

# metaSeq: Meta-analysis of RNA-seq count data

Koki Tsuyuzaki<sup>1</sup>, and Itoshi Nikaido<sup>2</sup>.

September 26, 2013

<sup>1</sup>Department of Medical and Life Science, Tokyo University of Science.

<sup>2</sup>Bioinformatics Research Unit, Advanced Center for Computing and Communication,  
RIKEN.

`k.t.the-answer@hotmail.co.jp`

## Contents

<b>1</b>	<b>Introduction</b>	<b>2</b>
<b>2</b>	<b>RSE: Read-Size Effect</b>	<b>2</b>
<b>3</b>	<b>Robustness against RSE</b>	<b>3</b>
<b>4</b>	<b>Getting started</b>	<b>5</b>
<b>5</b>	<b>Meta-analysis by non-NOISeq method</b>	<b>6</b>
<b>6</b>	<b>Setup</b>	<b>11</b>

# 1 Introduction

This document provides the way to perform meta-analysis of RNA-seq data using *metaSeq* package. Meta-analysis is a attempt to integrate multiple data in different studies and retrieve much reliable and reproducible result. In our package, the probability of one-sided *NOISeq* [1] is applied in each study. This is because the numbers of reads are often different depending on its study and *NOISeq* is robust method against its difference (see the next section).

## 2 RSE: Read-Size Effect

In many cases, the number of reads are depend on study. For example, here we prepared multiple RNA-Seq count data designed as Breast Cancer cell lines vs Normal cells measured in 4 different studies (this data is also accessible by `data(BreastCancer)`).

ID in this vignette	Accession (SRA / ERA Accession)	Experimental Design
StudyA	SRP008746	Breast Cancer (n=3) vs Normal (n=2)
StudyB	SRP006726	Breast Cancer (n=1) vs Normal (n=1)
StudyC	SRP005601	Breast Cancer (n=7) vs Normal (n=1)
StudyD	ERP000992	Breast Cancer (n=2) vs Normal (n=1)

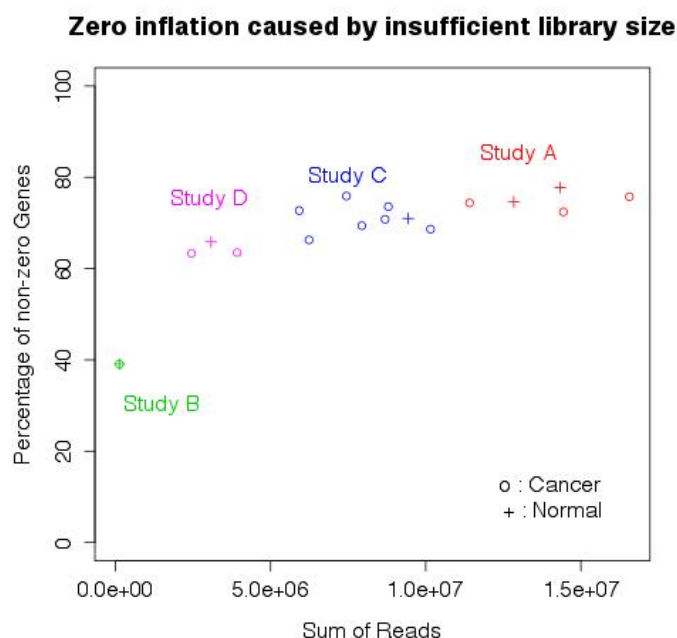


Figure 1: Difference of the number of reads

As shown in the figure 1, the number of reads in StudyA, B, C, and D are relatively different. Generally, statistical test is influenced by the number of reads; the more the number of reads is large, the more the statistical tests are tend to be significant (see the next section). Therefore, in meta-analysis of RNA-seq data, data may be suffered from this bias. Here we call this bias as RSE (Read Size Effect).

### 3 Robustness against RSE

In the point of view of robustness against RSE, we evaluated five widely used method in RNA-seq; *DESeq* [2], *edgeR* [3], *baySeq* [4], and *NOISeq* [1]. Here we used only StudyA data. All counts in the matrix are repeatedly down-sampled in accordance with distributions of binomial (the probability equals 0.5). 1 (original), 1/2, 1/4, 1/8, 1/16, and 1/32-fold data are prepared as low read size situation. In each read size, four methods are conducted (figure 2.A, this data is also accessible by data(StudyA) and data(pvals)), then we focussed on how top500 genes of original data in order of significance will change its members, influenced by low read size (figure 2.B).

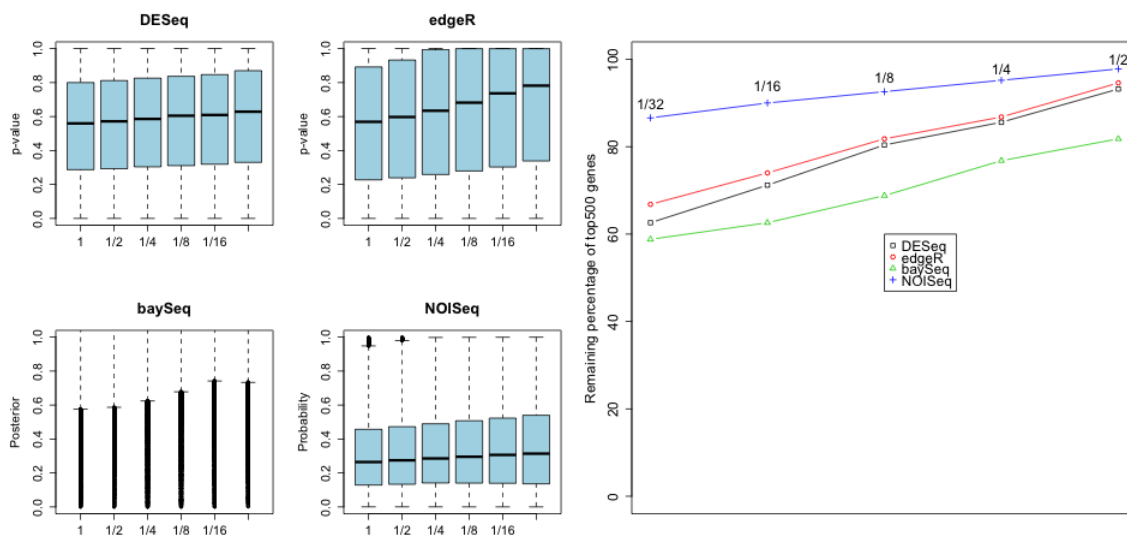


Figure 2: A(left): RSE in each RNA-Seq method, B(right): Top 500 genes in order of significance

Ideal method will returns same result regardless of read size, because same data was used. As shown in figure 2, *NOISeq* is not almost affected by the number of reads and robustly detects same genes as DEGs. Therefore, we concluded that *NOISeq* is suitable method at least in the point of view of meta-analysis. Note that probability of *NOISeq* is not equal to p-value; it is the probability that a gene is differentially expressed [1]. Our package integrates its probability by Fisher's method [5] or Stouffer's method (inverse normal method) [6]. In regard to Stouffer's method, weighting by the number of replicates (sample size) is used.

## 4 Getting started

At first, install and load the *metaSeq* and *snow*.

```
> library("metaSeq")
> library("snow")
```

The RNA-seq expression data in breast cancer cell lines and normal cells is prepared. The data is measured from 4 different studies. The data is stored as a matrix (23368 rows  $\times$  18 columns).

```
> data(BreastCancer)
```

We need to prepare two vectors. First vector is for indicating the experimental condition (e.g., 1: Cancer, 2: Normal) and second one is for indicating the source of data (e.g., A: StudyA, B: StudyB, C: StudyC, D: StudyD).

```
> flag1 <- c(1,1,1,0,0, 1,0, 1,1,1,1,1,1,1,0, 1,1,0)
> flag2 <- c("A","A","A","A","A", "B","B", "C","C","C","C","C","C","C","C", "D","D","D")
```

Then, we use `meta.readData` to create R object for `meta.oneside.noiseq`.

```
> cds <- meta.readData(data = BreastCancer, factor = flag1, studies = flag2)
```

`oneside.noiseq` is performed in each studies and each probabilities are summarized as member of list object.

```
> ## This is very time consuming step.
> # cl <- makeCluster(4, "SOCK")
> # result <- meta.oneside.noiseq(cds, k = 0.5, norm = "tmm", replicates = "biological",
> # factor = flag1, conditions = c(1, 0), studies = flag2, cl = cl)
> # stopCluster(cl)
>
> ## Please load pre-calculated result (Result.Meta)
> ## by data function instead of scripts above.
> data(Result.Meta)
> result <- Result.Meta
```

Fisher's method and Stouffer's method can be applied to the result of `meta.oneside.noiseq`.

```
> F <- Fisher.test(result)
> S <- Stouffer.test(result)
```

These outputs are summarized as list whose length is 3. First member is the probability which means a gene is upper-regulated genes, and Second member is lower-regulated genes. Weight in each study is also saved as its third member (weight is used only by Stouffer's method).

```
> head(F$Upper)
```

```

1/2-SBSRNA4      A1BG      A1BG-AS1      A1CF      A2LD1
  0.3842542    0.5316118    0.5325544      NA    0.1358559
      A2M
  0.2252807

```

```
> head(F$Lower)
```

```

1/2-SBSRNA4      A1BG      A1BG-AS1      A1CF      A2LD1
  0.8420357    0.6078896    0.4047202      NA    0.3661371
      A2M
  0.6197968

```

```
> F$Weight
```

```

Study 1 Study 2 Study 3 Study 4
      5      2      8      3

```

```
> head(S$Upper)
```

```

1/2-SBSRNA4      A1BG      A1BG-AS1      A1CF      A2LD1
  0.3709297    0.2663748    0.2711745      NA    0.2957139
      A2M
  0.2996707

```

```
> head(S$Lower)
```

```

1/2-SBSRNA4      A1BG      A1BG-AS1      A1CF      A2LD1
  0.6290703    0.7336252    0.7288255      NA    0.7042861
      A2M
  0.7003293

```

```
> S$Weight
```

```

Study 1 Study 2 Study 3 Study 4
      5      2      8      3

```

## 5 Meta-analysis by non-NOISeq method

For some reason, we may want to use non-NOISeq method like *DESeq*, *edgeR*, or even *cuffdiff* [7]. We prepared `other.oneside.noiseq` as optional function for such methods. Returned object can be directly applied for `Fisher.test` and `Stouffer.test`.

```

> ## Assume this matrix as one-sided p-values
> ## generated by non-NOISeq method (e.g., cuffdiff)
> upper <- matrix(runif(300), ncol=3, nrow=100)
> lower <- 1 - upper
> rownames(upper) <- paste0("Gene", 1:100)

```

```

> rownames(lower) <- paste0("Gene", 1:100)
> weight <- c(3,6,8)
> ## other.oneside.pvalues function return a matrix
> ## which can input Fisher.test or Stouffer.test
> result <- other.oneside.pvalues(upper, lower, weight)
> ## Fisher's method (without weighting)
> F <- Fisher.test(result)
> str(F)

```

List of 3

```

$ Upper : Named num [1:100] 0.2689 0.7954 0.361 0.6798 0.0563 ...
..- attr(*, "names")= chr [1:100] "Gene1" "Gene2" "Gene3" "Gene4" ...
$ Lower : Named num [1:100] 0.652 0.147 0.606 0.584 0.97 ...
..- attr(*, "names")= chr [1:100] "Gene1" "Gene2" "Gene3" "Gene4" ...
$ Weight: Named num [1:3] 3 6 8
..- attr(*, "names")= chr [1:3] "Exp 1" "Exp 2" "Exp 3"

```

> F

\$Upper

	Gene1	Gene2	Gene3	Gene4	Gene5	Gene6
0.26887544	0.79539709	0.36104357	0.67981904	0.05630875	0.82815743	
	Gene7	Gene8	Gene9	Gene10	Gene11	Gene12
0.79533495	0.55382732	0.21703959	0.46588908	0.15349836	0.27396489	
	Gene13	Gene14	Gene15	Gene16	Gene17	Gene18
0.39098451	0.21673232	0.39577063	0.22993103	0.68103680	0.40779439	
	Gene19	Gene20	Gene21	Gene22	Gene23	Gene24
0.65653452	0.54909973	0.43785146	0.56620293	0.04103333	0.77685804	
	Gene25	Gene26	Gene27	Gene28	Gene29	Gene30
0.40403225	0.86068346	0.32396153	0.36745680	0.73127190	0.75932992	
	Gene31	Gene32	Gene33	Gene34	Gene35	Gene36
0.45392988	0.56336136	0.05003923	0.30386104	0.85535751	0.57161305	
	Gene37	Gene38	Gene39	Gene40	Gene41	Gene42
0.10057188	0.34761241	0.45937419	0.33926033	0.93159012	0.85095826	
	Gene43	Gene44	Gene45	Gene46	Gene47	Gene48
0.21624892	0.29152807	0.83852383	0.83659417	0.45638506	0.12427522	
	Gene49	Gene50	Gene51	Gene52	Gene53	Gene54
0.08333039	0.36126924	0.08134459	0.60074658	0.66907227	0.74178398	
	Gene55	Gene56	Gene57	Gene58	Gene59	Gene60
0.83135209	0.72240983	0.58649083	0.46409292	0.29577362	0.15142499	
	Gene61	Gene62	Gene63	Gene64	Gene65	Gene66
0.29691563	0.03566004	0.27928016	0.10460455	0.33510078	0.28685043	
	Gene67	Gene68	Gene69	Gene70	Gene71	Gene72
0.43248838	0.85950235	0.21733554	0.92370472	0.39227056	0.51664371	
	Gene73	Gene74	Gene75	Gene76	Gene77	Gene78
0.78890398	0.49982276	0.19614476	0.69264709	0.09286959	0.80558715	

Gene79	Gene80	Gene81	Gene82	Gene83	Gene84
0.64543693	0.68563273	0.88759542	0.37649581	0.10647869	0.75005451
Gene85	Gene86	Gene87	Gene88	Gene89	Gene90
0.11184655	0.71920611	0.09242819	0.20858119	0.47627566	0.94984431
Gene91	Gene92	Gene93	Gene94	Gene95	Gene96
0.78124740	0.70633272	0.44185256	0.10880750	0.74041350	0.83573005
Gene97	Gene98	Gene99	Gene100		
0.59406789	0.26432207	0.78431898	0.39466776		

\$Lower

Gene1	Gene2	Gene3	Gene4	Gene5
0.652413882	0.146501514	0.606131358	0.583686392	0.970398775
Gene6	Gene7	Gene8	Gene9	Gene10
0.383833923	0.371901383	0.396517272	0.630097959	0.502001556
Gene11	Gene12	Gene13	Gene14	Gene15
0.657769400	0.596121163	0.655262341	0.940133552	0.596550388
Gene16	Gene17	Gene18	Gene19	Gene20
0.881225316	0.371145361	0.699595754	0.129714674	0.734718137
Gene21	Gene22	Gene23	Gene24	Gene25
0.403311827	0.620632192	0.469715168	0.478700011	0.512975131
Gene26	Gene27	Gene28	Gene29	Gene30
0.111416168	0.878635760	0.570622626	0.178330353	0.448206702
Gene31	Gene32	Gene33	Gene34	Gene35
0.740550638	0.484900496	0.577012268	0.678935220	0.239536749
Gene36	Gene37	Gene38	Gene39	Gene40
0.239644653	0.592008374	0.639015089	0.624452038	0.432982805
Gene41	Gene42	Gene43	Gene44	Gene45
0.005413383	0.338316468	0.352161821	0.852883888	0.245333232
Gene46	Gene47	Gene48	Gene49	Gene50
0.213024645	0.678752560	0.958092437	0.901033505	0.087289165
Gene51	Gene52	Gene53	Gene54	Gene55
0.933160978	0.675950433	0.283400160	0.030181120	0.237060195
Gene56	Gene57	Gene58	Gene59	Gene60
0.440804184	0.631823149	0.367323141	0.389547653	0.924748936
Gene61	Gene62	Gene63	Gene64	Gene65
0.302949943	0.980232102	0.848895645	0.501601400	0.722554083
Gene66	Gene67	Gene68	Gene69	Gene70
0.392372034	0.634429095	0.333286720	0.232123514	0.229197237
Gene71	Gene72	Gene73	Gene74	Gene75
0.823821663	0.305617637	0.381460004	0.719090100	0.901096726
Gene76	Gene77	Gene78	Gene79	Gene80
0.398397364	0.940038105	0.351058446	0.571450551	0.208109906
Gene81	Gene82	Gene83	Gene84	Gene85
0.043055097	0.727524598	0.867068540	0.407251559	0.788133895
Gene86	Gene87	Gene88	Gene89	Gene90



0.421128012	0.974720009	0.550183129	0.729664104	0.145350901
Gene91	Gene92	Gene93	Gene94	Gene95
0.333817298	0.483575022	0.630163430	0.977988047	0.276865051
Gene96	Gene97	Gene98	Gene99	Gene100
0.380463488	0.711097281	0.910599913	0.321765232	0.833695436

\$Weight

Exp 1	Exp 2	Exp 3
3	6	8

```
> ## Stouffer's method (with weighting by sample-size)
> S <- Stouffer.test(result)
> str(S)
```

List of 3

```
$ Upper : Named num [1:100] 0.5157 0.9336 0.5047 0.6159 0.0481 ...
..- attr(*, "names")= chr [1:100] "Gene1" "Gene2" "Gene3" "Gene4" ...
$ Lower : Named num [1:100] 0.4843 0.0664 0.4953 0.3841 0.9519 ...
..- attr(*, "names")= chr [1:100] "Gene1" "Gene2" "Gene3" "Gene4" ...
$ Weight: Named num [1:3] 3 6 8
..- attr(*, "names")= chr [1:3] "Exp 1" "Exp 2" "Exp 3"
```

```
> S
```

\$Upper

Gene1	Gene2	Gene3	Gene4	Gene5	Gene6
0.51566590	0.93360429	0.50465740	0.61594795	0.04809289	0.78095817
Gene7	Gene8	Gene9	Gene10	Gene11	Gene12
0.73879607	0.66526340	0.15247161	0.55364801	0.47305001	0.16876705
Gene13	Gene14	Gene15	Gene16	Gene17	Gene18
0.50900267	0.13788049	0.50147844	0.11058371	0.75395398	0.46938711
Gene19	Gene20	Gene21	Gene22	Gene23	Gene24
0.93011113	0.36975078	0.70782562	0.37963374	0.32689138	0.59348708
Gene25	Gene26	Gene27	Gene28	Gene29	Gene30
0.65730946	0.93395906	0.19188184	0.55999159	0.76472885	0.68237794
Gene31	Gene32	Gene33	Gene34	Gene35	Gene36
0.45396734	0.68144710	0.11795119	0.25030442	0.88362585	0.78519020
Gene37	Gene38	Gene39	Gene40	Gene41	Gene42
0.08082026	0.19222602	0.56405405	0.34578329	0.96950421	0.72584227
Gene43	Gene44	Gene45	Gene46	Gene47	Gene48
0.73480594	0.15282272	0.71718408	0.83583474	0.40574373	0.12455979
Gene49	Gene50	Gene51	Gene52	Gene53	Gene54
0.06932575	0.60336704	0.05248166	0.50250789	0.51059388	0.73325549
Gene55	Gene56	Gene57	Gene58	Gene59	Gene60
0.86519273	0.54103207	0.55896221	0.31096706	0.29559913	0.08256385
Gene61	Gene62	Gene63	Gene64	Gene65	Gene66

0.76731010	0.06039796	0.23577510	0.14977388	0.41886799	0.52027545
Gene67	Gene68	Gene69	Gene70	Gene71	Gene72
0.54837777	0.75699537	0.38404579	0.81633138	0.24351709	0.54614700
Gene73	Gene74	Gene75	Gene76	Gene77	Gene78
0.61973364	0.47941238	0.18035215	0.55852283	0.12056559	0.75092693
Gene79	Gene80	Gene81	Gene82	Gene83	Gene84
0.44281398	0.89367247	0.97513270	0.43650262	0.04346274	0.56775677
Gene85	Gene86	Gene87	Gene88	Gene89	Gene90
0.31173151	0.66785495	0.07147003	0.50816911	0.41174256	0.88276660
Gene91	Gene92	Gene93	Gene94	Gene95	Gene96
0.77119172	0.50594707	0.33313595	0.06181756	0.70019218	0.75300966
Gene97	Gene98	Gene99	Gene100		
0.45246476	0.14707687	0.74644736	0.30893299		

\$Lower

Gene1	Gene2	Gene3	Gene4	Gene5	Gene6
0.48433410	0.06639571	0.49534260	0.38405205	0.95190711	0.21904183
Gene7	Gene8	Gene9	Gene10	Gene11	Gene12
0.26120393	0.33473660	0.84752839	0.44635199	0.52694999	0.83123295
Gene13	Gene14	Gene15	Gene16	Gene17	Gene18
0.49099733	0.86211951	0.49852156	0.88941629	0.24604602	0.53061289
Gene19	Gene20	Gene21	Gene22	Gene23	Gene24
0.06988887	0.63024922	0.29217438	0.62036626	0.67310862	0.40651292
Gene25	Gene26	Gene27	Gene28	Gene29	Gene30
0.34269054	0.06604094	0.80811816	0.44000841	0.23527115	0.31762206
Gene31	Gene32	Gene33	Gene34	Gene35	Gene36
0.54603266	0.31855290	0.88204881	0.74969558	0.11637415	0.21480980
Gene37	Gene38	Gene39	Gene40	Gene41	Gene42
0.91917974	0.80777398	0.43594595	0.65421671	0.03049579	0.27415773
Gene43	Gene44	Gene45	Gene46	Gene47	Gene48
0.26519406	0.84717728	0.28281592	0.16416526	0.59425627	0.87544021
Gene49	Gene50	Gene51	Gene52	Gene53	Gene54
0.93067425	0.39663296	0.94751834	0.49749211	0.48940612	0.26674451
Gene55	Gene56	Gene57	Gene58	Gene59	Gene60
0.13480727	0.45896793	0.44103779	0.68903294	0.70440087	0.91743615
Gene61	Gene62	Gene63	Gene64	Gene65	Gene66
0.23268990	0.93960204	0.76422490	0.85022612	0.58113201	0.47972455
Gene67	Gene68	Gene69	Gene70	Gene71	Gene72
0.45162223	0.24300463	0.61595421	0.18366862	0.75648291	0.45385300
Gene73	Gene74	Gene75	Gene76	Gene77	Gene78
0.38026636	0.52058762	0.81964785	0.44147717	0.87943441	0.24907307
Gene79	Gene80	Gene81	Gene82	Gene83	Gene84
0.55718602	0.10632753	0.02486730	0.56349738	0.95653726	0.43224323
Gene85	Gene86	Gene87	Gene88	Gene89	Gene90
0.68826849	0.33214505	0.92852997	0.49183089	0.58825744	0.11723340

Gene91	Gene92	Gene93	Gene94	Gene95	Gene96
0.22880828	0.49405293	0.66686405	0.93818244	0.29980782	0.24699034
Gene97	Gene98	Gene99	Gene100		
0.54753524	0.85292313	0.25355264	0.69106701		

\$Weight

Exp 1	Exp 2	Exp 3
3	6	8

## 6 Setup

This vignette was built on:

```
> sessionInfo()
```

R version 3.0.1 (2013-05-16)

Platform: x86\_64-apple-darwin10.8.0 (64-bit)

locale:

```
[1] ja_JP.UTF-8/ja_JP.UTF-8/ja_JP.UTF-8/C/ja_JP.UTF-8/ja_JP.UTF-8
```

attached base packages:

```
[1] splines    parallel  stats      graphics  grDevices  utils
[7] datasets  methods   base
```

other attached packages:

```
[1] metaSeq_0.99.0    snow_0.3-12      NOISeq_2.0.0
[4] Biobase_2.20.1    BiocGenerics_0.6.0
```

loaded via a namespace (and not attached):

```
[1] tools_3.0.1
```

## References

- [1] Tarazona, S. and Garcia-Alcalde, F. and Dopazo, J. and Ferrer, A. and Conesa, A. Genome Research *Differential expression in RNA-seq: A matter of depth*, 21(12): 2213-2223, 2011.
- [2] Simon Anders and Wolfgang Huber Genome Biology *Differential expression analysis for sequence count data.*, 11: R106, 2010.
- [3] Robinson, M. D. and McCarthy, D. J. and Smyth, G. K. Bioinformatics *edgeR: a Bioconductor package for differential expression analysis of digital gene expression data.*, 26: 139-140, 2010
- [4] Thomas J. Hardcastle R package version 1.14.1. *baySeq: Empirical Bayesian analysis of patterns of differential expression in count data.*, 2012.
- [5] Fisher, R. A. Statistical Methods for Research Workers, 4th edition, Oliver and Boyd, London, 1932.
- [6] Stouffer, S. A. and Suchman, E. A. and DeVinney, L. C. and Star, S. A. and Williams, R. M. Jr. The American Soldier, Vol. 1 - Adjustment during Army Life. Princeton, Princeton University Press, 1949
- [7] Trapnell, C. and Williams, B. A. and Pertea, G. and Mortazavi, A. and Kwan, G. and Baren, M. J. and Salzberg, S. L. and Wold, B. J. and Pachter, L. Nature biotechnology *Transcript assembly and quantification by RNA-Seq reveals unannotated transcripts and isoform switching during cell differentiation*, 28: 511-515, 2010.