

tugHall version 1.1: USER-GUIDE-tugHall

Requirements for tugHall simulation:

R version **3.3** or later

libraries: **stringr**

Note that the program has two different procedures in general: the first is the simulation and the second is the analysis of the simulation results. Please, pay attention that the requirements for these procedures are different. This User-Guide pertains to the **simulation procedure** alone.

Table of Contents

1. [Quick start guide](#)
2. [Structure of directories](#)
3. [Inputs](#)
4. [Outputs](#)
5. [How to run](#)

1. Quick start guide

The simplest way to run tugHall:

- Save the **/tugHall/** directory to the working folder;
- Run **tugHall.R**.

The code has its initial input parameters and input files in the **/Input/** folder. After the simulation the user can see results of the simulation (please, see **User-Guide-Analysis** for details) in the dialogue box, which will save to the **/Output/** and **/Figures/** folders. Note that the analysis procedure requires additional libraries and a higher version of R - 3.6.0.

2. Structure of directories

Root directory:

User-Guide-tugHall.Rmd - user guide for simulation in the Rmd format.

User-Guide-tugHall.html - user guide for simulation in the html format.

User-Guide-Analysis.Rmd - user guide for analysis and report generation in the Rmd format.

User-Guide-Analysis.html - user guide analysis and report generation in the html format.

dir **/tugHall/** - directory that contains the program.

/tugHall/ directory:

tugHall.R - program to run a simulation and define the parameters.

dir **/Code/** - folder with the code and the function library.

dir **/Input/** - folder with the input files.

dir **/Output/** - folder with the output files.

dir **/Figures/** - folder with the plot figures.

/Code/ directory:

CanSim.bib, **pic_lic.jpg** - files necessary files for the user guide.

tugHall_functions.R - file that contains the functions for the simulation / core of program.

Analysis.R - file to analyze the results of a simulation and plot figures.

Functions.R - file with the functions for the analysis of results.

/Input/ directory:

cellinit.txt - file with a list of initial cells with/without destroyed genes.

gene_cds2.txt - file with hallmark variables and weights.

/Output/ directory:

cellout.txt - file with simulation output.

geneout.txt - file with information about hallmark variables and the weights.

log.txt - file with information about all parameters.

Weights.txt - file with information about weights between hallmarks and genes.

Order_of_dysfunction.txt - see **USER-GUIDE-Analysis**.

VAf.txt - see **USER-GUIDE-Analysis**.

/Figures/ directory

In the **/Figures/** directory there are figures in *.jpg format, which appear after the analysis of the simulation results. See **USER-GUIDE-Analysis**.

3. Inputs

Input of hallmark variables and gene weights

The file **tugHall/Input/gene_cds2.txt** defines the hallmark variables and weights (only first 10 lines are presented here):

Table 1. Input file for genes. Example of input file for hallmarks and weights in the file **tugHall/Input/gene_cds2.txt**.

Genes	length CDS	Hallmark	Suppressor or Oncogene	Weights
APC	8532	apoptosis	s	0.1111111
APC	8532	growth	s	0.1666667
APC	8532	invasion	s	0.2857143
KRAS	567	apoptosis	o	0.3333333
KRAS	567	growth	o	0.3333333
KRAS	567	immortalization	o	0.5714286
KRAS	567	angiogenesis	o	0.4285714
KRAS	567	invasion	o	0.1428571
TP53	1182	apoptosis	s	0.1111111
TP53	1182	growth	s	0.1666667

1. **Genes** - name of gene, e.g., TP53, KRAS. The names must be typed carefully. The program detects all the unique gene names.
2. **length CDS** - length of CDS for each gene, e.g., 2724, 10804.
3. **Hallmark** - hallmark name, e.g., "apoptosis". Available names:
 - apoptosis
 - immortalization
 - growth
 - anti-growth
 - angiogenesis
 - invasion

Note that "growth" and "anti-growth" are related to the single hallmark "growth/anti-growth". Note that "invasion" is related to "invasion/metastasis" hallmark.

4. **Suppressor or oncogene**. - Distinction of oncogene/suppressor:
 - o: oncogene
 - s: suppressor
 - ?: unknown (will be randomly assigned)
5. **Weights** - Hallmark weights for genes, e.g., 0.333 and 0.5. For each hallmark, the program checks the summation of all the weights. If it is not equal to 1, then the program normalizes it to reach unity. Note that, if the gene belongs to more than one hallmark type, it must be separated into separate lines.

After that, the program defines all the weights, and all the **unknown weights** are set equal to 0. Program performs normalization so that the sum of all weights should be equal to 1 for each column. The **tugHall/Output/Weights.txt** file saves these final input weights for the simulation. Only the first 10 lines are presented here:

Table 2. Weights for hallmarks. Example of weights for hallmarks and genes from

tugHall/Output/Weights.txt file. Unknown values equal 0.

Genes	Apoptosis, H_a	Angiogenesis, H_b	Growth / Anti-growth, H_d	Immortalization, H_i	Invasion / Metastasis, H_{im}
APC	0.1590909	0.0000000	0.1666667	0.0000000	0.2195122
KRAS	0.4772727	0.4285714	0.3333333	0.5714286	0.1097561
TP53	0.1590909	0.4285714	0.1666667	0.4285714	0.3292683
PIK3CA	0.2045455	0.1428571	0.3333333	0.0000000	0.3414634

1. **Genes** - name of genes.
2. **Apoptosis**, H_a - weights of hallmark "Apoptosis".
3. **Angiogenesis**, H_b - weights of hallmark "Angiogenesis".
4. **Growth / Anti-growth**, H_d - weights of hallmark "Growth / Anti-growth".
5. **Immortalization**, H_i - weights of hallmark "Immortalization".
6. **Invasion / Metastasis**, H_{im} - weights of hallmark "Invasion / Metastasis".

Input the probabilities

The input of the probabilities used in the model is possible in the code for parameter value settings, "tugHall.R":

Probability variable and value	Description
E0 <- 2E-4	Parameter $E0$ in the division probability
F0 <- 1E0	Parameter $F0$ in the division probability
m <- 1E-6	Mutation probability m'
uo <- 0.5	Oncogene mutation probability u_o
us <- 0.5	Suppressor mutation probability u_s
s <- 10	Parameter in the sigmoid function s
k <- 0.1	Environmental death probability k'

Filename input

Also in the code "tugHall.R" user can define names of input and output files, and additional parameters of simulation:

Variables and file names	Description
genefile <- 'gene_cds2.txt'	File with information about weights
cellfile <- 'cellinit.txt'	Initial Cells
geneoutfile <- 'geneout.txt'	Gene Out file with hallmarks
celloutfile <- 'cellout.txt'	Output information of simulation
logoutfile <- 'log.txt'	Log file to save the input information of simulation
censore_n <- 30000	Max cell number where the program forcibly stops
censore_t <- 200	Max time where the program forcibly stops

Input of the initial cells

The initial states of cells are defined in "tugHall/Input/cellinit.txt" file:

Cell ID	List of mutated genes
1	" "
2	"APC"
3	"APC, KRAS"

4	"KRAS"
5	"TP53, KRAS"
...	...
1000	""

1. **Cell ID** - ID of cell, e.g., 1, 324.
2. **List of mutated genes** - list of mutated genes for each cell, e.g. "", "KRAS, APC". The values are comma separated. The double quotes (") indicate a cell without mutations.

4. Outputs

The output data consists of several files after the simulation. The "log.txt" and "geneout.txt" files contain the input information about variables and gene names. "Weights.txt" has information about the weights of genes for hallmarks (Please refer the section "Inputs"). "Cellout.txt" has information about the dynamics of cell evolution and all variables.

"log.txt" file

The file "log.txt" contains information about probabilities and file names. These variables are explained in the "Inputs".

Table 3. log.txt file. Example of log.txt file.

Variable	Value
genefile	Input/gene_cds2.txt
cellfile	Input/cellinit.txt
geneoutfile	Output/geneout.txt
celloutfile	Output/cellout.txt
logoutfile	Output/log.txt
E	5e-04
F	10
m	1e-05
uo	0.5
us	0.5
s	10
k	0.1
censore_n	30000
censore_t	100

"geneout.txt" file

The file "geneout.txt" contains input information about the weights that connect the hallmarks and genes, which are defined by the user. These variables also are explained in the "Inputs".

Table 4. geneout.txt file. Given below is an example of the geneout.txt file.

Gene_name	Hallmark_name	Weight	Suppressor_or_oncogene
APC	apoptosis	0.1590909	s
KRAS	apoptosis	0.4772727	o
TP53	apoptosis	0.1590909	s
PIK3CA	apoptosis	0.2045455	o
KRAS	immortalization	0.5714286	o
TP53	immortalization	0.4285714	s
APC	growth anti-growth	0.1666667	s
KRAS	growth anti-growth	0.3333333	o
TP53	growth anti-growth	0.1666667	s
PIK3CA	growth anti-growth	0.3333333	o

“cellout.txt” file

The file “cellout.txt” contains the results of the simulation and includes the evolution data: all the output data for each cell at each timestep (only the first 10 lines are presented):

Table 5. Output data. Example of output data for all cells. The names of columns are related to the description in the Tables 1,2 and *USER-GUIDE-Analysis's* figures.

Time	AvgOrindx	ID	ParentID.Birthday	c.	d.	i.	im.	a.	k.	E.	N	Nmax.	M	Ha	Him	Hi	Hd	Hb	type	mut_den	PosDriver.APC	PosDriver.KRAS	PosDriver.TP53	PosDriver.PIK3CA	PosPasngr.APC	PosPasngr.KRAS	PosPasngr.TP53	PosPasngr.PIK3CA	Clone.number	P
1	avg	-	-	0	1	1	0	0.0066929	0.1	5e-04	1000	2000	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	-	
1	1	1	0:0	0	1	1	0	0.0066929	0.1	5e-04	1000	2000	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
1	2	2	0:0	0	1	1	0	0.0066929	0.1	5e-04	1000	2000	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
1	3	3	0:0	0	1	1	0	0.0066929	0.1	5e-04	1000	2000	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
1	4	4	0:0	0	1	1	0	0.0066929	0.1	5e-04	1000	2000	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
1	5	5	0:0	0	1	1	0	0.0066929	0.1	5e-04	1000	2000	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
1	6	6	0:0	0	1	1	0	0.0066929	0.1	5e-04	1000	2000	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
1	7	7	0:0	0	1	1	0	0.0066929	0.1	5e-04	1000	2000	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
1	8	8	0:0	0	1	1	0	0.0066929	0.1	5e-04	1000	2000	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
1	9	9	0:0	0	1	1	0	0.0066929	0.1	5e-04	1000	2000	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	

1. **Time** - the time step, e.g., 1, 50.
2. **AvgOrindx** - “avg” or “index”: “avg” is for a line with averaged values across different (index) lines at the same time step; “index” shows the cell's index at the current time step, e.g., avg, 4,7.
3. **ID** - the unique ID of cell, e.g., 1, 50.
4. **ParentID.Birthday** - the first number is the parent ID, the second number is the birthday time step, e.g., 0:0, 45:5.
5. **c** - the counter of cell divisions for the cell.
6. **d** - the probability of division for the cell, e.g., 0.1, 0.8.
7. **i** - the probability of immortalization for the cell, e.g., 0.1, 0.8.
8. **im** - the probability of invasion/metastasis for the cell, e.g., 0.1, 0.8.
9. **a** - the probability of apoptosis for the cell, e.g., 0.1, 0.8.
10. **k** - the probability of death due to the environment, e.g., 0.1, 0.8.
11. **E** - the E coefficient for the function of the division probability, e.g., 10^4, 10^5.
12. **N** - the number of primary tumor cells at this time step, e.g., 134, 5432.
13. **Nmax** - the theoretically maximal number of primary tumor cells, e.g., 10000, 5000.
14. **M** - the number of metastasis cells at this time step, e.g., 16, 15439.
15. **Ha** - the value of the hallmark “Apoptosis” for the cell, e.g., 0.1, 0.4444.
16. **Him** - the value of the hallmark “Invasion / Metastasis” for the cell, e.g., 0.1, 0.4444.
17. **Hi** - the value of the hallmark “Immortalization” for the cell, e.g., 0.1, 0.4444.
18. **Hd** - the value of the hallmark “Growth / Anti-growth” for the cell, e.g., 0.1, 0.4444 .
19. **Hb** - the value of the hallmark “Angiogenesis” for the cell, e.g., 0.1, 0.4444 .
20. **type** - the type of the cell: “0” is primary tumor cell, “1” is the metastatic cell, e.g., 0, 1.
21. **mut_den** - the density of mutations (tumor mutation burden) for the cell, e.g., 0, 0.32.

The columns from 22 to 25 are related to names in the form **PosDriver. gene name**, where **gene name** is related to user defined genes. The number of columns equals the number of the genes. These columns show the position(s) of driver mutation(s) in a gene: the first number is the mutational site on the gene and the second number is the time step of the mutation, e.g., 3493:4, 4531:34.

22. **PosDriver.(Gene_1=“APC”)** - for the first gene.
23. **PosDriver.(Gene_2=“KRAS”)** - for the second gene.
24. **PosDriver.(Gene_...)** - ...
25. **PosDriver.(Gene_last=“PIK3CA”)** - for the last gene.

The columns from 26 to 29 are related to names in the form **PosPassngr. gene name**, where **gene name** is related to user defined genes. The number of columns equals the number of the genes. These columns show the position(s) of **passenger** mutation(s) in a gene: the first number is the mutational site on the gene and the second number is the time step of the mutation, e.g., 8952:43, 531:4.

26. **PosPassngr.(Gene_1=“APC”)** - for the first gene.
27. **PosPassngr.(Gene_2=“KRAS”)** - for the second gene.
28. **PosPassngr.(Gene_...)** - ...
29. **PosPassngr.(Gene_last=“PIK3CA”)** - for the last gene.

30. **Clone.number** - the clone number is calculated from the binary code of driver mutations. If a gene is mutated, then its binary code value is 1, and if not, it is 0. For example, the cells have only 4 genes in simulation, so "Clone.number" can have binary numbers from 0000 to 1111, which is related decimal numbers from 0 to 15, e.g., 15, 4.
 31. **Passengers.Clone.number** - same as for "Clone.number", but for passenger mutations, e.g., 15, 4.
 32. **Mix.Clone.number** - same as for "Clone.number", but for passenger and driver mutations together. In this case, the length of the binary number is two times larger than for the driver case, e.g., 35, 16.
-

5. How to run

In order to make the simulation, please follow the procedure:

1. Copy **/tugHall/** directory into the working directory.
2. CD to the **/tugHall/** directory.
3. Run the **tugHall.R** file, using the command line like

R --vanilla < tugHall.R

or using the line by line procedure in **R Studio**. In this case we have:

- `load library(stringr)` and `source(file = "Code/tugHall_functions.R");`
 - create the Output and Figures directories, if needed;
 - define the simulation parameters;
 - make the input file for initial cells, if needed;
 - run the `model()` function to simulate;
 - run `source("Code/Analysis.R")` in order to analyze the results and plot the figures in the dialogue box (see **User-Guide-Analysis**).
4. To make a report of the simulation, rename and change the **User-Guide-Analysis.RMD** (for more details, please, see "Writing reproducible reports in R" on the github <https://nicercode.github.io/guides/reports/>), which uses the output and input files in **/Input/**, **/Output/**, **/Figures/** directories that contain actual information. For more details, please, see **User-Guide-Analysis**.