

Practical activity - part 2: Transcriptomic, cistromic and functional databases

In the first phase of this exercise, using public omic and functional databases, you have to characterize the expression and the regulatory elements of the *DSCAM* and *DSCAM-AS1* genes. In the second phase you will analyse the functional role of *DSCAM* protein coding gene in a neuronal disease phenotype.

Gene and protein expression databases

Access to Gene Expression Atlas (<https://www.ebi.ac.uk/gxa/home>) and using the FAQ section answer the following questions:

- What type of data are analysed in this database?
- What are the expression units that are used in this database?
- What are the two main factors that are considered to compute this measure?
- What results are shown in the Baseline Expression section and how their are represented?
- What results are shown in the Differential Expression section and how their are represented?

Search *DSCAM-AS1* in the database and analyse the results from the Baseline Expression section.

- What are the tissues in which *DSCAM-AS1* is expressed and in which level?
- Using the filter at the top of the heatmap sort the data by expression rank and identify the cell lines characterized by the highest expression of this gene.

Using the Differential Expression section identify two comparisons in which *DSCAM-AS1* is down-regulated and two in which is up-regulated. Report the log₂ fold change measured in these experiments. Selecting the results from the “Diseases” section, report the phenotype characterized by the highest *DSCAM-AS1* expression.

Search for *DSCAM* protein coding gene in the home page search field and report the human tissues and the cell lines characterized by the highest expression of this gene.

Access to Human Protein Atlas database (<http://www.proteinatlas.org/>). Using the “About / Introduction” section report the main characteristics of the Atlas contained in this database. Using the help section identify how the protein expression level is scored, then search for the DSCAM protein coding gene.

- What are the cell lines and the tissues with the highest DSCAM RNA and protein expression?
- These results are consistent across the three different databases analysed (HPA, GTex, and Fantom5)?

Cistromic data analysis

Since DSCAM-AS1 is highly expressed in MCF-7 compared to neuronal tissue in (which DSCAM protein coding gene is expressed), it is pivotal to identify the regulatory elements that could be involved in the gene overexpression in breast cancer.

Access to Washu Genome Browser (<http://epigenomegateway.wustl.edu/browser/>) and select human hg19 genome. Using the public track hub load the tracks from “Encyclopedia of DNA Elements”. Search DSCAM-AS1 by clicking on the genomic coordinates on the top left. Using the “Tracks” option on the top right load the following datasets (in epigenetic marks, first column) for the MCF-7 cancer cell lines:

In Epigenetic mark (first column)

- USC ChipSeq MCF-7 H3K36me3B
- USC ChipSeq MCF-7 H3K27me3B
- USC ChipSeq MCF-7 H3K27ac
- USC ChipSeq MCF-7 H3K9me3
- USC ChipSeq MCF-7 Input (is a negative control ChIP-Seq experiment)
- UW ChipSeq MCF-7 H3K4me3 (both replicates)
- UW ChipSeq MCF-7 Input (is a negative control ChIP-Seq experiment)

In Transcriptional Regulators

- HudsonAlpha ChipSeq MCF-7 CTCF (both replicates)
- HudsonAlpha ChipSeq MCF-7 p300 (both replicates)
- HudsonAlpha ChipSeq MCF-7 RevX (is a negative control ChIP-Seq experiment)

- UT-A ChipSeq MCF-7 Pol2 (both replicates)
- UT-A ChipSeq MCF-7 Input (is a negative control ChIP-Seq experiment, you will find it in the Epigenetic mark section)

By right clicking on the yellow bar on the right, select “Option” and with the “+” symbol expand the signal tracks. Click on “Assay” on the right bar to bring similar experiments in proximity. Use the “-1” or “-1/3” button on the top to gain a clear view of upstream and downstream signals around the DSCAM-AS1 locus. Each track represents a ChIP-Seq coverage signal in which the intensity is proportional to the ChIP enrichment in a specific genomic position (expressed as number of mapped ChIP-Seq reads). The number reported on the left for each track represents the max number of reads counted for that experiment considering the analysed genomic region. By comparing these numbers, answer the following questions:

- Which epigenetic modifications have a signal on the DSCAM-AS1 locus higher than the corresponding negative control experiments?
- What is the modification with the higher signal at DSCAM-AS1 gene body?
- Which of the analysed proteins bound at DSCAM-AS1 promoter?
- Comparing the intensity of activatory and inhibitory epigenetic modifications you can confirm the transcriptionally active status of the region?

Transcriptional regulation in complex organisms is mainly driven by distal regulatory elements including enhancers and insulators. To search for DSCAM-AS1 distal regulatory regions expand the analysed genomic regions using two times the “-5” option at the top.

- Do you see any other signals mapped upstream of the DSCAM-AS1 regions?
- Which epigenetic modifications and transcriptional regulators ChIP-Seq signals are enriched in these regions?
- Using the arrow at the top of the screen, approximately what is the distance between these elements and the DSCAM-AS1 gene?

Genomic experiment like ChIA-PET and Hi-C can identify long-range chromatin interactions between genomic regions. Using the “Track / Public track hub” option, load the track “Long-range chromatin interaction experiments”. Select all the tracks “GIS-Ruan_ChiaPet_MCF-7” annotated in the third column of the track table section.

- Since each purple arc represents an experimentally validated long range interaction, the data confirm the interactions between the regulatory regions and the DSCAM-AS1 promoter?
- Finally explore the DSCAM protein coding gene promoter. Does the long range chromatin interactions between these enhancers involved this regulatory region?

Using the section “Apps / Screenshot” take a screenshot of the analysed genomic regions. The session can be saved also using the “Apps / Session” section.

Cancer data integration databases

Since DSCAM-AS1 is a gene regulated by Estrogen Receptor alpha, would be interesting to investigate in which subtype of breast cancer the Estrogen Receptor alpha coding gene (*ESR1*) is over-expressed or altered by mutational or copy number variation events.

To analyse multiple genomic and transcriptomic data of primary tumors from large cohorts of patients, the website cBioPortal can be used. Access to cBioPortal (<http://www.cbioportal.org/>). How many studies are analysed? How many of them concern breast cancer disease?

Select the study “Breast Invasive Carcinoma (TCGA, Cell 2015)” and view the study summary by clicking on the pie chart on the right. How many patients were involved in this study? What is the gene that is more frequently mutated?

Return on the main page, select the same study and all the options in the field “Select Genomic Profiles”. Digit *ESR1* gene in the form on the bottom and then press on “Submit the query”. In how many patients *ESR1* gene is altered? What is the main alteration of the gene? There is difference in survival time between patients carrying or not an altered *ESR1* gene (see the section “Survival”)?

Functional databases

Given the large amount of public biological experiments it is now easy to explore the function of a gene if belong to functional annotated gene sets. Access to Enrichr web tools (<http://amp.pharm.mssm.edu/Enrichr/#>). From “Help / Background information” section report the definition of gene set, gene set library, enrichment analysis and enrichment term. From the

“Libraries” section report the list of the gene-set libraries collected in Enrichr. Using the “Find a gene” section you can verify the annotation of a gene in these libraries. Then, searching for DSCAM gene report:

- Three biological processes involving the gene product (Ontologies / GO Biological Processes 2017b)
- Three brain region in which DSCAM gene is down-regulated (Cell types / Allen Brain Atlas)
- One disease phenotype in human in which the gene is up-regulated (Crowd / Disease Perturbation from GEO up). Report the associated Gene Expression Omnibus (GEO) identifier (GSE...).

Access to GEO (<https://www.ncbi.nlm.nih.gov/geo/>) and using the previous identifier search the information about the experiment. In the “Summary” and “Overall design” section report the information about the tissues on which the experiment was performed. What is the name of the high-throughput technology that was used for this experiment (“Platform” section)?

The following list of gene was obtained by the analysis of this experiment and it contains genes altered in stem cells from people with autism.

<https://drive.google.com/a/unito.it/file/d/0B1cpL92euPUDNkFfZ0F3WTAxRDQ/view?usp=sharing>

Copy this list of gene into the search form of the Enrichr web tool to identify the enriched functional gene sets. Exploring the Enrichr results answer the following questions by clicking on the indicated gene set library and selecting the table section. Consider results associated with a p-value lower than 0.001.

- What is a candidate transcriptional regulator of these genes? (Transcription / ChEA 2016)
- What are the main biological processes involving a significant fraction of these genes? (GO Biological Process 2017b)
- What is a brain region in which a significant fraction of these genes is up-regulated and a regions in which a significant fraction is down-regulated? (Cell Types / Allen Brain Atlas Up - Down)
- Excluding autism disorder, is there an other neurological disease in which a significant number of these genes is down-regulated? (Crowd / Disease Perturbations from GEO down).