

S7 Text: Gene Set Enrichment Analysis.

SCRaPL: A Bayesian hierarchical framework for detecting technical associates in single cell multiomics data.

Christos Maniatis^{1*}, Catalina A. Vallejos^{2,3*}, Guido Sanguinetti^{4,1*}

1 School of Informatics, The University of Edinburgh, Edinburgh, UK

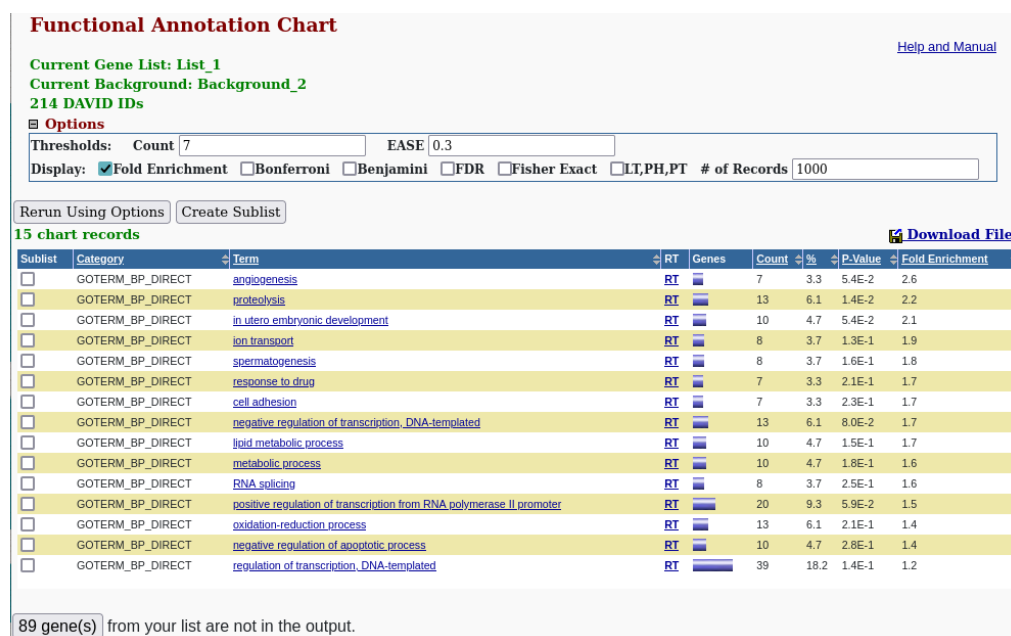
2 The Alan Turing Institute, London, UK

3 MRC Human Genetics Unit, Institute of Genetics and Cancer, Western General Hospital, The University of Edinburgh, , Edinburgh, UK

4 International School for Advanced Studies (SISSA-ISA), Trieste, Italy

* s1315538@sms.ed.ac.uk(CM);* catalina.vallejos@ed.ac.uk(CAV);* gsanguin@sissa.it(GS)

To get a round understanding of the genes with strong regulatory action identified by alternative approaches, we carry our Gene Set Enrichment Analysis (GSEA) analysis using DAVID [1]. This process helps us to identify important biological processes by looking at over-represented genes in a starting pool of genetic markers. Using SCRaPL correlation outcomes we can determine a lists of biological outcomes that practitioners would discover. The γ threshold for SCRaPL was set to 0.205 or 90% quantile of the folded distribution constructed by correlation samples of the permuted dataset.



(A)

Fig A. GO analysis with features detected by SCRaPL.

Functional Annotation Chart

[Help and Manual](#)

Current Gene List: **imp_gene_sp**

Current Background: **all_gene**

85 DAVID IDs

Options

Thresholds: Count EASE
Display: Fold Enrichment Bonferroni Benjamini FDR Fisher Exact LT,PH,PT # of Records

[Rerun Using Options](#) [Create Sublist](#)

7 chart records

[Download File](#)

Sublist	Category	Term	RT	Genes	Count	%	P-Value	Fold Enrichment
<input type="checkbox"/>	GOTERM_BP_DIRECT	meiotic cell cycle	RT		10	11.8	1.5E-8	15.0
<input type="checkbox"/>	GOTERM_BP_DIRECT	spermatid development	RT		7	8.2	2.4E-5	11.8
<input type="checkbox"/>	GOTERM_BP_DIRECT	negative regulation of transcription, DNA-templated	RT		8	9.4	4.8E-2	2.4
<input type="checkbox"/>	GOTERM_BP_DIRECT	multicellular organism development	RT		9	10.6	6.8E-2	2.0
<input type="checkbox"/>	GOTERM_BP_DIRECT	negative regulation of transcription from RNA polymerase II promoter	RT		9	10.6	1.3E-1	1.8
<input type="checkbox"/>	GOTERM_BP_DIRECT	regulation of transcription from RNA polymerase II promoter	RT		11	12.9	1.9E-1	1.5
<input type="checkbox"/>	GOTERM_BP_DIRECT	cell differentiation	RT		7	8.2	2.0E-1	1.8

49 gene(s) from your list are not in the output.

(A)

Fig B. GO analysis with features detected by Spearman correlation.

The detected biological processes presented in Figs A,B have at 7 or more genes associated with them and maximum p-value 0.3. Then processes were sorted based on their enrichment score (ie. how many times larger is the detected set compared to a random subset of the pool of genes linked to a particular biological process). SCRaPL has detected many processes directly linked to regulation of transcription and regulation of transcription in promoter regions (something expected as we are looking at promoters of a methylation-expression pair at early development where methylation plays a crucial role). Apart from them there are also biological processes linked to development, like in utero embryonic development and angiogenesis with high enrichment score. Spearman also detects processes linked to transcription regulation like negative regulation of transcription or meiotic cell cycle as seen in Fig B. However no process with immediate links to development appears. For the exact same filtering parameter and genes detected with Pearson correlation, the enrichment would link genes to "regulation of transcription" with enrichment score 1.5.

References

1. Sherman BT, Lempicki RA, et al. Systematic and integrative analysis of large gene lists using DAVID bioinformatics resources. Nature protocols. 2009;4(1):44.