# Gene Mutation Web Application (MUTABLE)

Data Science Institute
COLUMBIA UNIVERSITY

Authors: Junyi Yao, Yi Duan, Yiquan Li, Zhe Hou, Zining Chen
Faculty Mentor: Yufeng Shen

Data Science Capstone Project with Professor Yufeng Shen, Department of Systems Biology and Biomedical Informatics

## Intuition and Project Scope

From a micro perspective, human contains enormous amounts of genes within cells of different type, and mutations of genes are highly related to human diseases. With the large dataset of cells and mutations, we build a website that aim to curate and visualize genomic data to help exploratory analysis of genetic mutations in human diseases.

- Visualization of gene expression is important in genetic analysis. Disease relevant mutations are likely to have specific patterns in their location, therefore we use lollipop plots to visualize the mutations.
- Main learning objectives contain machine learning methods applied in genomics and genetics including data engineering, hierarchical clustering, statistical analysis, etc.

## Single Cell Data Analysis

- Curated single cell expression data of *AnnData* type (Table 1) are provided containing info about tissue, gene, cell type, expression level, etc.
- Based on the data, single cell expression profile (Figure 4) for each gene are generated including two parts:
  - A table showing gene expression with respect to cell type, in which we display average expression level, percentage of positive expression, and number of cells
  - A bar plot (Figure 1) showing positive expression percentage of each gene in each cell type. For x-axis, the cell types are ordered according to hierarchical clustering (Figure 2) so that similar cell types are put together.

| Citation | tissue | location on md22 | number of cells | number of cell types |
|---|---|---|---|---|
| La Manno, et al. 2016 | Fetal Midbrain | /fisher/Projects/SingleCellData/LaManno_2016_midbrain | 1,977 | 26 |
| Zhong, et al. 2018 | Prefrontal cortex | /fisher/Projects/SingleCellData/Zhong_2018_PrefrontalContex | 8,686 | 6 |
| Cao, et al. 2020 | Cerebellum | /fisher/Projects/SingleCellData/Cao_2020 | 1,080,771 | 9 |
| same as above | Cerebrum | same as above | 1,711,950 | 9 |
| same as above | Heart | same as above | 96,622 | 16 |
| same as above | Lung | same as above | 214,387 | 13 |

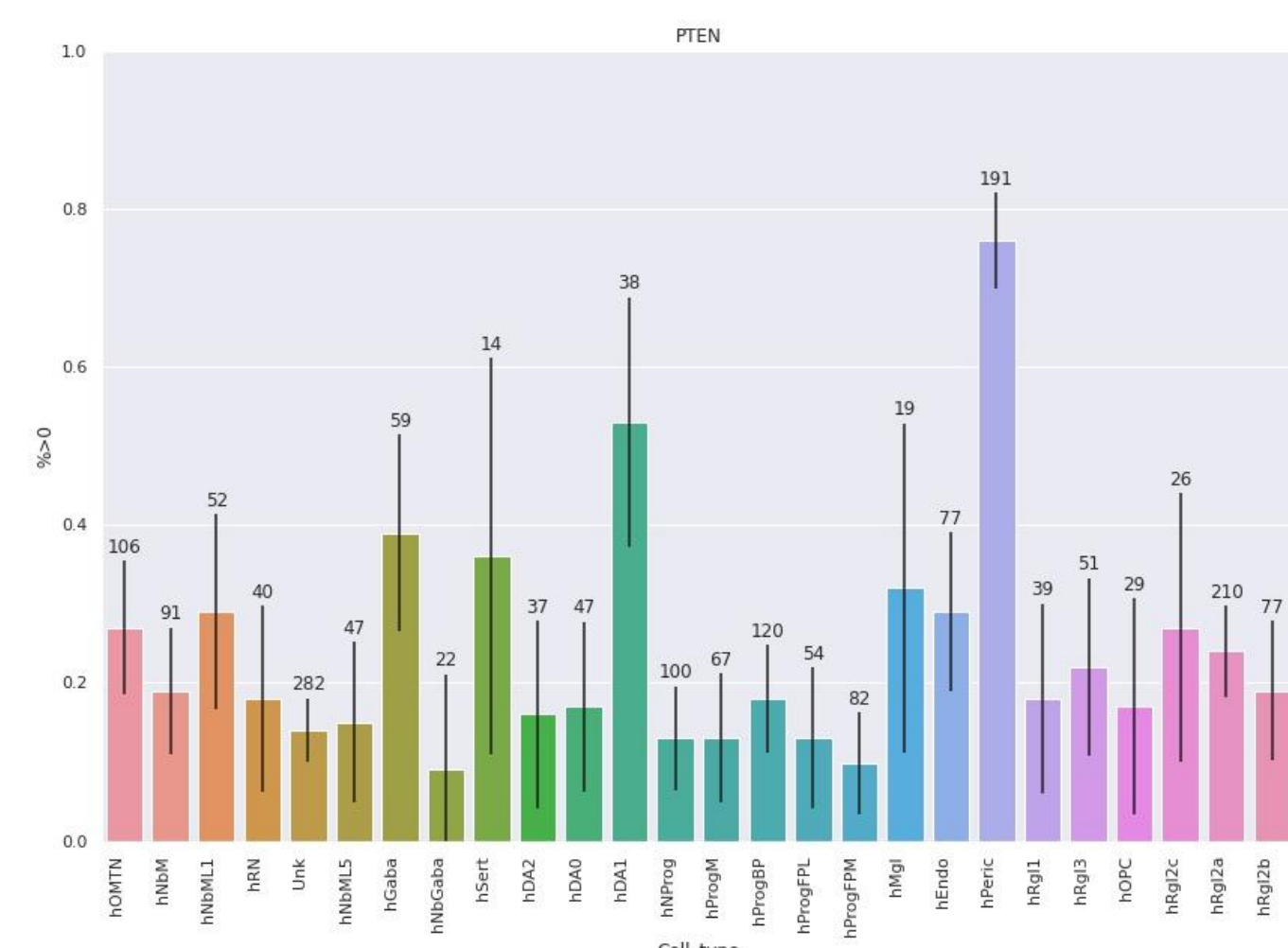Table 1. Main data sources of single cell data



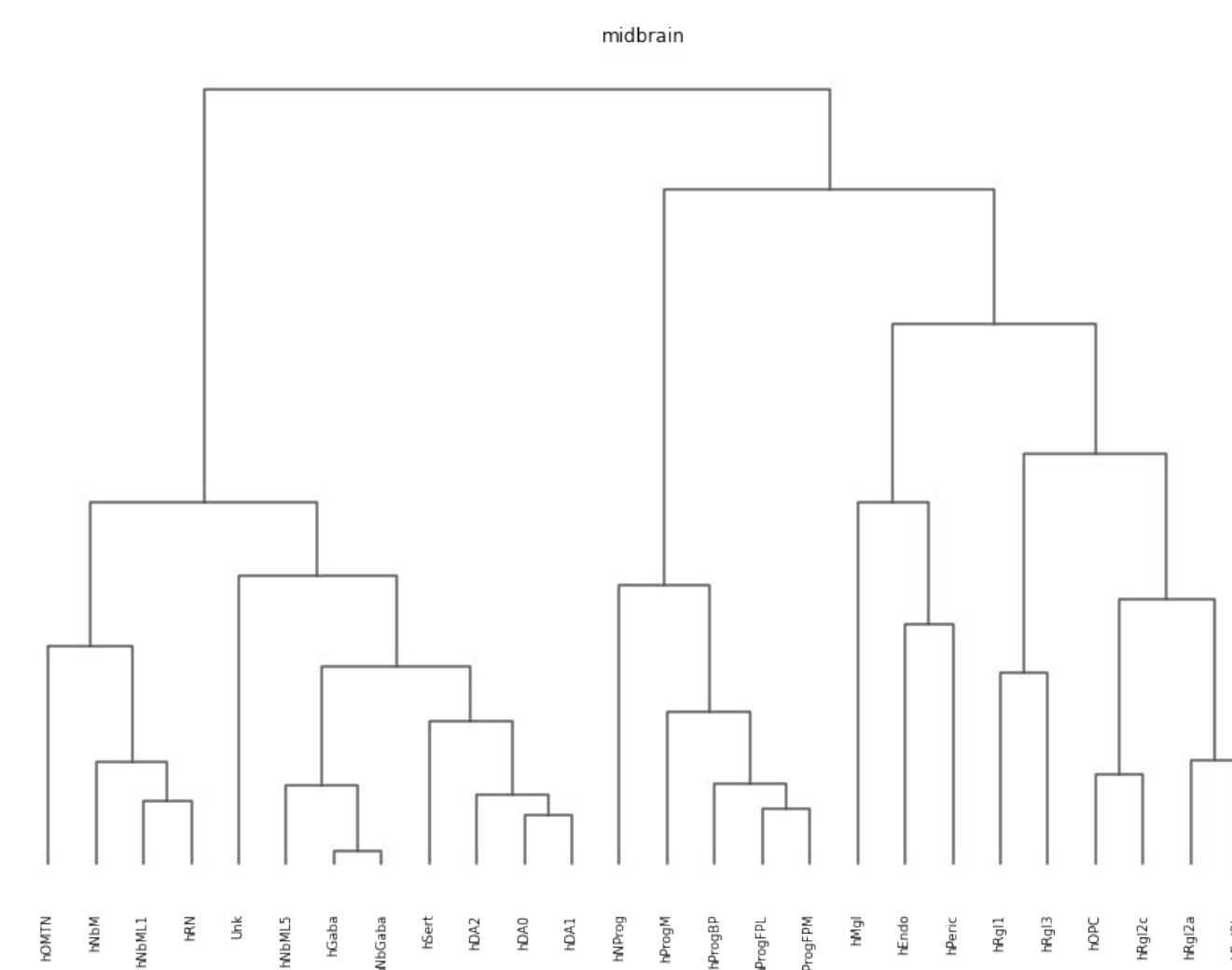Figure 1. Bar Plot of Gene PTEN in Fetal Midbrain



Figure 2. Dendrogram of Fetal Midbrain

## Acknowledgments

## Lollipop Plots for Gene Mutations

To visualize the mutations of each gene, lollipop plots are generated. Each mutation is presented as a lollipop plotted along the protein sequence. The stacked number of lollipops indicates the number of mutations on the same sequence. And the x-axis represents the protein sequence with *PFAM* protein blocks colored. Mutations are categorized into different groups based on impact (Figure 3). All mutation data and attributes are retrieved from de novo mutations (Table 2).
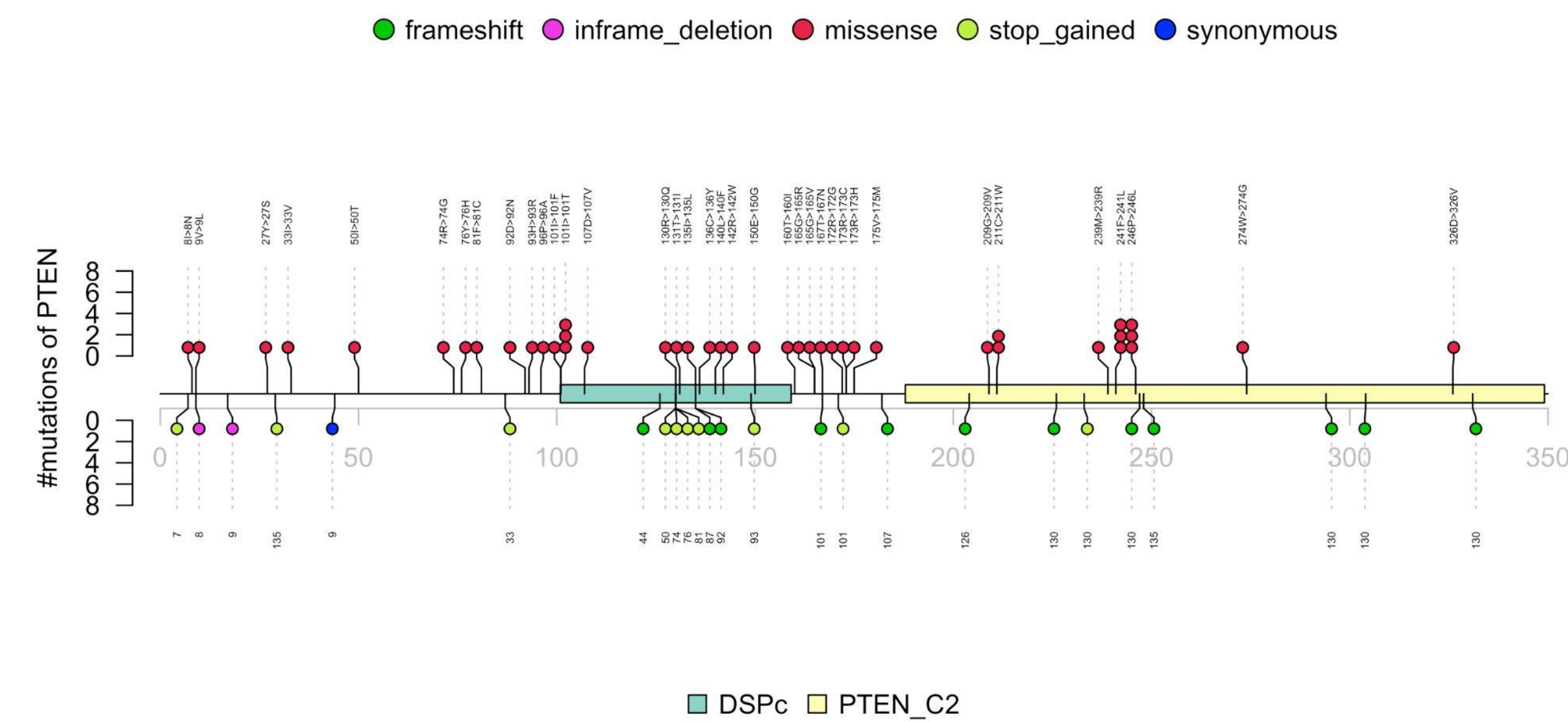


Figure 3. Lollipop Plot for Gene PTEN

| Condition | Sample Size |
|---|---|
| Autism | 23400 |
| Neurodevelopmental Disorders | 31565 |
| Congenital Diaphragmatic hernia | 595 |
| Congenital Heart Disease | 3841 |
| Esophageal Atresia | 141 |

Table 2. Data source for mutations

## Website Layout (plots excluded)

The website supports searching a specific gene in different tissues for corresponding mutation data, single cell info, bar plot & lollipop plot (plots excluded in figures below).
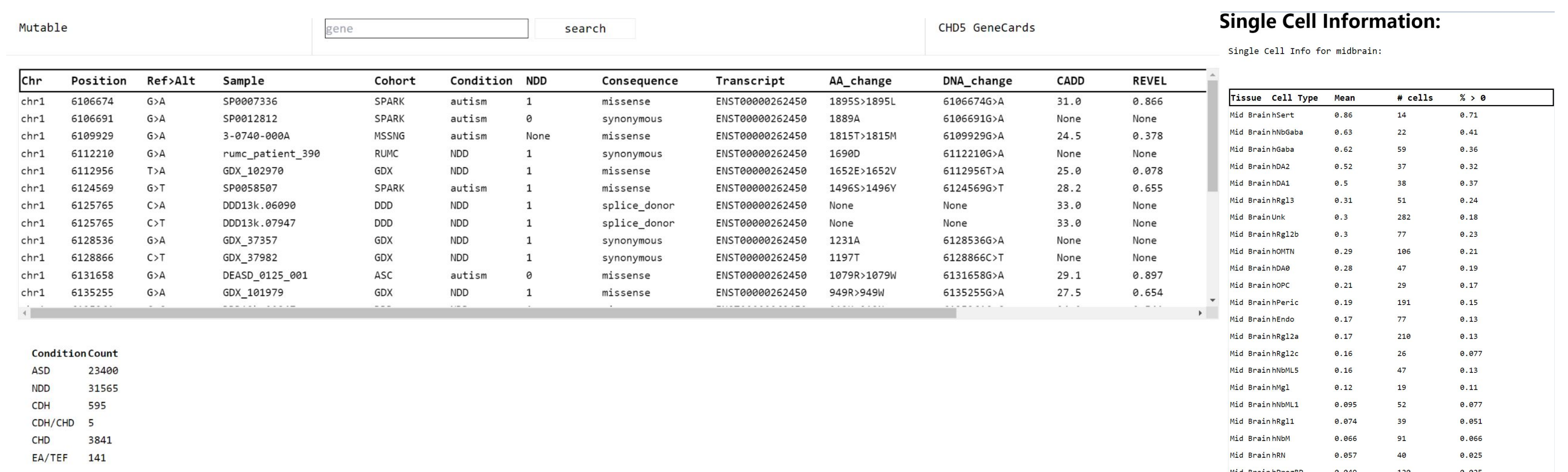


Figure 4. Website Layout when Searching for Gene CHD5

## Conclusion

Genes highly expressed in cell types relevant to a condition or diseases are more likely to be a risk gene. On the flip side, cell types in which known risk genes are highly and specifically expressed are likely to be relevant to a disease. We create a gene centric view in a web engine to support the relevant exploratory analysis of genomic data in human disease studies. The website is built reliable and sustainable for future research in Shen Lab at Columbia University.

## References

The Tabula Muris Consortium., Overall coordination., Logistical coordination. et al. Nature 562, 367–372 (2018). https://doi.org/10.1038/s41586-018-0590-4