users might also wish to subset samples into groups by phenotypic traits known to affect patients' responses to chemotherapy, such as the stages of their cancers. To do that through the moderated t-statistic strategy, we simply need to provide a design matrix in the lmFit step.

```

X <- model.matrix(
  ~0 + Stage,
  data = sample_meta %>%
    dplyr::filter(treatment == "post-NACT") %>%
    .[match(colnames(pert_matrix), .$patient_id),] %>%
    mutate(
      Stage = ifelse(Stage == "IVA", "IVA", "IIIC/IVB")
    ) %>%
  set_colnames(str_remove_all(colnames(.), "Stage")) %>%
  .[,colSums(.) != 0]
fit_kg2 <- lmFit(pert_matrix, design = X)

```

### Visualise perturbed pathways as networks

sSNAPPY provides various visualisation functions to assist in the interpretation of results. Since biological pathways are not independent of each other and often contain redundant genes, visualising pathway analysis results as a network is a powerful way to not only intuitively summarise the results but also to facilitate the interpretations of the underlying biology. The `plot_gs_network()` function allows users to easily convert a list of significantly perturbed biological pathways to a network where edges between pathway nodes represent overlapping genes. Defined by the `colorBy` parameter, pathway nodes can be coloured by either the pathways' predicted direction of changes or the significance levels (Figure 1). The returned plot is a `ggplot2` [22] object, meaning that components of its theme could be customized just as any other `ggplots` using the `ggplot2::theme` function.

```

# Extract significantly perturbed pathway. To colour the nodes by directions
# of changes, the column name of the average perturbation scores must be robustZ.
sigPathway <- topTable(fit_kg, number = Inf) %>%
  dplyr::filter(adj.P.Val < 0.05) %>%
  rownames_to_column("gs_name") %>%
  dplyr::rename(
    robustZ = AveExpr,
    pvalue = P.Value
  )
# Plot the network structure
p1 <- sSNAPPY::plot_gs_network(
  normalisedScores = sigPathway,
  gsTopology = gsTopology,
  colorBy = "robustZ"
) +

```

```

    theme(
      panel.border = element_blank(),
      panel.background = element_blank()
    )
p2 <- sSNAPPY::plot_gs_network(
  normalisedScores = sigPathway,
  gsTopology = gsTopology,
  # or color nodes by significance levels
  colorBy = "pvalue"
) +
  theme(
    panel.border = element_blank(),
    panel.background = element_blank()
  )
plot_grid(
  p1, p2,
  labels = c("A", "B")
)

```

[Figure 1 about here.]

By examining the network structure, we can see that many of the highly connected pathways playing a central role in the network are immune-related. To confirm it and further condense the information, we can summarise the network structures into key biological groups by performing community detection. Widely used in network analysis, community detection is a technique used to identify groups of nodes that are more densely connected than to any other nodes in the network[23]. sSNAPPY's `plot_community()` function is a one-stop shop for applying a community detection algorithm of the user's choice to the network structure and annotating identified communities by the most common pathway category, denoting the main biological processes perturbed in that community. Retrieved directly from the KEGG website, we have curated the most recent categorisations of KEGG pathways and included it as part of the sSNAPPY package. Annotation of KEGG pathway communities will be automatically completed by calling the in-built database, while analyses involving other pathway databases require user-provided pathway categorisations.

The defaulted Louvain method was applied to the network of significantly perturbed biological pathways and revealed two community structures, where one was annotated to be endocrine system related and the other one was infectious disease-related (Figure 2).

```

sSNAPPY::plot_community(
  normalisedScores = sigPathway,
  gsTopology = gsTopology,
  colorBy = "robustZ",
) +
  theme(
    panel.border = element_blank(),
    panel.background = element_blank()
  )

```

[Figure 2 about here.]

Inferred directly from the expression matrix, a key advantage of sSNAPPY is that it does not require the prior definition of differentially expressed genes, which is not always detectable in clinical datasets. However, knowing the genes that are implicated in perturbing biological pathways, particularly those that affect multiple gene-sets, can provide valuable insights for biologists seeking to formulate hypotheses about underlying biological mechanisms. Therefore, sSNAPPY presents another visualisation feature called `plot_gs2gene`, which enables the inclusion of pathway genes in network structures. Although defaulted to colour all gene nodes in grey, users can provide a vector of logFCs to colour the genes by their changes in expression. As pathways are often made of hundreds of genes, it is recommended not to plot all genes included in perturbed pathways but filter for genes more likely to be playing a significant role, achieved through only providing logFCs of genes to plot. In this example dataset, we chose to only include genes with the top 500 magnitudes of mean logFCs (Figure 3).

```

# Calculate gene-wise mean logFCs
meanFC <- apply(weightedFC$logFC, 1, mean )
# Extract the top 500 meanFCs

```



```
top500_FC <- meanFC %>%
  abs() %>%
  sort(decreasing = TRUE, ) %>%
  .[1:500]
```

Since pathway topologies were retrieved in Entrez IDs, by default, genes' Entrez IDs will be used to annotate gene nodes in the plot. However, users can provide a data.frame mapping Entrez IDs to their chosen identifiers through the mapRownameTo parameter. A data.frame converting Entrez IDs to ensemble gene names has been made available as part of the package.

```
# Read in built-in data.frame entrez2name that matches genes'
# Entrez IDs to gene names
load(system.file("extdata", "entrez2name.rda", package = "sSNAPPY"))
head(entrez2name)
```

```
## # A tibble: 6 x 3
##   gene_id      mapTo   entrezid
##   <chr>      <chr>   <chr>
## 1 ENSG00000223972 DDX11L1 ENTREZID:84771
## 2 ENSG00000223972 DDX11L1 ENTREZID:727856
## 3 ENSG00000223972 DDX11L1 ENTREZID:100287102
## 4 ENSG00000223972 DDX11L1 ENTREZID:100287596
## 5 ENSG00000223972 DDX11L1 ENTREZID:102725121
## 6 ENSG00000227232 WASH7P  ENTREZID:653635
```

```
# Plot the pathway-gene network for genes with top 500
#magnitudes of mean FCs and label gene nodes by gene names
sSNAPPY::plot_gs2gene(
  normalisedScores = sigPathway,
  gsTopology = gsTopology,
  colorGS_By = "robustZ",
  mapEntrezID = entrez2name,
  geneFC = top500_FC,
  edgeAlpha = 0.3,
  GsName_size = 4
) +
  theme(
    panel.border = element_blank(),
    panel.background = element_blank()
  )
```

[Figure 3 about here.]

### Identify hub genes contributing to pathway perturbation

If we would like to further investigate a specific pathway and elucidate the key genes that contributed to its perturbation, such as the activation of the “p53 signalling pathway,” we can employ a heatmap to display the gene-level perturbation scores of all the genes within the pathway and annotate each column (ie. each sample) by the direction of pathway perturbation in that sample or any other sample metadata using the plot\_gene\_contribution function (Figure 4).

```
plot_gene_contribution(
  genePertScore = genePertScore,
  gsToPlot = "p53 signaling pathway",
  metadata = dplyr::filter(sample_meta, treatment == "post-NACT") %>%
  dplyr::select(sample = patient_id, Stage),
  annotation_attribute = c("pathwayPertScore", "Stage"),
  pathwayPertScore = ssPertScore,
  mapEntrezID = entrez2name
)
```

[Figure 4 about here.]

From this heatmap we can easily identify that the hub genes making the biggest contribution to the activation of p53 signalling pathway upon chemotherapy were gene ART and gene ATM. The Ataxia-telangiectasia mutated (ATM) gene is a well-established oncosuppressor[24], mutation of which has been observed in many types of cancers[25]. Also involved in DNA damage repair, the Ataxia telangiectasia and RAD3-related protein kinase (ATR) gene has been shown to be a promising therapeutic target for HGSOC[26].

## Discussion

In conclusion, the paper showcased an R/Bioconductor package that offers a novel single-sample pathway perturbation testing approach. sSNAPPY utilizes pathway topology information to compute perturbation scores that predict pathways' potential directions of changes in individual samples. This approach addresses the limitations of current strategies that fail to account for gene-gene interactions encoded by pathway topologies or predict the directionality of pathway activities. By applying sSNAPPY to a public scRNA-seq data collected before and after HGSOC patients were subjected to chemotherapy, we demonstrated its ability to detect significant pathway perturbations of various interesting biological processes beyond what were shown in the original studies. Overall, sSNAPPY presents a promising strategy for single sample-based pathway analysis in RNA-seq data.

## Data availability

The dataset analysed in this manuscript are stored in the data directory of this GitHub repository.

## Software availability

- Software available from: <https://bioconductor.org/packages/release/bioc/html/sSNAPPY.html>
- Source code available from: <https://github.com/Wenjun-Liu/sSNAPPY>
- Archived source code at time of publication: [DOI (found on right hand side of a Zenodo record)]
- License: MIT

## Competing interests

No competing interests were disclosed

## Grant information

Any grants that supported the work must be listed here, including the grant number.

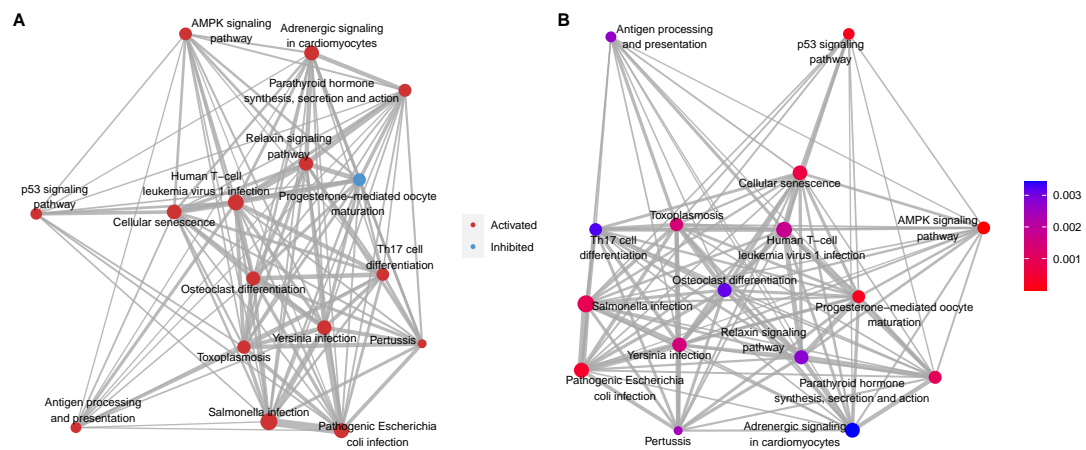
## Acknowledgements

- [1] Aravind Subramanian, Pablo Tamayo, Vamsi K Mootha, Sayan Mukherjee, Benjamin L Ebert, Michael A Gillette, Amanda Paulovich, Scott L Pomeroy, Todd R Golub, Eric S Lander, and Jill P Mesirov. Gene set enrichment analysis: a knowledge-based approach for interpreting genome-wide expression profiles. *Proc. Natl. Acad. Sci. U. S. A.*, 102(43):15545–15550, October 2005.
- [2] Matthew D Young, Matthew J Wakefield, Gordon K Smyth, and Alicia Oshlack. Gene ontology analysis for RNA-seq: accounting for selection bias. *Genome Biol.*, 11(2):R14, February 2010.
- [3] Farhad Maleki, Katie Ovens, Daniel J Hogan, and Anthony J Kusalik. Gene set analysis: Challenges, opportunities, and future research. *Front. Genet.*, 11:654, June 2020.
- [4] Sarah Mubeen, Alpha Tom Kodamullil, Martin Hofmann-Apitius, and Daniel Domingo-Fernández. On the influence of several factors on pathway enrichment analysis. *Brief. Bioinform.*, 23(3), May 2022.
- [5] Hiroyuki Ogata, Susumu Goto, Kazushige Sato, Wataru Fujibuchi, Hidemasa Bono, and Minoru Kanehisa. KEGG: Kyoto Encyclopedia of Genes and Genomes.
- [6] Marvin Martens, Ammar Ammar, Anders Riutta, Andra Waagmeester, Denise N Slenter, Kristina Hanspers, Ryan A. Miller, Daniela Digles, Elisson N Lopes, Friederike Ehrhart, Lauren J Dupuis, Laurent A Winckers, Susan L Coort, Egon L Willighagen, Chris T Evelo, Alexander R Pico, and Martina Kutmon. WikiPathways: connecting communities. *Nucleic Acids Research*, 49(D1):D613–D621, January 2021. ISSN 0305-1048, 1362-4962. doi: 10.1093/nar/gkaa1024.
- [7] Adi Laurentiu Tarca, Sorin Draghici, Purvesh Khatri, Sonia S Hassan, Pooja Mittal, Jung-Sun Kim, Chong Jai Kim, Juan Pedro Kusanovic, and Roberto Romero. A novel signaling pathway impact analysis. *Bioinformatics*, 25(1):75–82, January 2009.

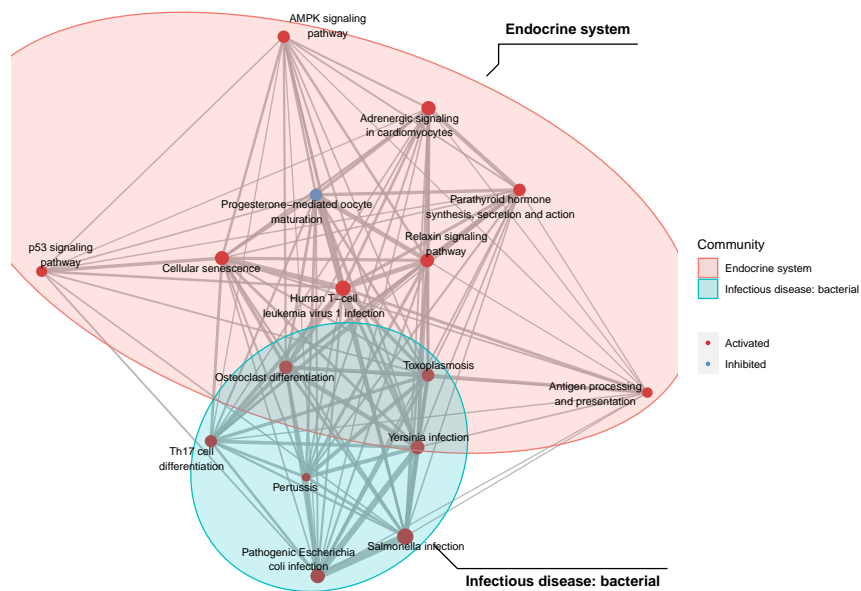
- [8] Laurent Jacob, Pierre Neuvial, and Sandrine Dudoit. More power via graph-structured tests for differential expression of gene networks. *The Annals of Applied Statistics*, 6(2):561 – 600, 2012. doi: 10.1214/11-AOAS528. URL <https://doi.org/10.1214/11-AOAS528>.
- [9] Jing Ma, Ali Shojaie, and George Michailidis. Network-based pathway enrichment analysis with incomplete network information. *Bioinformatics*, 32(20):3165–3174, 06 2016. ISSN 1367-4803. doi: 10.1093/bioinformatics/btw410. URL <https://doi.org/10.1093/bioinformatics/btw410>.
- [10] Maysson Al-Haj Ibrahim, Sabah Jassim, Michael Anthony Cawthorne, and Kenneth Langlands. A topology-based score for pathway enrichment. *Journal of Computational Biology*, 19(5):563–573, 2012. doi: 10.1089/cmb.2011.0182. URL <https://doi.org/10.1089/cmb.2011.0182>. PMID: 22468678.
- [11] Di Wu and Gordon K. Smyth. Camera: a competitive gene set test accounting for inter-gene correlation. *Nucleic Acids Research*, 40(17):e133–e133, 05 2012. ISSN 0305-1048. doi: 10.1093/nar/gks461. URL <https://doi.org/10.1093/nar/gks461>.
- [12] Tuan-Minh Nguyen, Adib Shafi, Tin Nguyen, and Sorin Draghici. Identifying significantly impacted pathways: a comprehensive review and assessment. *Genome Biol.*, 20(1):203, October 2019.
- [13] Jing Ma, Ali Shojaie, and George Michailidis. A comparative study of topology-based pathway enrichment analysis methods. *BMC Bioinformatics*, 20(1):546, December 2019. ISSN 1471-2105. doi: 10.1186/s12859-019-3146-1.
- [14] Sonja Hännelmann, Robert Castelo, and Justin Guinney. GSVA: gene set variation analysis for microarray and RNA-Seq data. *BMC Bioinformatics*, 14(1):7, December 2013. ISSN 1471-2105. doi: 10.1186/1471-2105-14-7.
- [15] Charity W Law, Yunshun Chen, Wei Shi, and Gordon K Smyth. voom: precision weights unlock linear model analysis tools for RNA-seq read counts. *Genome Biol*, 15(2):R29, 2014. ISSN 1465-6906. doi: 10.1186/gb-2014-15-2-r29.
- [16] Adi Laurentiu Tarca, Sorin Draghici, Purvesh Khatri, Sonia S. Hassan, Pooja Mittal, Jung-sun Kim, Chong Jai Kim, Juan Pedro Kusanovic, and Roberto Romero. A novel signaling pathway impact analysis. *Bioinformatics*, 25(1):75–82, January 2009. ISSN 1367-4811, 1367-4803. doi: 10.1093/bioinformatics/btn577.
- [17] Gabriele Sales, Enrica Calura, Duccio Cavalieri, and Chiara Romualdi. graphite - a Bioconductor package to convert pathway topology to gene network. *BMC Bioinformatics*, 13(1):20, December 2012. ISSN 1471-2105. doi: 10.1186/1471-2105-13-20.
- [18] Gordon K Smyth. Linear models and empirical bayes methods for assessing differential expression in microarray experiments. *Statistical Applications in Genetics and Molecular Biology*, 3(1), 2004. doi: doi:10.2202/1544-6115.1027. URL <https://doi.org/10.2202/1544-6115.1027>.
- [19] Matthew E Ritchie, Belinda Phipson, Di Wu, Yifang Hu, Charity W Law, Wei Shi, and Gordon K Smyth. limma powers differential expression analyses for RNA-sequencing and microarray studies. *Nucleic Acids Research*, 43(7):e47, 2015. doi: 10.1093/nar/gkv007.
- [20] Kaiyang Zhang, Erdogan Pekcan Erkan, Sanaz Jamalzadeh, Jun Dai, Noora Andersson, Katja Kaipio, Tarja Lamminen, Naziha Mansuri, Kaisa Huhtinen, Olli Carpén, Sakari Hietanen, Jaana Oikonen, Johanna Hynninen, Anni Virtanen, Antti Häkkinen, Sampsa Hautaniemi, and Anna Vähärautio. Longitudinal single-cell RNA-seq analysis reveals stress-promoted chemoresistance in metastatic ovarian cancer. *Sci. Adv.*, 8(8):eabm1831, February 2022. ISSN 2375-2548. doi: 10.1126/sciadv.abm1831.
- [21] K. D. Hansen, R. A. Irizarry, and Z. Wu. Removing technical variability in RNA-seq data using conditional quantile normalization. *Biostatistics*, 13(2):204–216, April 2012. ISSN 1465-4644, 1468-4357. doi: 10.1093/biostatistics/kxr054.
- [22] Hadley Wickham. *ggplot2: Elegant Graphics for Data Analysis*. Springer New York, New York, NY, 2009. ISBN 978-0-387-98140-6. doi: 10.1007/978-0-387-98141-3. URL <https://link.springer.com/10.1007/978-0-387-98141-3>.
- [23] M. E. J. Newman and M. Girvan. Finding and evaluating community structure in networks. *Phys. Rev. E*, 69(2):026113, February 2004. ISSN 1539-3755, 1550-2376. doi: 10.1103/PhysRevE.69.026113.
- [24] Masoumeh Moslemi, Yousef Moradi, Hojat Dehghanbanadaki, Hamed Afkhami, Mansoor Khaledi, Najmeh Sedighimehr, Javad Fathi, and Ehsan Sohrabi. The association between ATM variants and risk of breast cancer: a systematic review and meta-analysis. *BMC Cancer*, 21(1):27, December 2021. ISSN 1471-2407. doi: 10.1186/s12885-020-07749-6.
- [25] Michael Choi, Thomas Kipps, and Razelle Kurzrock. ATM Mutations in Cancer: Therapeutic Implications. *Molecular Cancer Therapeutics*, 15(8):1781–1791, August 2016. ISSN 1535-7163, 1538-8514. doi: 10.1158/1535-7163.MCT-15-0945.
- [26] Siyu Li, Tao Wang, Xichang Fei, and Mingjun Zhang. ATR Inhibitors in Platinum-Resistant Ovarian Cancer. *Cancers*, 14(23):5902, November 2022. ISSN 2072-6694. doi: 10.3390/cancers14235902.

## List of Figures

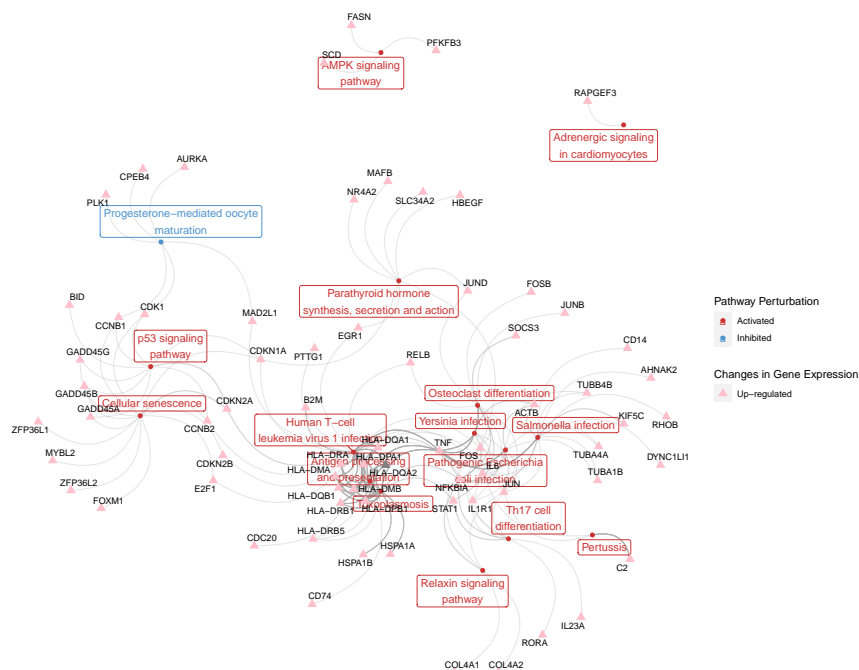
- 1 Significantly perturbed KEGG pathways identified among post-chemotherapy samples using sS-NAPPY, colored by (A) pathways' predicted directions of changes and (B) pathways'  $-\log_{10}(p\text{-values})$ . Pathways with a FDR < 0.05 in the moderated t-test were included. . . . . 13
- 2 Significantly perturbed KEGG pathways identified among post-chemotherapy samples using sS-NAPPY, colored by (A) pathways' predicted directions of changes and (B) pathways'  $-\log_{10}(p\text{-values})$ . Pathways with a FDR < 0.05 in the moderated t-test were included. . . . . 14
- 3 Significantly perturbed KEGG pathways identified among post-chemotherapy samples using sS-NAPPY, annotated by community structure identified with the Louvain algorithm. 2 communities were formed, both of which were annotated by the pathway category that the majority of the pathways belong to. . . . . 15
- 4 Gene-level perturbation scores of all genes included in the "p53 signalling pathway" pathway in each sample where columns of samples were annotated by pathway-level perturbation and the stages of cancers. Genes ATR and ATM were the key driver of the activation of p53 signalling pathway. . . . . 16



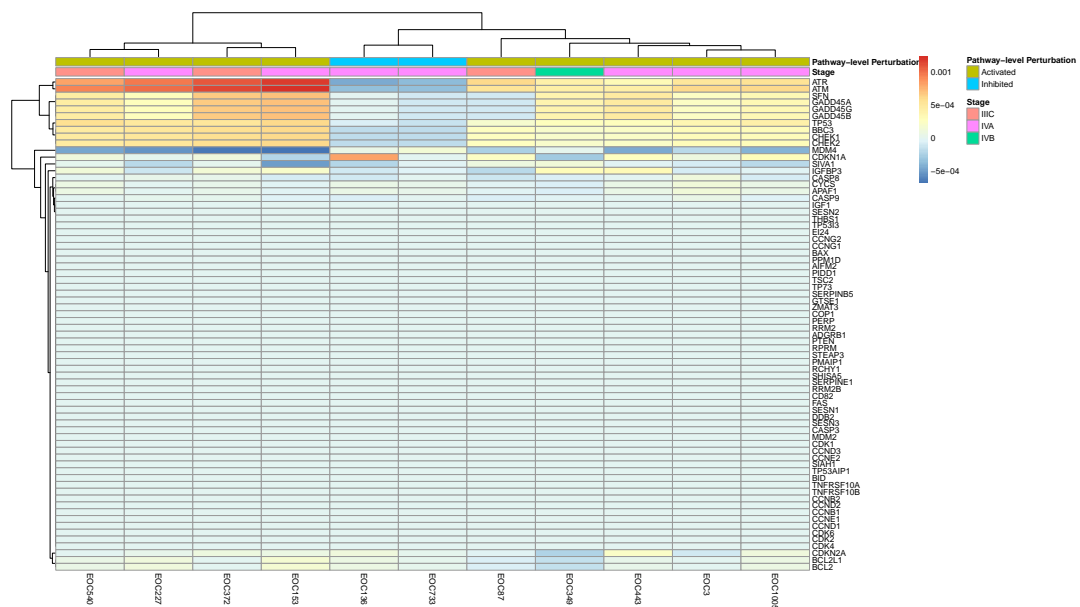
**Figure 1.** Significantly perturbed KEGG pathways identified among post-chemotherapy samples using sSNAPPY, colored by (A) pathways' predicted directions of changes and (B) pathways'  $-\log_{10}(p\text{-values})$ . Pathways with a FDR < 0.05 in the moderated t-test were included.



**Figure 2.** Significantly perturbed KEGG pathways identified among post-chemotherapy samples using sSNAPPY, colored by (A) pathways' predicted directions of changes and (B) pathways'  $-\log_{10}(\text{p-values})$ . Pathways with a  $\text{FDR} < 0.05$  in the moderated t-test were included.



**Figure 3.** Significantly perturbed KEGG pathways identified among post-chemotherapy samples using sSNAPPY, annotated by community structure identified with the Louvain algorithm. 2 communities were formed, both of which were annotated by the pathway category that the majority of the pathways belong to.



**Figure 4.** Gene-level perturbation scores of all genes included in the "p53 signalling pathway" pathway in each sample where columns of samples were annotated by pathway-level perturbation and the stages of cancers. Genes ATR and ATM were the key driver of the activation of p53 signalling pathway.



## List of Tables

1	Results of KEGG pathways identified among post-chemotherapy samples using sSNAPPY . . . .	18
2	. Detection of stress-associated and proliferative DNA repair signatures, as reported by Zhang et al., 2022, in significantly perturbed pathways identified by sSNAPPY. . . . .	19

**Table 1. Results of KEGG pathways identified among post-chemotherapy samples using sSNAPPY**

Gene-set Name	Change in Perturbation Score	PValue	FDR	Direction	Significant
AMPK signaling pathway	0.8339	0.0000053	0.0011315	Activated	TRUE
p53 signaling pathway	0.8218	0.0001878	0.0183811	Activated	TRUE
Progesterone-mediated oocyte maturation	-0.7096	0.0002577	0.0183811	Inhibited	TRUE
Pathogenic Escherichia coli infection	0.6356	0.0003715	0.0198775	Activated	TRUE
Cellular senescence	0.6659	0.0006706	0.0287013	Activated	TRUE
Salmonella infection	0.6165	0.0008808	0.0302502	Activated	TRUE
Parathyroid hormone synthesis, secretion and action	0.6150	0.0009895	0.0302502	Activated	TRUE
Toxoplasmosis	0.6350	0.0014170	0.0357963	Activated	TRUE
Yersinia infection	0.5952	0.0015055	0.0357963	Activated	TRUE
Human T-cell leukemia virus 1 infection	0.5954	0.0018025	0.0385739	Activated	TRUE
Pertussis	0.5647	0.0024215	0.0458335	Activated	TRUE
Antigen processing and presentation	0.5594	0.0026594	0.0458335	Activated	TRUE
Relaxin signaling pathway	0.5105	0.0027843	0.0458335	Activated	TRUE
Osteoclast differentiation	0.5906	0.0031175	0.0459781	Activated	TRUE
Th17 cell differentiation	0.6178	0.0032525	0.0459781	Activated	TRUE
Adrenergic signaling in cardiomyocytes	0.6312	0.0034376	0.0459781	Activated	TRUE
Rheumatoid arthritis	0.5555	0.0044125	0.0509248	Activated	FALSE
Human immunodeficiency virus 1 infection	0.5380	0.0046017	0.0509248	Activated	FALSE
Leukocyte transendothelial migration	0.5401	0.0047306	0.0509248	Activated	FALSE
B cell receptor signaling pathway	0.5605	0.0047593	0.0509248	Activated	FALSE
Fluid shear stress and atherosclerosis	0.5416	0.0050961	0.0515160	Activated	FALSE
Chagas disease	0.4553	0.0052960	0.0515160	Activated	FALSE
Amphetamine addiction	0.4969	0.0066948	0.0600218	Activated	FALSE
Parkinson disease	0.6322	0.0067314	0.0600218	Activated	FALSE
Chemical carcinogenesis - reactive oxygen species	0.4914	0.0076541	0.0636785	Activated	FALSE
Prion disease	0.5744	0.0077366	0.0636785	Activated	FALSE
Transcriptional misregulation in cancer	0.5247	0.0083880	0.0664828	Activated	FALSE
Amyotrophic lateral sclerosis	0.5598	0.0095735	0.0722516	Activated	FALSE
Autophagy - other	0.5089	0.0097911	0.0722516	Activated	FALSE
Endocrine resistance	0.4517	0.0103565	0.0738765	Activated	FALSE
cAMP signaling pathway	0.4566	0.0111216	0.0767749	Activated	FALSE
Estrogen signaling pathway	0.5307	0.0115901	0.0775091	Activated	FALSE
Leishmaniasis	0.5675	0.0130130	0.0843876	Activated	FALSE
C-type lectin receptor signaling pathway	0.6259	0.0145868	0.0891971	Activated	FALSE
Growth hormone synthesis, secretion and action	0.4164	0.0145883	0.0891971	Activated	FALSE
Chemokine signaling pathway	0.5603	0.0154961	0.0921156	Activated	FALSE
Aldosterone-regulated sodium reabsorption	0.4137	0.0192526	0.1113530	Activated	FALSE
T cell receptor signaling pathway	0.4436	0.0205149	0.1155315	Activated	FALSE
Mitophagy - animal	0.4634	0.0224493	0.1231833	Activated	FALSE
Non-alcoholic fatty liver disease	0.4035	0.0270647	0.1447960	Activated	FALSE
Longevity regulating pathway	0.4070	0.0283340	0.1478898	Activated	FALSE
Shigellosis	0.4729	0.0309009	0.1574472	Activated	FALSE
IL-17 signaling pathway	0.4055	0.0352484	0.1754221	Activated	FALSE
Axon guidance	-0.3872	0.0403841	0.1947390	Inhibited	FALSE
EGFR tyrosine kinase inhibitor resistance	-0.3843	0.0409498	0.1947390	Inhibited	FALSE
Olfactory transduction	0.4286	0.0477795	0.2222785	Activated	FALSE
FoxO signaling pathway	-0.4077	0.0508525	0.2315413	Inhibited	FALSE
Colorectal cancer	-0.3727	0.0531799	0.2370936	Inhibited	FALSE
Fc epsilon RI signaling pathway	-0.3842	0.0558730	0.2440167	Inhibited	FALSE
Bladder cancer	-0.3708	0.0648394	0.2775128	Inhibited	FALSE
Spinocerebellar ataxia	0.3255	0.0674275	0.2829310	Activated	FALSE
Viral protein interaction with cytokine and cytokine receptor	0.4174	0.0705270	0.2902456	Activated	FALSE
Signaling pathways regulating pluripotency of stem cells	0.3469	0.0733669	0.2940451	Activated	FALSE
Endometrial cancer	-0.3182	0.0753600	0.2940451	Inhibited	FALSE
Chronic myeloid leukemia	-0.3675	0.0774357	0.2940451	Inhibited	FALSE
Coronavirus disease - COVID-19	0.3324	0.0776344	0.2940451	Activated	FALSE
Cushing syndrome	0.3756	0.0789040	0.2940451	Activated	FALSE
Necroptosis	-0.3850	0.0799111	0.2940451	Inhibited	FALSE
Rap1 signaling pathway	0.3348	0.0810685	0.2940451	Activated	FALSE
Fanconi anemia pathway	0.3371	0.0891791	0.3145572	Activated	FALSE
Non-small cell lung cancer	-0.3177	0.0896635	0.3145572	Inhibited	FALSE
Chemical carcinogenesis - receptor activation	0.2794	0.0923548	0.3187729	Activated	FALSE
Vasopressin-regulated water reabsorption	0.3015	0.0952416	0.3235191	Activated	FALSE
Ferroptosis	0.3148	0.0974137	0.3255310	Activated	FALSE
Apoptosis	0.3510	0.0988762	0.3255310	Activated	FALSE
Insulin signaling pathway	-0.3454	0.1005605	0.3260598	Inhibited	FALSE
Gap junction	-0.2715	0.1052140	0.3360568	Inhibited	FALSE
ErbB signaling pathway	-0.3217	0.1070099	0.3367664	Inhibited	FALSE
Prolactin signaling pathway	0.3137	0.1090413	0.3381862	Activated	FALSE
Natural killer cell mediated cytotoxicity	-0.3447	0.1143208	0.3396919	Inhibited	FALSE
Longevity regulating pathway - multiple species	0.2782	0.1152363	0.3396919	Activated	FALSE
Oocyte meiosis	-0.2866	0.1168138	0.3396919	Inhibited	FALSE
Systemic lupus erythematosus	0.2844	0.1169634	0.3396919	Activated	FALSE
Ras signaling pathway	0.3548	0.1174636	0.3396919	Activated	FALSE
Legionellosis	0.3318	0.1206566	0.3416993	Activated	FALSE
Regulation of actin cytoskeleton	-0.3046	0.1219023	0.3416993	Inhibited	FALSE
Epstein-Barr virus infection	-0.3856	0.1230569	0.3416993	Inhibited	FALSE
Choline metabolism in cancer	-0.3105	0.1245446	0.3416993	Inhibited	FALSE
Apelin signaling pathway	0.2871	0.1330388	0.3603836	Activated	FALSE
Adherens junction	0.3023	0.1362206	0.3643900	Activated	FALSE
Alcoholism	0.2735	0.1457567	0.3850857	Activated	FALSE
Alcoholic liver disease	0.3097	0.1518312	0.3962424	Activated	FALSE
Central carbon metabolism in cancer	-0.2967	0.1583378	0.4037397	Inhibited	FALSE
Bacterial invasion of epithelial cells	0.2818	0.1584772	0.4037397	Activated	FALSE
Aldosterone synthesis and secretion	0.2538	0.1613727	0.4062796	Activated	FALSE
AGE-RAGE signaling pathway in diabetic complications	0.2664	0.1659032	0.4128290	Activated	FALSE
Hippo signaling pathway - multiple species	-0.2509	0.1702716	0.4188290	Inhibited	FALSE
NF-kappa B signaling pathway	0.3130	0.1759813	0.4279546	Activated	FALSE

**Table 2.** . Detection of stress-associated and proliferative DNA repair signatures, as reported by Zhang et al., 2022, in significantly perturbed pathways identified by ssNAPPY.

[illegible]