

1 **cellxgene VIP unleashes full power of interactive visualization, plotting and**
2 **analysis of scRNA-seq data in the scale of millions of cells**

3

4 Authors: Kejie Li^{1‡}, Zhengyu Ouyang^{2‡}, Dongdong Lin¹, Michael Mingueneau¹, Will Chen¹, David Sexton¹,
5 Baohong Zhang^{1*}

6

7 Affiliations:

8 ¹Research Department, Biogen, Inc., 225 Binney St, Cambridge, MA 02142, USA

9 ²BioInfoRx, Inc., 510 Charmany Dr, Suite 275A, Madison, WI 53719, USA

10

11

12 Emails:

13 KL: kejie.li@biogen.com

14 ZO: oyoung@bioinforx.com

15 DL: dongdong.lin@biogen.com

16 MM: michael.mingueneau@biogen.com

17 WC: wwchen@post.harvard.edu

18 DS: david.sexton@biogen.com

19 BZ: baohong.zhang@biogen.com

20

21 ‡ These authors contributed equally to the work

22 * Correspondence to Baohong Zhang, baohong.zhang@biogen.com

23

Abstract

24 To meet the growing demands from scientists to effectively extract deep insights from single
25 cell RNA-seq datasets, we developed cellxgene VIP, a frontend interactive visualization plugin
26 to cellxgene framework, which directly interacts with in-memory data to generate a
27 comprehensive set of plots in high resolution, perform advanced analysis, and make data
28 downloadable for further analysis. It makes large scale scRNA-seq data visualization and
29 analysis more accessible and reproducible with the potential to become an ecosystem for the
30 scientific community to contribute even more modules to the Swiss knife of scRNA-seq data
31 exploration tool.

32

Introduction

33 Since the first single-cell RNA sequencing (scRNA-seq) study was debuted in 2009¹, over
34 1080 scRNA-seq studies have been published to date. At least 33 studies have reported
35 profiles in 200k or more cells². The largest scRNA-seq study reported 2.5 million mouse cells³.
36 It is foreseen that there will be a trend of increasing cell size of 500k or more in scRNA-seq
37 studies. The sheer amount of data had brought challenges in visualizing and exploring such
38 big data set interactively for scientists, even computational biologists.

39 Cellxgene⁴ is a leading open source scRNA-Seq data visualization tool recommended in a
40 recent evaluation⁵, which scales well in millions of cells and scores high in user experience by
41 leveraging modern web techniques with interactive features. Tested by a large community of
42 Biogen biologists, cellxgene works the best in meeting the need of handling data themselves
43 except lack of some essential plotting seen in scRNA-seq related publications and analysis
44 functions that biologists are used to, hinders its utility and limits scientists from taking
45 advantage of ever accumulating scRNA-seq data in the public to its full potential.

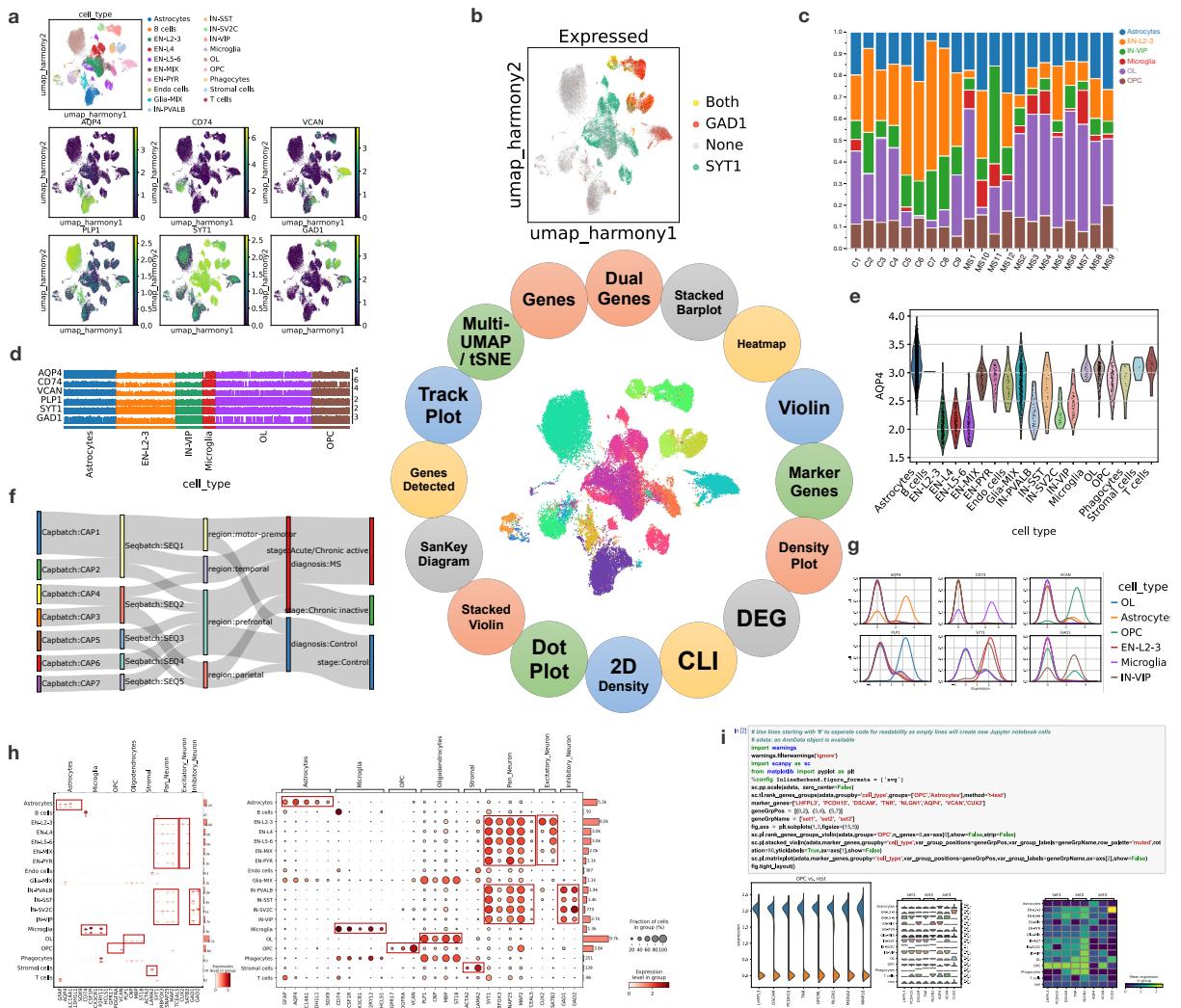
46

Result

47 To fill the gap, we developed a plugin of cellxgene named Visualization in Plugin, in short
48 VIP, to address the urgent needs of such essential functions for interactive visual exploration
49 and generation of publication-ready figures. Notably, it enhanced plotting functions significantly
50 to generate violin, stacked violin, stacked bar, heatmap, volcano, embedding, dot, track,
51 density, 2D density, sankey and dual-gene plot in high-resolution by calling server-side
52 scanpy's⁶ plotting functions and general plotting libraries as illustrated in Figure 1 and
53 Supplementary Tutorial. Innovative split plots in embedding and density plots by annotations
54 such as cell types make it easier to see which genes are more strongly expressed in which
55 cells (Supplementary Fig. S10 and Fig. S13). With the goal of becoming a platform to bridge
56 users and various tools, cellxgene VIP not only leverages python analytical capabilities, but
57 also opens to another popular language R demonstrated by volcano plot (Supplementary Fig.
58 S19). In addition, interactive Sankey diagram and stacked barplot based on yet another
59 language, JavaScript enable investigation of relationship between metadata and gene
60 expression (Supplementary Fig. S16 and Fig S17).

61 Beyond the visualization functions, we integrated diffxpy⁷ to provide analytical capabilities in
62 VIP. In addition to simple t-test and limited number of differentially expressed (DE) genes in the
63 original tool, users can get an unlimited number of DE genes by all of the methods
64 implemented in diffxpy. Volcano plot with fold change and p-value is also provided to enable
65 global view of biological perturbation (Fig. S19). To facilitate reproducible research, sessions
66 can be saved and loaded to preserve previous analysis settings with option to pre-set default
67 embedding and color-by option when the session is loaded (Supplementary for more details).

68 VIP provides easy-to-use server-side installation script of cellxgene with customized minimal
69 modifications. It took advantage of the cached data by cellxgene server to avoid data reloading
70 that is quite time consuming for large data set in gigabytes. On the client-side, a non-



71

72 **Figure 1 | cellxgene VIP serves as an ecosystem of plugins which provide essential functions for**
 73 **publication-ready, interactive visualization, as well as Command Line Interface (CLI) for analytics.** (svg files
 74 were assembled by figureComposer⁸ with zoomable version available at <https://bit.ly/2QqdMg3> that is best viewed
 75 by Chrome) (a) **Multi-tSNE/UMAP plot** visually highlights which cells expressing cell markers on selected
 76 embedding (UMAP based on harmony batch correction in this example). (b) **Dual-gene plot** highlights cells express
 77 SYT1 and GAD1 (green SYT1 only, red GAD1 only, yellow co-expression of SYT1 and GAD1), expression cutoff
 78 2.2. (c) **Stacked barplot** demonstrates the fraction of each major cell type across each sample (C are Control and
 79 MS are MS patients). (d) **Trackplot** shows expression of lineage marker genes across individual cells in annotated
 80 clusters. (e) **Violin plot** shows the AQP4 gene expression across cell types. (f) **Sankey diagram** (a.k.a. Riverplot)
 81 provides quick and easy way to explore the inter-dependent relationship of variables in the MS snRNAseq dataset⁹.
 82 (g) **Density plots** shows expression of marker genes across annotated clusters and split across cell types. (h)
 83 **Stacked violin** and **Dot plot** are the key visualizations of selected cell markers across cell types. They highlight
 84 their selective expression and validates the scRNASeq approach and visualization method. (i) **CLI** exposed by mini
 85 Jupyter Notebook to provide maximal flexibility of performing various analytics on the whole or sliced single cell
 86 dataset.

87 invasive JavaScript panel was elegantly plugged in to bridge the communications between the
88 user and server through web browser store (Supplementary Methods). With extensibility in
89 mind, we decoupled the development of the plugin and cellxgene as much as possible, so it is
90 easy to incorporate additional tools in Python or any other languages like R on server side
91 while taking full benefits of interactive user experience offered by both cellxgene and the
92 plugin. The VIP inherited all cellxgene functions and goes beyond through real time
93 communication with cellxgene by co-localization on the same web page. Within such design,
94 developers can easily add new functions to the plugin without touching cellxgene code at all.

95 Loading half-million-cell dataset even with sparse matrix compression applied would take up
96 approximately 20G of memory. Standard Mac/PC laptop would not be able to load such a
97 dataset into memory, not to mention doing additional interactive manipulation and plotting.
98 Within cellxgene VIP, all heavy lifting functions such as differential gene expression analysis
99 are executed on the server end. It greatly reduces the burden on the client browser side.

100 To further empower advanced users with programming skills, a Command Line Interface
101 (CLI) is provided to offer unlimited analytical potentials by calling any Python and/or R
102 packages such as Seurat¹⁰ (sample vignette in GitHub) installed on the server directly through
103 a mini Jupyter notebook environment (Supplementary Fig. S21).

104 In addition to scalability, versatility and flexibility, we paid great attention to user-friendliness
105 by crafting the plugin to allow users to change the look-n-feel of VIP, set figure options like font
106 size, image format and resolution, color palette for embedding (Supplementary Fig. S5),
107 initialize cellxgene embedding, coloring and brief description of the data (Supplementary
108 Methods), sync cell selection, gene input from main window to VIP automatically, save and
109 load sessions, and create combinatorial annotations and abbreviations (Supplementary Fig.
110 S22). Further, we made cellxgene VIP open source and free of charge to facilitate broader
111 scientific collaborations on large scale scRNA-Seq data sets.

112 In summary, our web-based interactive tool built upon cellxgene but greatly extended its

113 plotting and analytical capabilities by integrating state-of-the-art tools in this field. It allows
114 users with no programming experience to rapidly explore scRNA-seq data and create high-
115 resolution figures commonly seen in high-profile publications. Furthermore, it is the first tool to
116 the author's knowledge to allow computational biologists to write their own code to
117 communicate with the hosting server via a mini Jupyter notebook like interface. It opens up
118 unlimited capabilities even beyond the rich set of plotting functions provided in the tool.

119 **Code availability**

120 The demo of the tool is available at <https://cellxgenevip-ms.bxgenomics.com> while the source
121 code and detailed installation guide is provided at
122 https://github.com/interactivereport/cellxgene_VIP. Cellxgene VIP is released under the MIT
123 License.

124 **Data availability**

125 The demo data in h5ad format can be reproduced by following the Jupyter notebook
126 <https://bit.ly/2CeUHtQ> or downloaded from <https://bit.ly/2F5RUnQ>.
127

128 **Acknowledgements**

129 The authors are grateful and indebted to the cellxgene team from Chan Zuckerberg Initiative
130 for the development of the base framework. We would also like to thank L. Schirmer and D.H.
131 Rowitch for allowing us to use scRNA-seq data on Multiple Sclerosis for the demonstration of
132 cellxgene VIP.

133 **Author contributions**

134 K.L., Z.O. and B.Z. designed and implemented the software. All authors contributed testing of
135 the tool. K.L., D.L., M.M. and B.Z. wrote the manuscript with input from all authors. All authors
136 read and approved the manuscript.

137 Competing interests

138 The authors declare no competing interests.

139 References

- 140 1. Tang, F. et al. mRNA-Seq whole-transcriptome analysis of a single cell. *Nature Methods* **6**, 377-
141 382 (2009).
- 142 2. Svensson, V., da Veiga Beltrame, E. & Pachter, L. A curated database reveals trends in single-cell
143 transcriptomics. *bioRxiv*, 742304 (2019).
- 144 3. Rodrigues, S.G. et al. Slide-seq: A scalable technology for measuring genome-wide expression at
145 high spatial resolution. *Science* **363**, 1463 (2019).
- 146 4. Chan Zuckerberg Initiative chanzuckerberg/cellxgene: An interactive explorer for single-cell
147 transcriptomics data. Available at <https://chanzuckerberg.github.io/cellxgene/> ([Accessed: 5 May 2020]).
- 148 5. Cakir, B. et al. Comparison of visualization tools for single-cell RNAseq data. *NAR Genomics and*
149 *Bioinformatics* **2** (2020).
- 150 6. Wolf, F.A., Angerer, P. & Theis, F.J. SCANPY: large-scale single-cell gene expression data analysis.
151 *Genome Biology* **19**, 15 (2018).
- 152 7. David S. Fischer, F.R.H. Fast and scalable differential expression analysis on single-cell RNA-seq
153 data. Available at <https://github.com/theislab/diffxpy> ([Accessed: 8 May 2020]).
- 154 8. Li, K. et al. figureComposer: A web-based interactive multi-panel bio-infographic designing tool.
155 *bioRxiv*, 2020.2003.2004.976589 (2020).
- 156 9. Schirmer, L. et al. Neuronal vulnerability and multilineage diversity in multiple sclerosis. *Nature*
157 **573**, 75-82 (2019).
- 158 10. Stuart, T. et al. Comprehensive integration of single-cell data. *Cell* **177**, 1888-1902. e1821 (2019).

Supplementary Materials of cellxgene VIP

160	WEB URLs OF TOOLS AND ONLINE USER GUIDE	9
161	SUP. TUTORIAL – HOW TO USE CELLXGENE VIP.....	10
162	1. <i>Graphical user interface of cellxgene and VIP</i>	10
163	2. <i>Cell selection by categorial annotations</i>	12
164	3. <i>Cell selection by brushing on distribution of continues variables</i>	13
165	4. <i>Free hand Lasso selection on dots representing cells</i>	14
166	5. <i>VIP – Figure Option</i>	15
167	6. <i>VIP – Add Genes / Gene Sets.....</i>	16
168	7. <i>VIP – Violin Plot.....</i>	17
169	8. <i>VIP – Stacked Violin</i>	18
170	9. <i>VIP – Heatmap.....</i>	19
171	10. <i>VIP – UMAP/tSNE</i>	20
172	11. <i>VIP – Dot Plot.....</i>	21
173	12. <i>VIP – Track Plot.....</i>	22
174	13. <i>VIP – Density Plot.....</i>	23
175	14. <i>VIP – 2D Density Plot</i>	24
176	15. <i>VIP – Dual Genes.....</i>	25
177	16. <i>VIP – Sankey Diagram.....</i>	26
178	17. <i>VIP – Stacked Barplot.....</i>	27
179	18. <i>VIP – Gene Detected</i>	28
180	19. <i>VIP – DEG (Differential Expressed Genes).....</i>	29
181	20. <i>VIP – Marker Genes</i>	30
182	21. <i>VIP – Command Line Interface.....</i>	31
183	22. <i>VIP – Comb. & Abbr.</i>	32
184	23. <i>VIP – Other Functions</i>	33
185	SUP. METHOD 1 – CLIENT-SIDE INTEGRATION BY A JSPanel WINDOW (VIP)	34
186	SUP. METHOD 2 – SERVER-SIDE INTEGRATION.....	35
187	SUP. METHOD 3 – COMMUNICATION BETWEEN VIP AND CELLXGENE WEB GUI	36
188	SUP. METHOD 4 – DIFFXPY INTEGRATION.....	37
189	SUP. METHOD 5 – CREATE H5AD FILE FROM SEURAT OBJECT.....	38
190	HELPFUL TIPS	39
191	Handle nulls in categorical annotation	39
192	Display full traceback stack for debugging in VIP.....	39
193	Pitfall of using special characters.....	39
194	Potential use for bulk or pseudo bulk sample dataset	39

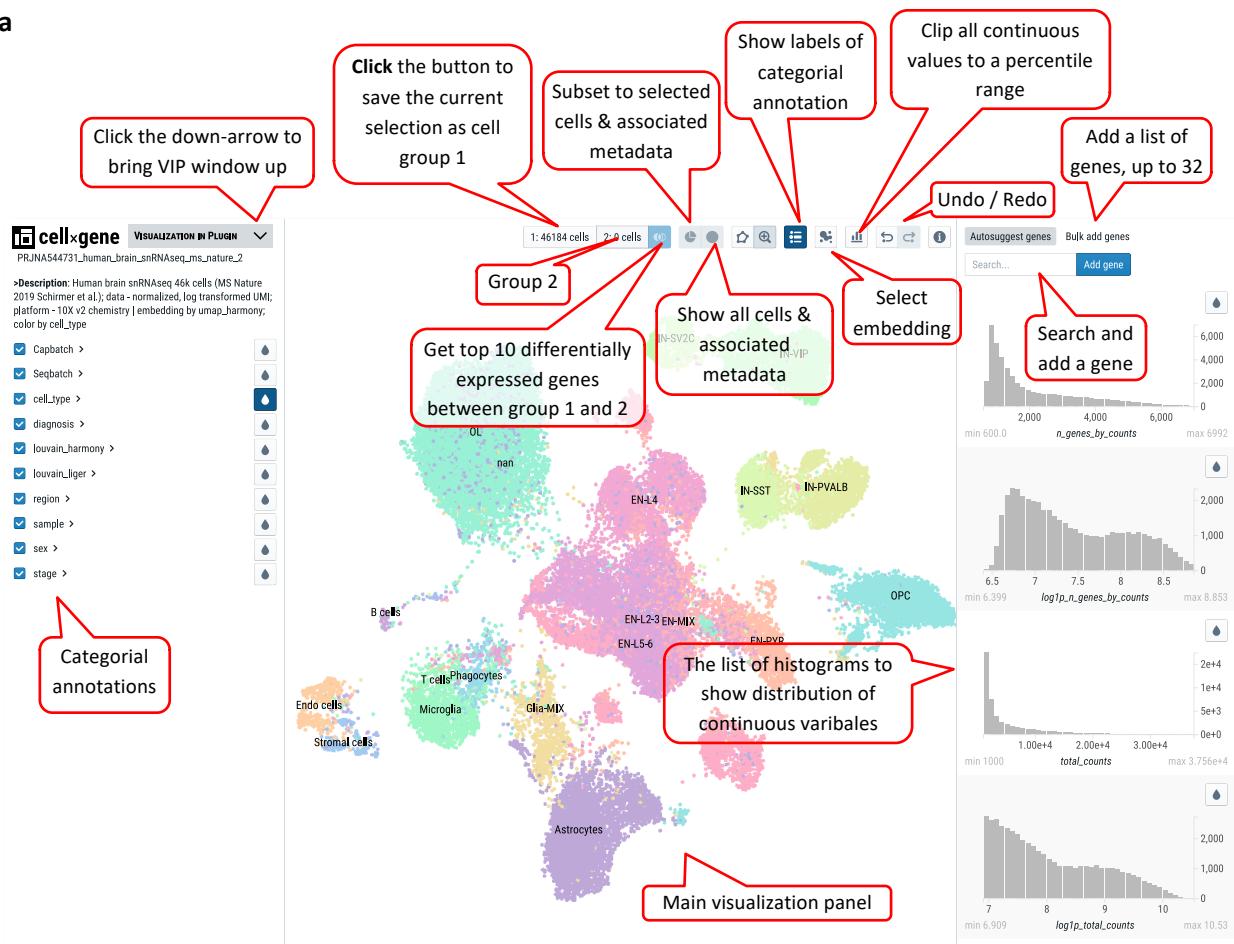
195 **Web URLs of tools and online user guide**
196
197 cellxgene VIP source code: https://github.com/interactivereport/cellxgene_VIP
198 cellxgene VIP demo sites: <https://cellxgenevip-ms.bxgenomics.com> , Schirmer / Rowitch MS,
199 Human brain snRNAseq 46k cells, Nature 2019 Schirmer et al.
200
201 cellxgene: <https://github.com/chanzuckerberg/cellxgene>
202 cellxgene tutorial: <https://cellgeni.readthedocs.io/en/latest/visualisations.html>
203 cellxgene features: <https://chanzuckerberg.github.io/cellxgene/posts/gallery>
204 cellxgene data preparation: <https://chanzuckerberg.github.io/cellxgene/posts/prepare.html>
205
206 scanpy: <https://github.com/theislab/scanpy>
207 scanpy plots: <https://scanpy-tutorials.readthedocs.io/en/latest/visualizing-marker-genes.html>
208 diffxpy: <https://github.com/theislab/diffxpy>
209 diffxpy tests: <https://diffxpy.readthedocs.io/en/latest/tutorials.html#differential-testing>
210
211 Benchmarking of interactive data visualization of single-cell RNAseq data — Batuhan Cakir :
212 https://www.youtube.com/watch?v=3nH2xi_Ni6I
213 sceasy: convertor of scRNA-Seq data formats <https://github.com/cellgeni/sceasy>
214 jupytertext: <https://jupytertext.readthedocs.io> , Jupyter Notebooks as Markdown Documents, Julia,
215 Python or R Scripts
216 nbconvert: <https://nbconvert.readthedocs.io> , Convert a Jupyter .ipynb notebook document file
217 into another static format including HTML, LaTeX, PDF, Markdown, and more

218 **Sup. Tutorial – How to use cellxgene VIP**

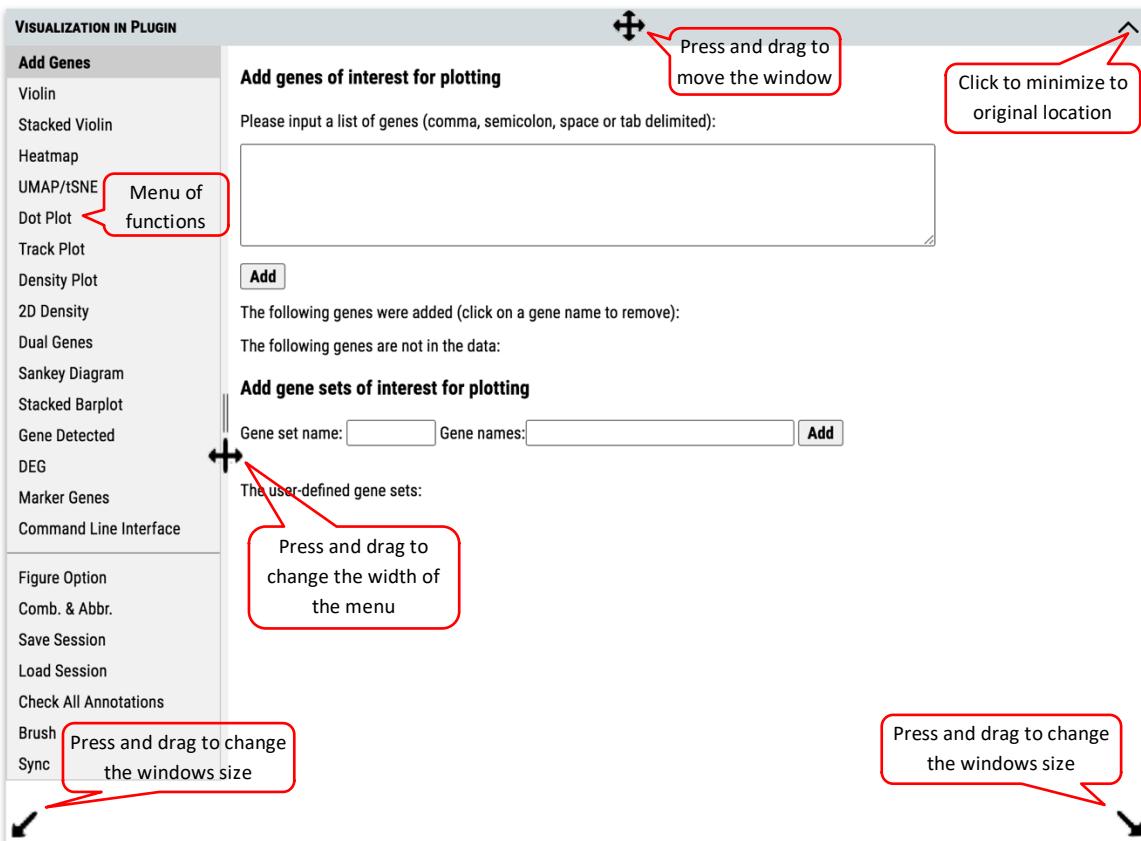
219 **1. Graphical user interface of cellxgene and VIP**

220 The main window of cellxgene is divided into three regions, the left panel mainly displays categorial
221 annotations, brief description of the data set and initial graphics setting, specifically embedding and
222 coloring of cells. On the right panel, it hosts continuous variables, such as qc metrics shown in histogram
223 with x, y corresponding to values of a measurement and numbers of cells, respectively. More
224 importantly, cells shown as individual dots are presented in the center panel based on a selected
225 embedding and colored by either categorial annotations or continuous variables, which is indicated by
226 pressed rain drop icon, e.g., cell_type in Fig. S1a.

227 **a**



232 b

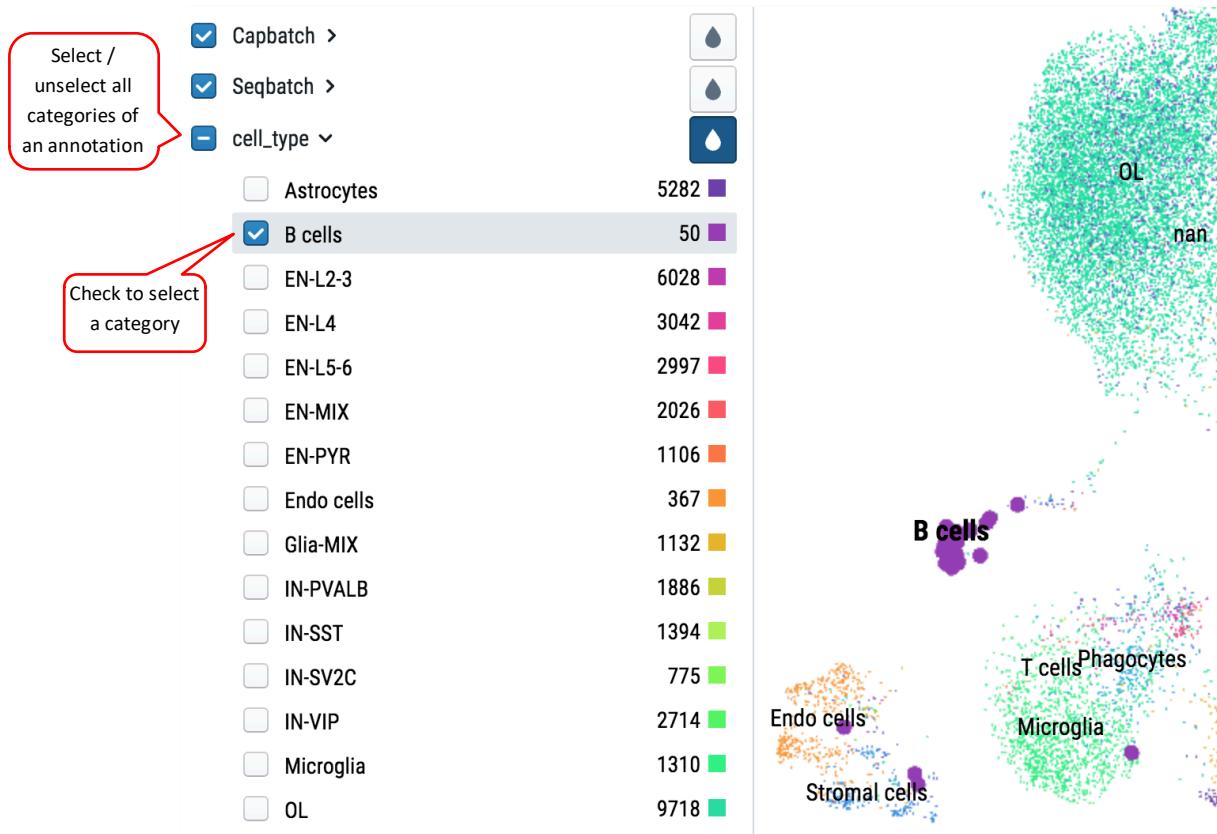


233

234 Fig. S1 | (a) Cellxgene main window, functional icons and minimized VIP bar next to cellxgene logo. (b)
235 VIP (Visualization in Plugin) window and controls of user interface. The cursor will change to
236 corresponding icon when mouse hovers over control anchors inside the window. In the case of missing
237 title bar after operation, changing the size of outside browser window (not VIP window) will always
238 bring the VIP window back to the original location near the cellxgene logo.

239 2. Cell selection by categorial annotations

240



241

242 Fig. S2 | Cell selection by categorial annotation. Selected B cells are shown in bold dots and highlighted
243 in purple color when hovering mouse over the cluster.

244 It is an overlap operation when categories from multiple annotations are checked to make the final
245 selection. E.g., if male from sex is also checked besides B cells, it means cells from B cells cluster of male
246 samples are selected.

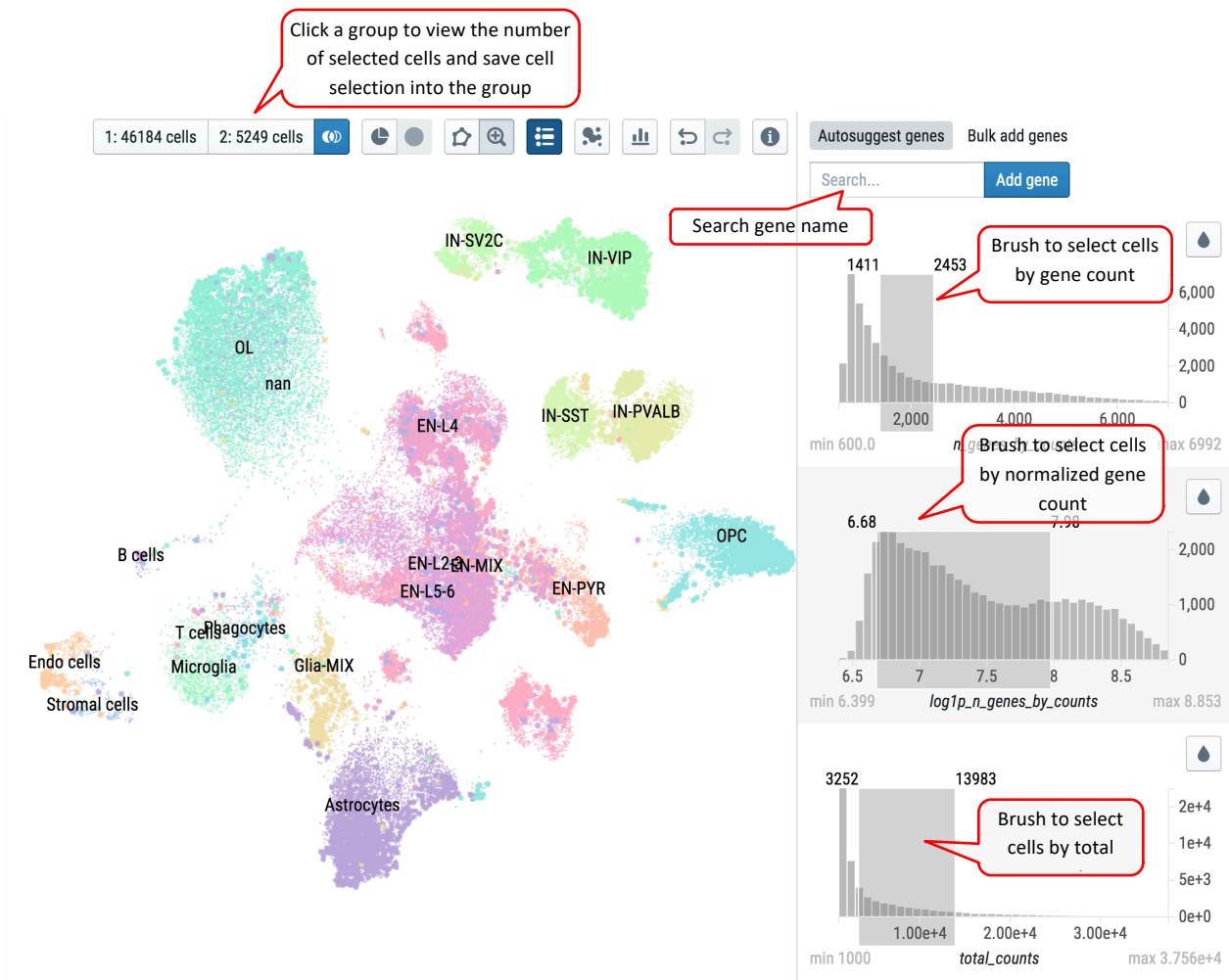
247 Note: Click "1:" or "2:" button to save cell selection into group 1 or 2.

248 3. Cell selection by brushing on distribution of continues variables

249

250

251



252

253 Fig. S3 | Cell selection by brushing the ranges of continuous variables. Low- and high-end values are
254 shown at top corners of brushing boxes in dark gray.

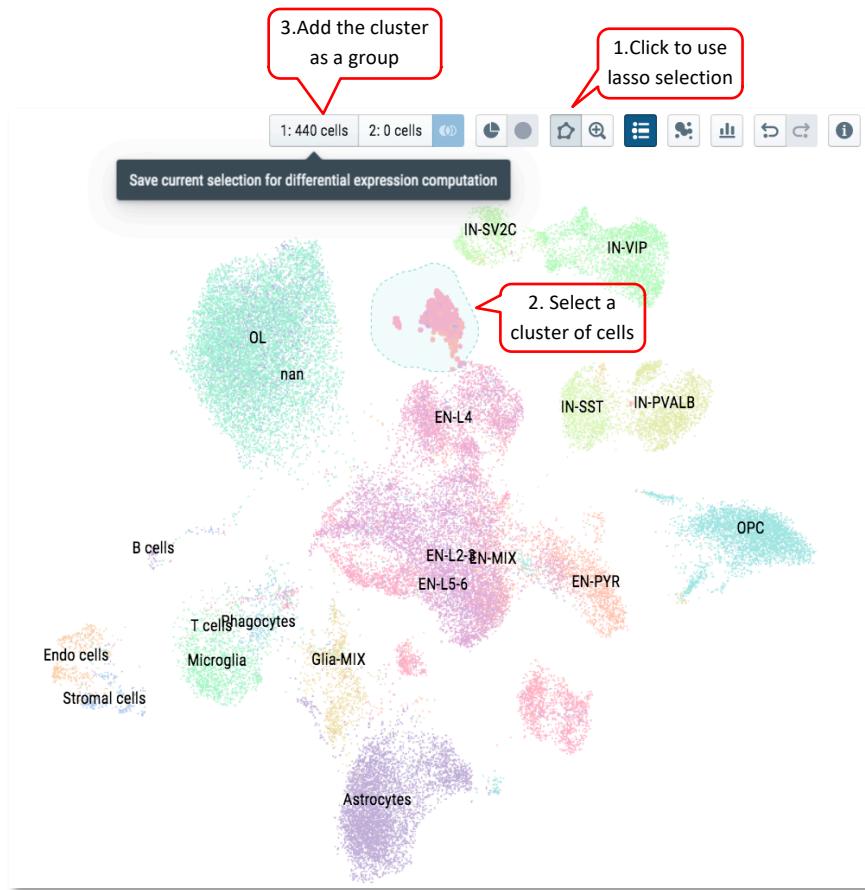
255 Note: Histograms of expression values of genes can be brushed as well to get cells expressing certain
256 genes in the range.

257 4. Free hand Lasso selection on dots representing cells

258 From the cell visualization panel, user can freely select a cluster of cells of interest by using 'Lasso'
259 selection tool. The selected cluster of cells can also be added as a group for downstream analysis in
260 cellxgene VIP.

261

262



263

264 Fig. S4 | Select cells by using free hand Lasso selection tool and add these cells as a group for further
265 analysis in cellxgene VIP.

266 Note: Please try to draw as close as possible to the starting point in the end to make an enclosed shape
267 to ensure successfully lasso selection.

268 5. VIP – Figure Option

269 User can set parameters for figure plotting that control plotting functions except CLI. ‘split_show’
270 branch of Scanpy offers better representation of Stacked Violin and Dot Plot comparing to master
271 branch.



272

273 Fig. S5 | Setting parameters for figure plotting.

274 Scaled data have zero mean and unit variance per gene. This was performed by calculating z-scores of
275 the expression data using Scanpy’s scale function. (Scanpy [pp.scale](#) function: Scale data to unit variance
276 and zero mean.)

277 We provide flexibility to allow 1) scale to unit variance or not; 2) Zero centered or not; 3) Capped at max
278 value after scaling.

279 We recommend using scaled data for plotting/visualization while using non-scaled data for differential
280 gene expression analysis.

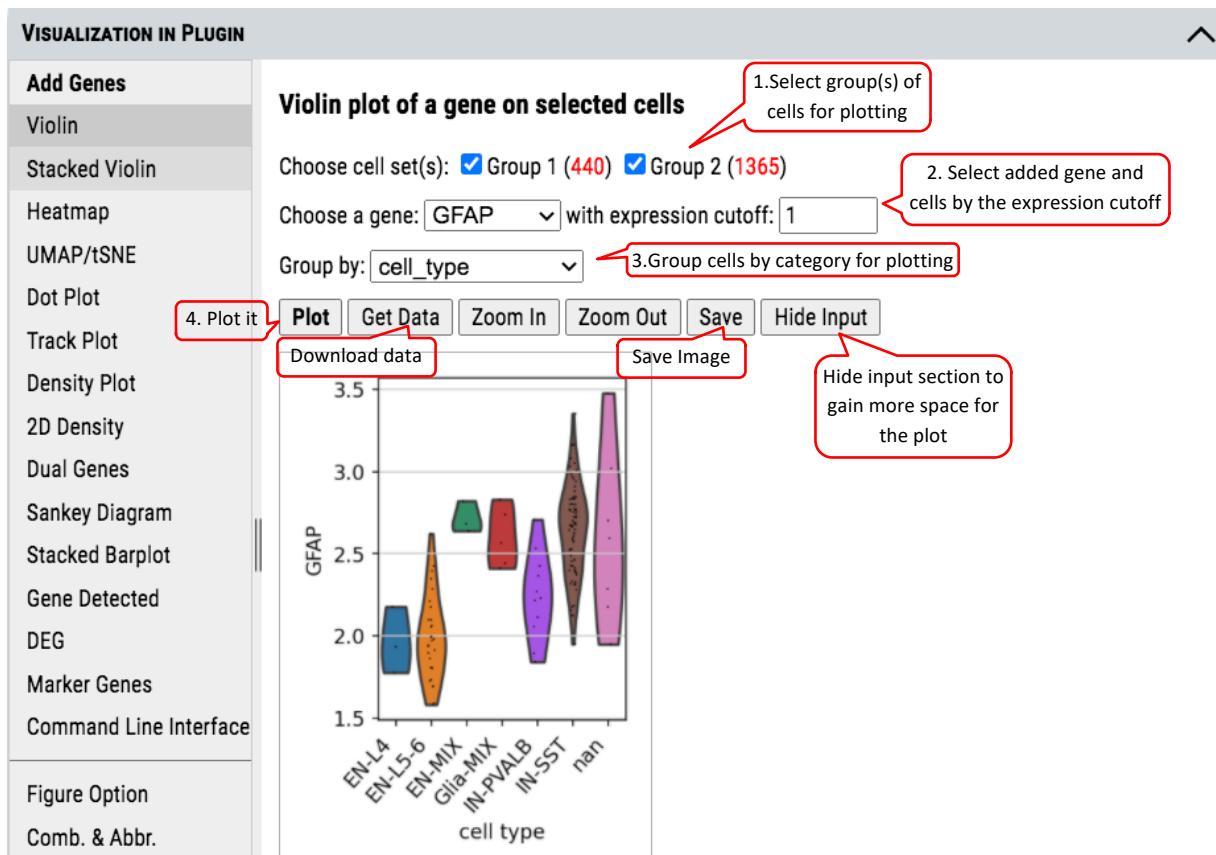
281 Note: Dot plot is one exception in visualization category which uses non-scaled data for meaningful
282 interpretation.

283 6. VIP – Add Genes / Gene Sets
284 Cellxgene VIP allows user to add any genes or gene sets for extensive exploration and visualization. User
285 can either type a list of gene in the textbox or create sets of genes to be grouped together in plots. Then
286 the genes will be automatically listed for plotting in other functional modules after checking availability
287 in the dataset.

The screenshot shows the 'Add Genes' section of the Cellxgene VIP visualization interface. On the left, a sidebar lists various plotting options: Violin, Stacked Violin, Heatmap, UMAP/tSNE, Dot Plot, Track Plot, Density Plot, 2D Density, Dual Genes, Sankey Diagram, Stacked Barplot, Gene Detected, DEG, Marker Genes, and Command Line Interface. The main area has two sections: 'Add genes of interest for plotting' and 'Add gene sets of interest for plotting'. In the first section, there is a text input field with placeholder text 'Please input a list of genes (comma, semicolon, space or tab delimited)'. Below it is an 'Add' button. A red callout bubble points to this button with the text 'Add the list of genes for plotting'. In the second section, there is a 'Gene set name:' input field containing 'set1', a 'Gene names:' input field containing 'GFAP, AQP4, SLC14A1, ALDH1L1, SOX9', and another 'Add' button. A red callout bubble points to this 'Add' button with the text 'Add the list of genes as set for plotting'. Below these fields, a message says 'The following genes were added (click on a gene name to remove): GFAP AQP4 SLC14A1 ALDH1L1 SOX9 CD74 CSF1R CX3CR1 P2RY12 HCLS1 PLP1 CNP MB ST18'. At the bottom of the 'Add gene sets of interest for plotting' section, there is a 'Delete' button followed by the text 'set1: GFAP,AQP4,SLC14A1,ALDH1L1,SOX9'. A red callout bubble points to this 'Delete' button with the text 'Delete a gene set'.

288
289 Fig. S6 | Add gene or gene sets for plotting.
290 Note: The cursor will turn to cross icon while hovering over a gene name, then click to delete the gene.

- 291 7. VIP – Violin Plot
- 292 To plot expression of gene among categories of an annotation, e.g., cell type, sex, or batch etc.
- 293 Step 1. User needs to select the group(s) of cells for plotting. These groups can be created by using
- 294 selection tools illustrated in tutorial section 2, 3 and/or 4. Initially, all of cells are gathered in ‘Group 1’
- 295 by default.
- 296 Step 2. Select a gene from the gene list which could be added as shown in section 6. An expression level
- 297 cutoff can be set to further filter out cells with low level expression of such gene.
- 298 Step 3. Select the annotation to group cells for plotting.
- 299 Step 4. Execute plotting, get plotting data (i.e., gene expression), manipulate image (e.g., zoom in/out)
- 300 or save the image.



- 301
- 302 Fig. S7 | Violin plot of gene expression values of a gene grouped by cell type.
- 303 Note: Figure resolution and format can be set in “Figure Option” tab as shown in tutorial section 5.

304 8. VIP – Stacked Violin

305 Beyond plotting expression values of a gene, stacked violin allows plotting of multiple genes together.



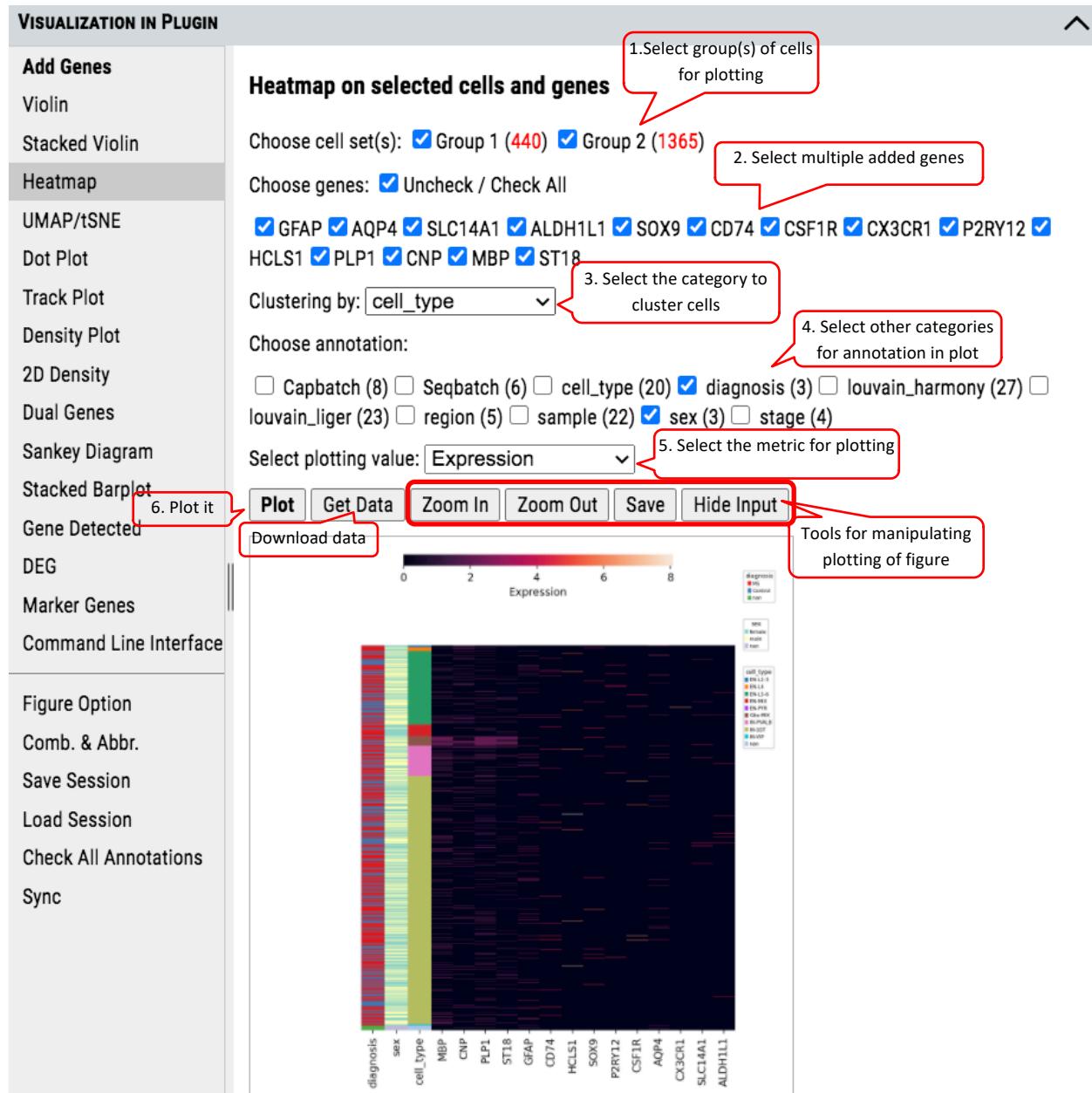
306

307 Fig. S8 | Stacked Violin plot of multiple genes and/or gene set.

308 Note: If collapsing of gene sets is set to 'Yes', average gene expression of genes in a set is used for
309 plotting.

310 9. VIP – Heatmap

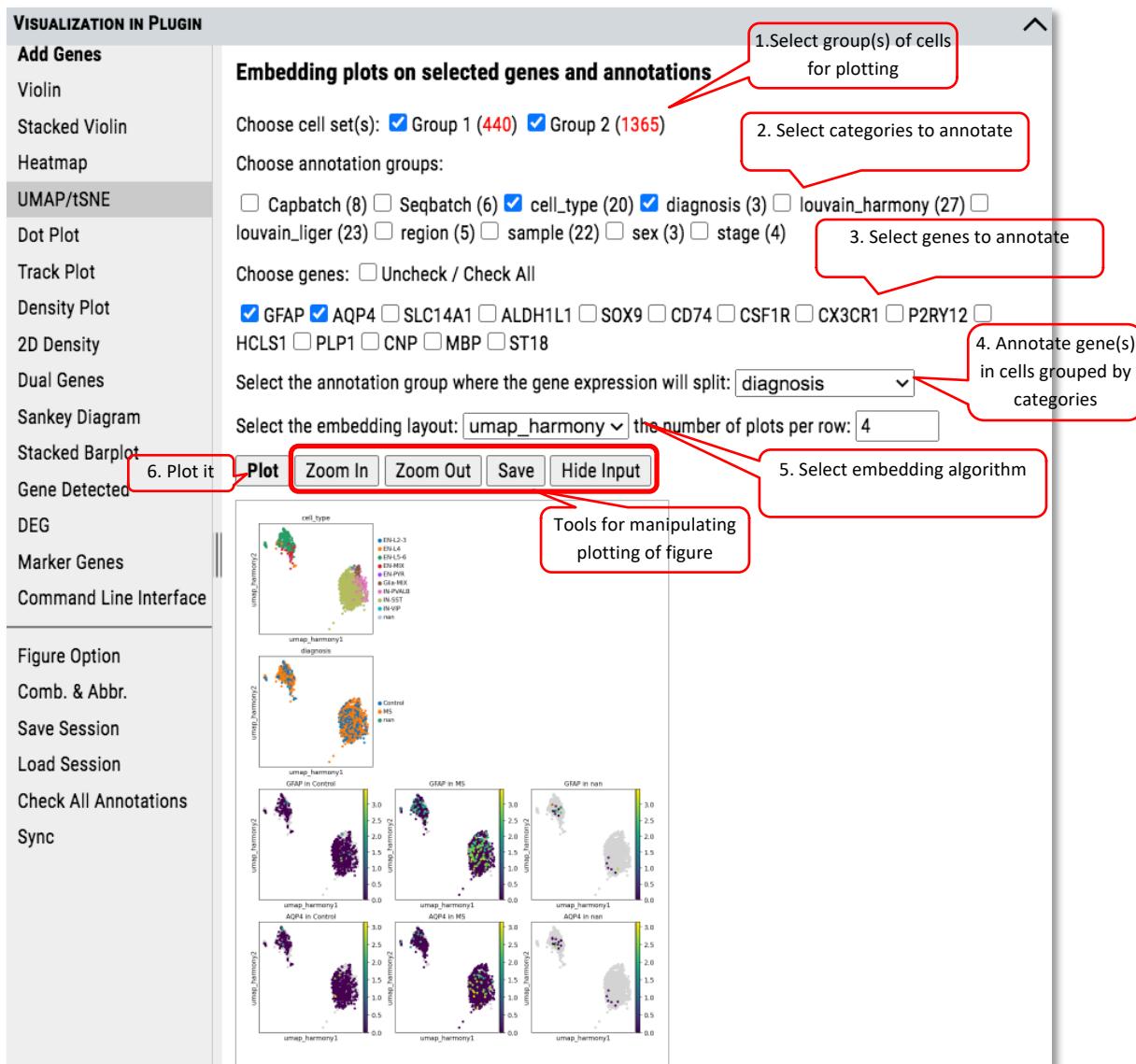
311 To show or compare the expression level (i.e., expression value or expression Z-score) of multiple genes
 312 among the selected group of cells.



313

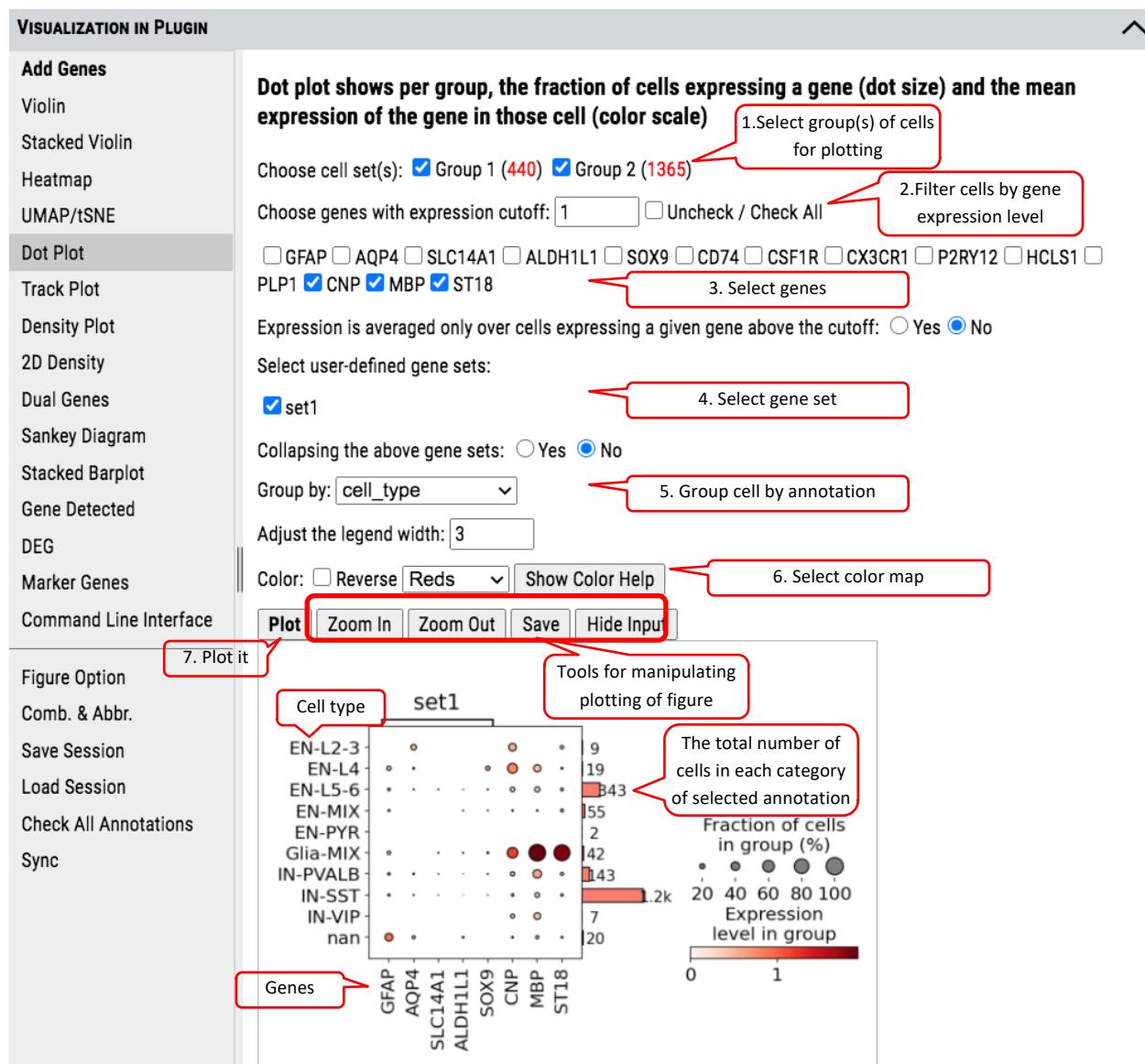
314 Fig. S9 | Heatmap of gene expression in cells grouped by annotations.

- 315 10. VIP – UMAP/tSNE
- 316 To plot the embedding of cells in the selected group(s). One of pre-computed and loaded embeddings
317 can be selected.
- 318 User can color cells in the embedding plots by multiple annotations (e.g., cell_type, diagnosis).
- 319 Besides coloring cells by annotations, user can color cells based on gene expression level of selected
320 genes in the embedding plots.
- 321



- 322
- 323 Fig. S10 | Embedding plotting of expression level of genes or gene set in the cells split by categories of
324 an annotation.

325 11. VIP – Dot Plot
 326 To show the fraction of cells (annotated by dot size) expressing a gene in each group and the averaged
 327 expression level of the gene (annotated by color intensity) in the group.

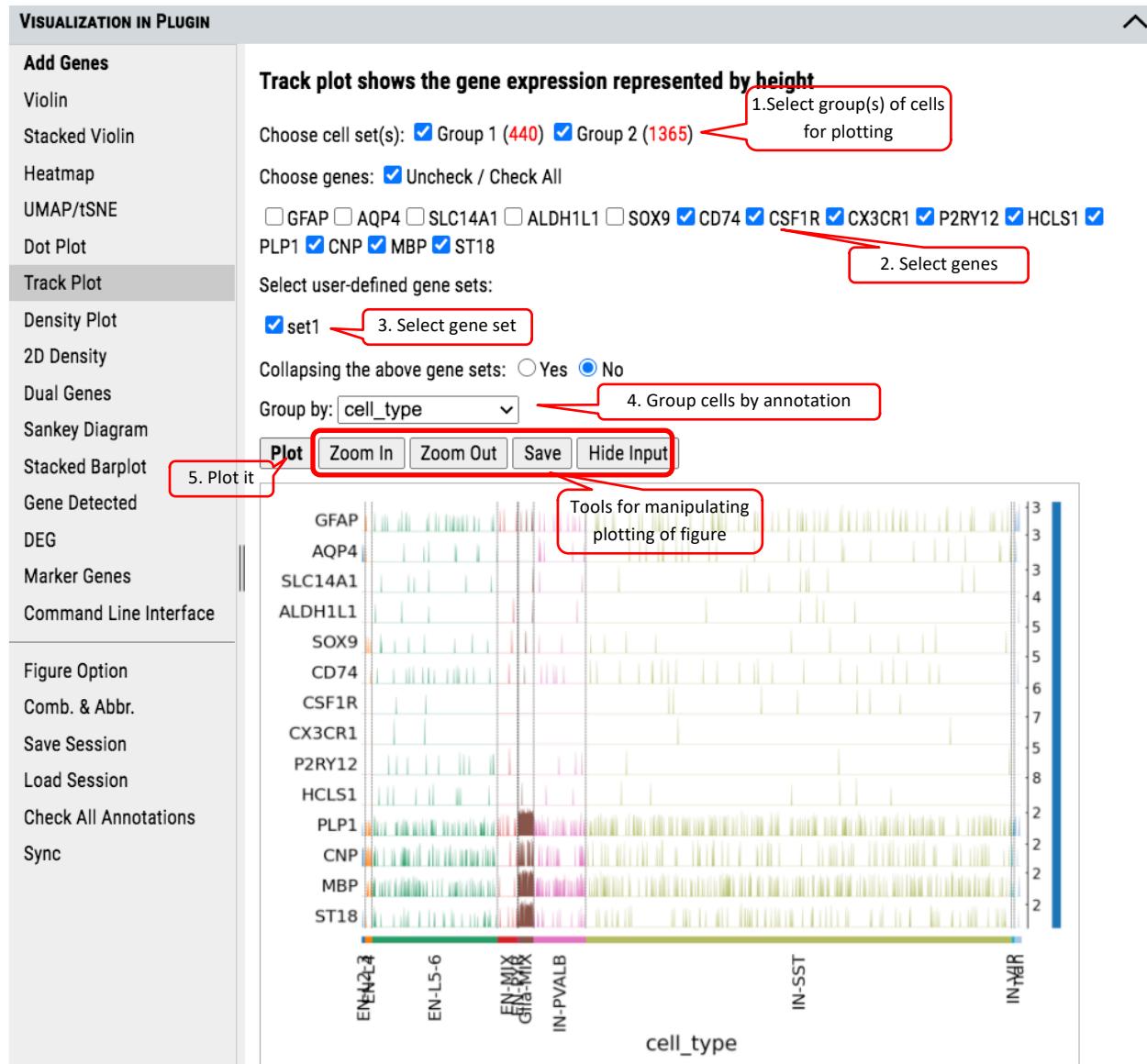


328
 329 Fig. S11 | Dot plotting of the fraction of cells expressing genes above a cutoff in each categorie of the
 330 selected annotation.

331 Note: The number of cells represented by side bar chart are always numbers of cells distributed in each
 332 category of certain annotation without filtering. It will give accurate estimate of number of cells in each
 333 bubble in the plot. The use of the plot is only meaningful when the counts matrix contains zeros
 334 representing no gene counts. Its visualization does not work for scaled or corrected matrices in which
 335 zero counts had been replaced by other values, see <https://scipy-tutorials.readthedocs.io/en/multiomics/visualizing-marker-genes.html#Dot-plots>.

337 12. VIP – Track Plot

338 To show the expression of gene(s) of individual cells as vertical lines grouped by the selected annotation
339 on x-axis. Instead of a color scale, the gene expression is represented by height.



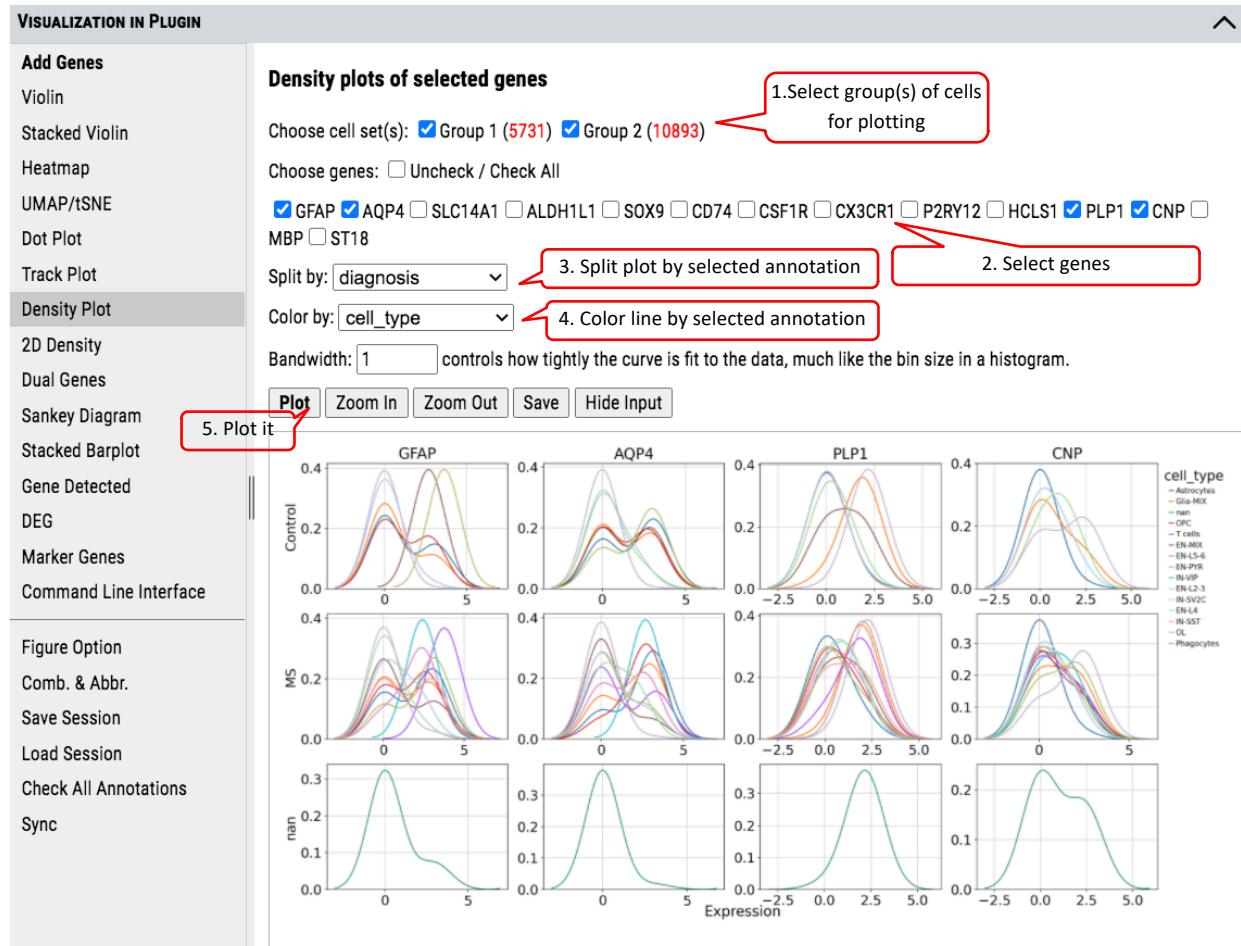
340

341 Fig. S12 | Track plotting of expression of genes or gene set in each category of the selected annotation.
342 Gene expression levels are represented by the heights of vertical lines.

343 13. VIP – Density Plot

344 To show the density of gene(s) expression in the cells annotated by category in the selected group(s) of
 345 cells. A density plot is a representation of the distribution of a numeric variable. It uses a kernel density
 346 estimate to show the probability density function of the variable (see more). It is a smoothed version of
 347 the histogram and is used in the same concept.

348 The bandwidth defines how close to a value point the distance between two points must be to influence
 349 the estimation of the density at the point. A small bandwidth only considers the closest values, so the
 350 estimation is close to the data. A large bandwidth considers more points and gives a smoother
 351 estimation.



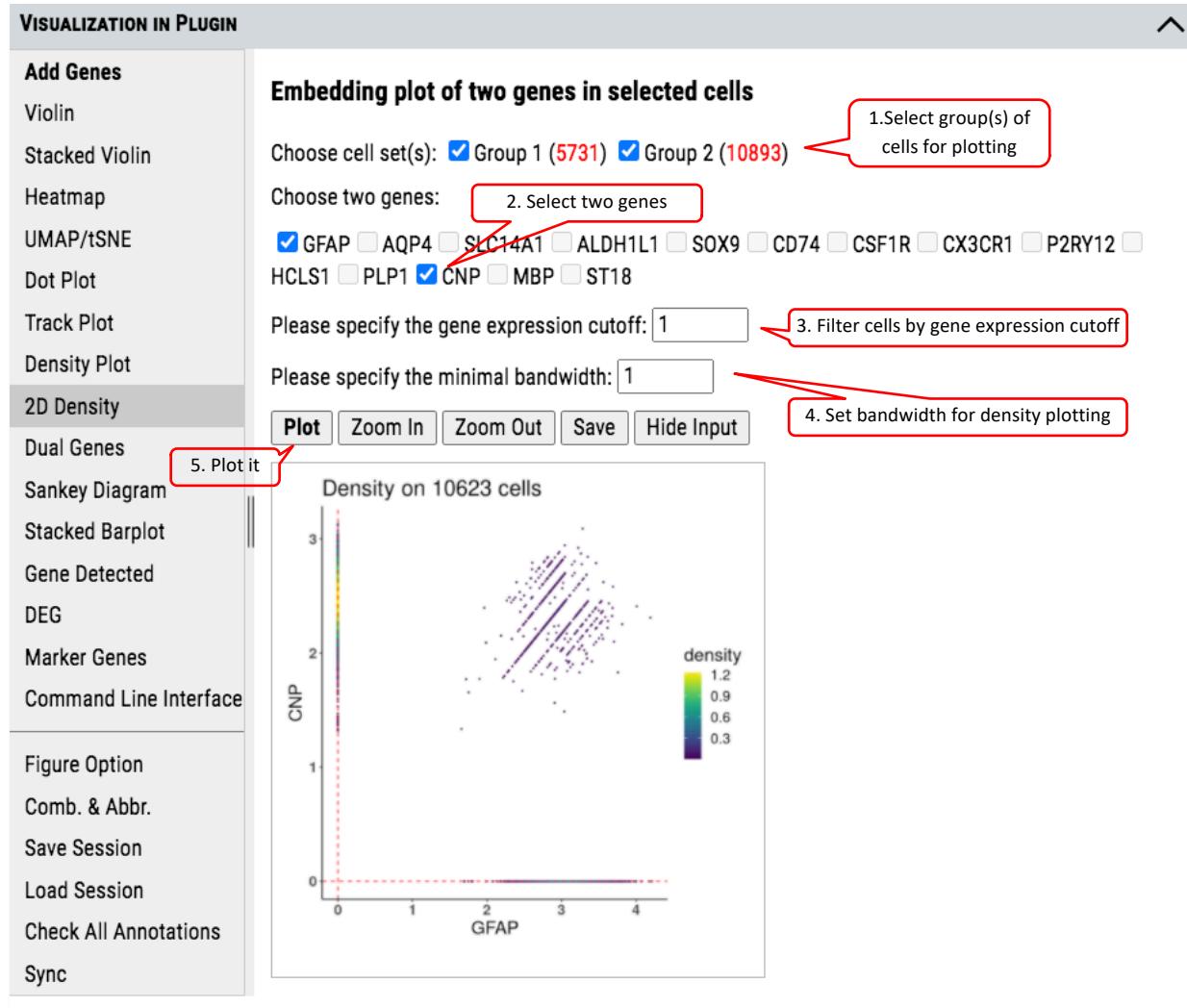
352

353 Fig. S13 | Density plotting of the expression of genes in each group split by one annotation while colored
 354 by another.

355 14. VIP – 2D Density Plot

356 Besides plotting of expression density of single gene, 2D density plot allows to explore the joint
357 expression density of two genes in the cells expressing both genes above a cutoff.

358



359

360 Fig. S14 | 2D Density plotting of expression of two genes in the selected cells.

- 361 15. VIP – Dual Genes
- 362 To view the relationship of expression levels of two genes in selected cells. It is based on the embedding plot of cells while coloring cells according to the expression levels of gene(s) in each cell.



- 364
- 365 Fig. S15 | Embedding plotting of the expression of two genes in the selected cell group(s).

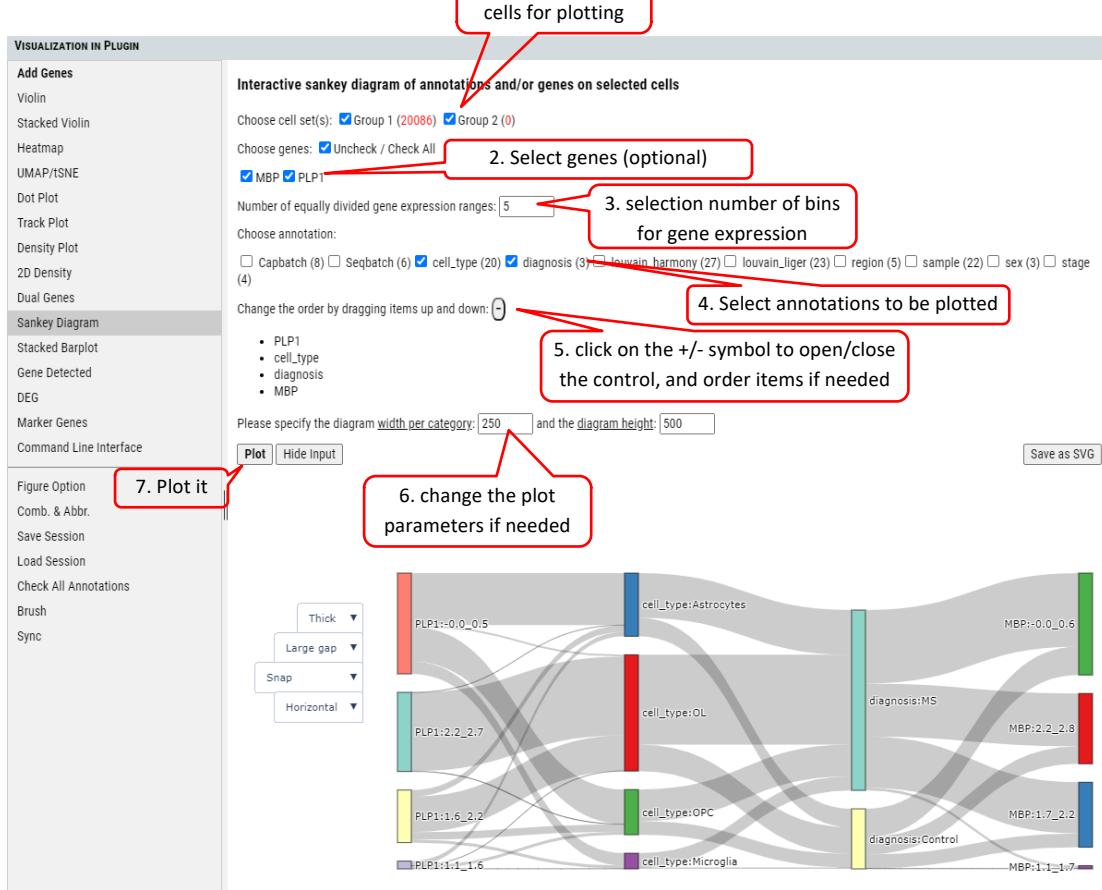
366 16. VIP – Sankey Diagram

367 Sankey diagram shows the flow of gene expression and annotations linked by cells. Gene expression is
368 divided equally into bins so user can view distribution relationship between gene expression and
369 annotations.

370 The diagram is also shown in an interactive way that user can change the layout by selecting several
371 items (e.g., thin or thick on color bar, small or large space) from the panel. Also, user can drag these
372 small boxes on the plot to get preferred layout and save it as high resolution SVG figure.

373 In addition, when you hover over mouse on a box, you can get detailed information about the source
374 and target of flow.

375



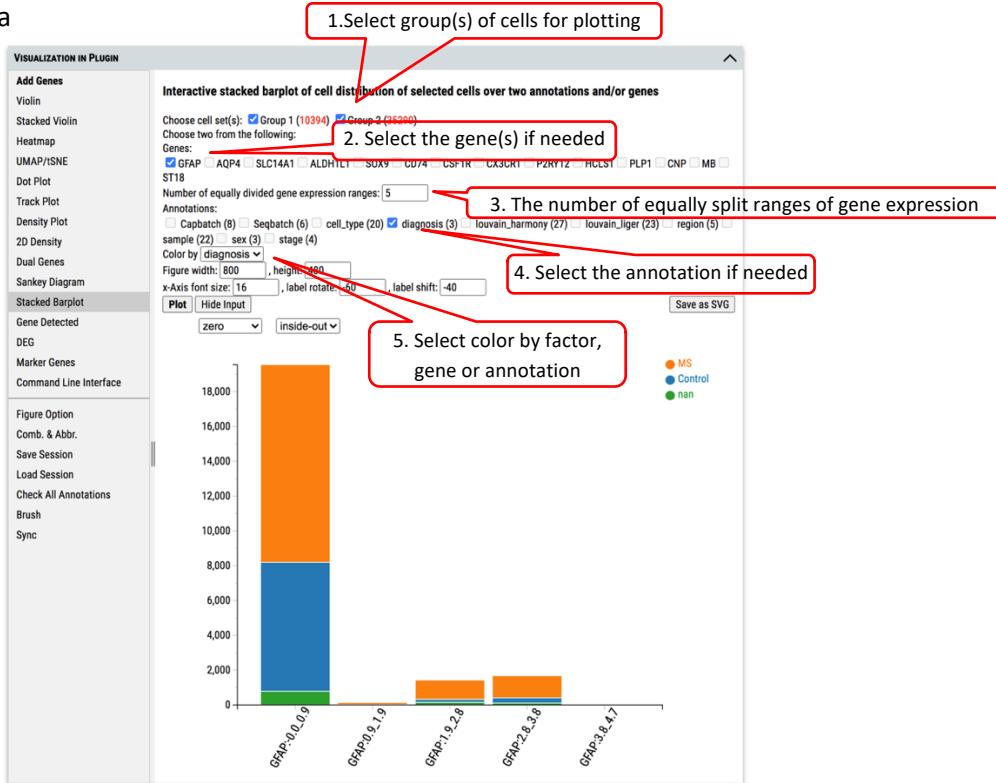
376

377 Fig S16 | Sankey diagram provides quick and easy way to explore the inter-dependent relationship of
378 variables.

379 17. VIP – Stacked Barplot

380 To show the distribution of cells among categories of an annotation and/or ranges of expression of a
 381 gene. Only two factors from annotations or genes can be chosen. The plot allows user to explore the
 382 distribution of cells in different views interactively as shown in Fig. S17b and Fig. S17c.

383 a



384
385 b

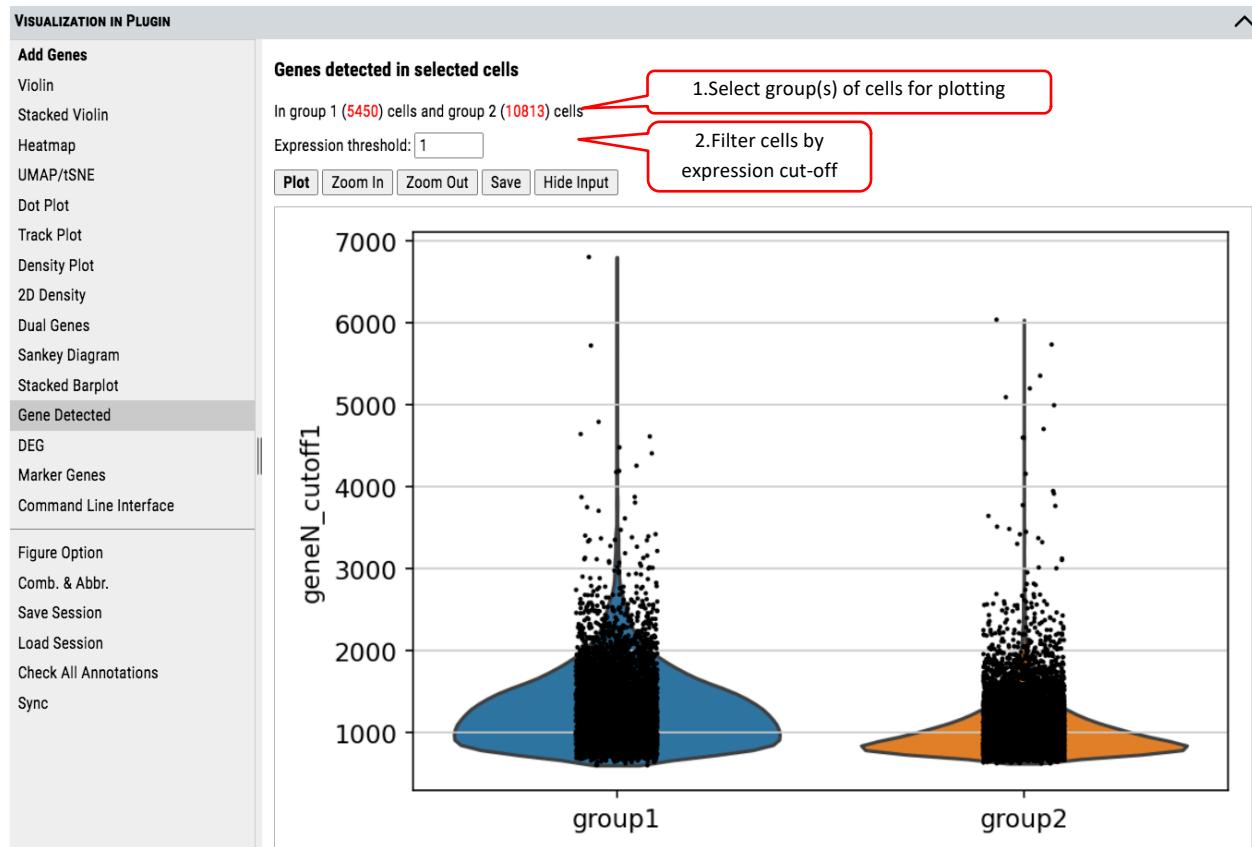


386

387 Fig. S17 | (a) The distribution of cells in the selected group(s) regarding categories of an annotation and
 388 expression ranges of a gene. (b) Expand view to show percentage instead of numbers of cells. (c) Mouse
 389 over a bar to show details. Y-axis is re-located for better view and the number of cells represented by
 390 the bar is also shown in legend.

391 18. VIP – Gene Detected

392 To show the number of genes expressed above the specified expression cut-off in the selected group(s)
393 of cells.



397 19. VIP – DEG (Differential Expressed Genes)

398 Besides plotting functions, cellxgene VIP also provides differential analysis between two selected groups

399 of cells to identify differential expressed genes.

400 Three differential analysis statistical test methods are provided including Welch's t-test, Wilcoxon rank

401 test and Wald's test. The statistical test results are presented in a table format including log2 Fold

402 change, p-value and padj value (i.e. FDR value). Please note, ***we provide users with simple test methods***

403 ***for quick exploration within the interactive framework. However, there would be covariates need to***

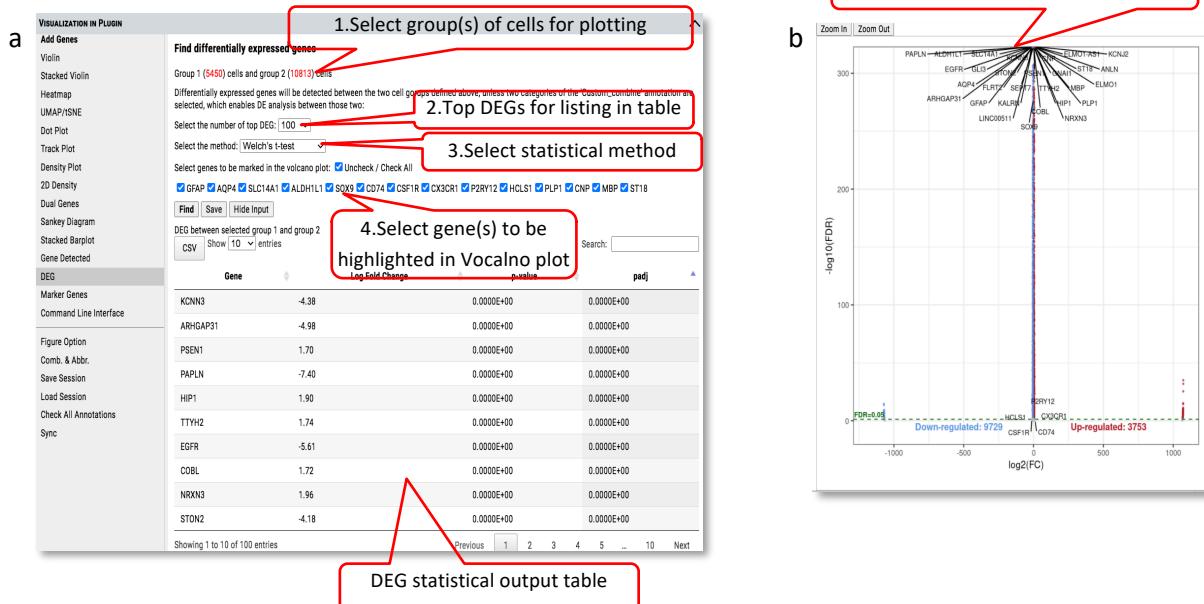
404 ***be considered in a proper statistical test. Please consult your stats experts for appropriate test by***

405 ***using the right test method and right model.***

406 Volcano plotting is also provided to show the log2FC vs. -log10(FDR) relationship for all genes. User can

407 select the gene(s) from the pre-selected gene list to be highlighted with text in the volcano plot.

408



409

410

411 Fig. S19 | (a) DEG analysis between the selected group(s) and (b) Volcano plotting of identified DEGs.

412 Note: The data used by DEG is unscaled (please refer to description of the dataset to find out what

413 preprocessing was done on the data). Scaling control in the Figure Option does not apply to DEG. The

414 three methods: 'Welch's t-test' uses t-test (assuming underlying data with normal distributions) this

415 uses **cellxgene** t-test implementation, 'Wilcoxon rank test' uses Wilcoxon rank-sum test (does not

416 assume known distributions, non-parametric test) and 'Wald's test' uses Wald Chi-Squared test which is

417 based on maximum likelihood. 'Wilcoxon rank test' and 'Wald's test' use diffxpy's implementation.

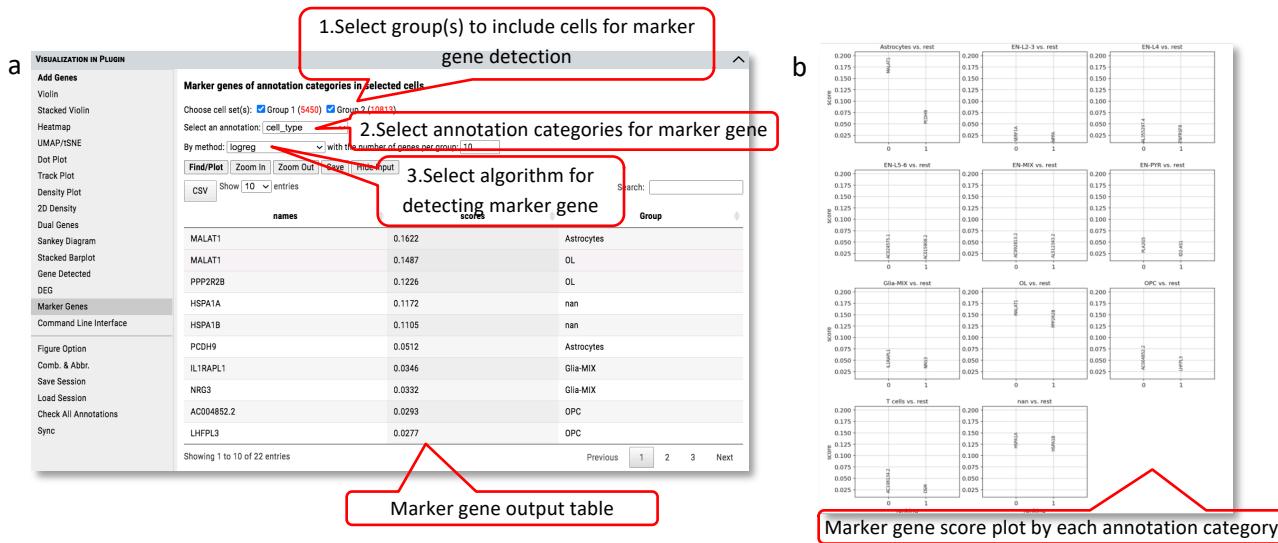
418 20. VIP – Marker Genes

419 This functional module allows user to identify marker genes in the selected group(s) (more than 2, if 2
 420 groups, please use DEG) of cells by annotation categories.

421 Four methods are provided for detecting marker genes including logreg, t-test, Wilcoxon, and t-
 422 testoverest_var. For each identified marker gene, the gene name, scores (the z-score underlying the
 423 computation of a p-value for each gene for each group) and assigned group are listed in the output
 424 table.

425 In each annotation category, top ranked marker genes (this example shows top 2) will be plotted by
 426 score in comparison to the rest of the categories.

427



428

429 Fig. S20 | (a) Marker genes detection in the selected group(s) of cells regarding to the selected
 430 annotation categories and (b) plotting of identified marker genes by each category.

431 Note: The four methods implementations by calling scanpy.tl.rank_genes_groups function: 'logreg' uses
 432 logistic regression, 't-test' uses t-test, 'wilcoxon' uses Wilcoxon rank-sum, and 't-test_overestim_var'
 433 overestimates variance of each group.

434 21. VIP – Command Line Interface

435 Although cellxgene VIP provides a rich set of visualization modules as shown above, command line
 436 interface is also built to allow unlimited visualization and analytical capabilities by power user who know
 437 how to program in Python / R languages.

Don't select any gene unless you operate on these genes only. In that case, it will speed up the execution of code in CLI by creating a much smaller AnnData

VISUALIZATION IN PLUGIN

Add Genes

Violin

Stacked Violin

Heatmap

UMAP/tSNE

Dot Plot

Track Plot

Density Plot

2D Density

Dual Genes

Sankey Diagram

Stacked Barplot

Gene Detected

DEG

Marker Genes

Command Line Interface

Figure Option

Comb. & Abbr.

Save Session

Load Session

Check All Annotations

Brush

Sync

Choose cell set(s): Group 1 (46184) Group 2 (0)

Choose genes: Uncheck / Check All

GFAP AQP4 SLC14A1 ALDH1L1 SOX9 CD74 CSF1R CX3CR1 P2RY12 HCLS1 PLP1 CNP MB ST18

Choose annotation:

Capbatch (8) Seqbatch (6) cell_type (20) diagnosis (3) louvain_harmony (27) louvain_liger (23) region (5) sample (22) sex (3) stage (4)

Choose layout:

pca_harmony tsne_harmony tsne_liger umap_harmony umap_liger

```
2 # Use lines starting with '#' to separate code for readability as empty lines will create new Jupyter notebook cells
3 # adata: an AnnData object is available
4 import warnings
5 warnings.filterwarnings('ignore')
6 import scanpy as sc
7 from matplotlib import pyplot as plt
8 %config InlineBackend.figure_formats = ['svg']
9 sc.pp.scale(adata, zero_center=False)
10 sc.tl.rank_genes_groups(adata, groupby='cell_type', groups=['OPC', 'Astrocytes'], method='t-test')
11 marker_genes=[LHFPL3, PCDH15, DSCAM, TNR, NLGN1, AQP4, VCAN, CUX2]
12 geneGrpPos = [(0,2), (3,4), (5,7)]
13 geneGrpName = [set1, set2, set3]
14 fig,axs = plt.subplots(1,3,figsize=(15,5))
15 sc.pl.rank_genes_groups_violin(adata, groups='OPC', n_genes=8, ax=axs[0], show=False, strip=False)
16 sc.pl.stack_violin(adata, marker_genes, groupby='cell_type', var_group_positions=geneGrpPos, var_group_labels=geneGrpName, row_palette='muted', rotation=90, ticklabels=True, ax=axs[1], show=False)
17 sc.pl.matrixplot(adata, marker_genes, groupby='cell_type', var_group_positions=geneGrpPos, var_group_labels=geneGrpName, ax=axs[2], show=False)
18 fig.tight_layout()
```

Run **Hide Input**

In [2]:

```
# Use lines starting with '#' to separate code for readability as empty lines will create new Jupyter notebook cells
# adata: an AnnData object is available
import warnings
warnings.filterwarnings('ignore')
import scanpy as sc
from matplotlib import pyplot as plt
%config InlineBackend.figure_formats = ['svg']
sc.pp.scale(adata, zero_center=False)
sc.tl.rank_genes_groups(adata, groupby='cell_type', groups=['OPC','Astrocytes'], method='t-test')
marker_genes=[LHFPL3, PCDH15, DSCAM, TNR, NLGN1,AQP4, VCAN,CUX2]
geneGrpPos = [(0,2), (3,4), (5,7)]
geneGrpName = [set1, set2, set3]
fig,axs = plt.subplots(1,3,figsize=(15,5))
sc.pl.rank_genes_groups_violin(adata, groups='OPC', n_genes=8, ax=axs[0], show=False, strip=False)
sc.pl.stack_violin(adata, marker_genes, groupby='cell_type', var_group_positions=geneGrpPos, var_group_labels=geneGrpName, row_palette='muted', rotation=90, ticklabels=True, ax=axs[1], show=False)
sc.pl.matrixplot(adata, marker_genes, groupby='cell_type', var_group_positions=geneGrpPos, var_group_labels=geneGrpName, ax=axs[2], show=False)
fig.tight_layout()
```

The interface displays three plots: a violin plot titled 'OPC vs. rest' showing expression levels for various genes across three cell sets (set1, set2, set3); a heatmap showing gene expression across different cell types; and a matrix plot showing gene expression across cell types. A red box highlights the 'Run' button, and another red box highlights the 'Write code as you are in Jupyter notebook' text.

438

439 Fig. S21 | Command line interface for user to program for advanced plotting and statistical analysis.

440 Note: In CLI the AnnData (adata) object is available by default, and it is processed as 'Description' of the
 441 dataset states (i.e.: normalized and log transformed, but no scaled etc.). Settings in 'Figure Option' tab
 442 won't apply to CLI.

31

- 443 22. VIP – Comb. & Abbr.
- 444 User can combine multiple annotations to create combinatorial annotation to group cells in various of plotting, e.g., stacked violin and dot plot. Firstly, user ‘uncheck all annotations’ and, secondly go to the annotation panel in main window to select annotation categories to be combined, e.g., diagnosis (Control and MS) combined with sex (female and male). After clicking on ‘Create’ button, all possible combinatorial names will be listed and ‘Custom_combine’ will be automatically available as an option in ‘Group by’ drop down menu of many plotting functions.
- 450 User can also rename each annotation by creating abbreviations to shorten axis labels in figures.

The screenshot shows the cellxgene interface. On the left, there's a tree view of annotations: Capbatch, Seqbatch, cell_type, diagnosis (with Control, MS, nan), sex (with female, male, nan), and stage. A red box highlights the 'diagnosis' and 'sex' sections with the text 'Select "diagnosis" and "sex" as a new combination'. On the right, a modal window titled 'Create a combinatorial annotation' has two steps: 1. Uncheck All Annotations (with a red box around it) and 2. Select the annotations of interest on the left panel of main window, then Create a new annotation named "Custom_combine" includes the following categories: Control:female; Control:male; MS:female; MS:male. Below this is another modal titled 'Create abbreviations of annotations' with a table:

cell_type	diagnosis
Astrocytes	Astro.
B cells	B
EN-L2-3	EN-L2-3
EN-L4	EN-L4
EN-L5-6	EN-L5-6
EN-MIX	EN-MIX
EN-PYR	EN-PYR
Endo cells	Endo cells
Glia-MIX	Glia-MIX
IN-PVALB	IN-PVALB
IN-SST	IN-SST
IN-SV2C	IN-SV2C
IN-VIP	IN-VIP
Microglia	Microglia
OL	OL
OPC	OPC
Phagocytes	Phagocytes
Stromal cells	Stromal cells

A red box highlights the 'Rename "Astrocytes" to "Astro."'

- 451
- 452 Fig. S22 | Comb. & Abbr. function allows user to create new annotation by combining multiple annotations and abbreviations to shorten axis labels in figures especially when custom combinatorial names are used.

455 23. VIP – Other Functions

456 There are other convenient functions available to user, such as ‘Save’ or ‘Load’ session, ‘Check All
457 Annotations’ and ‘Brush’.

458 ‘Save’ or ‘Load’ session are used to save the current cell selections and parameter settings to text file or
459 load previously saved session file in the tool for visualization.

460 ‘Check All Annotations’ is used to check all categorical selection boxes of annotations on the left panel.

461 ‘Brush’ is to display exactly these selected ranges from histograms of variables on the right panel in a
462 nice table that is not available in original cellxgene.



463

464 Fig. S23 | Other functions allow user to ‘Save’ or ‘Load’ session information, ‘Check All Annotations’ and
465 show values of brushed ranges on histograms.

466 Sup. Method 1 – Client-side Integration by a jsPanel Window (VIP)

467

468 Following section in config.sh file.

```
469
470 read -d '' insertL << EOF
471 <script src="static/jquery.min.js"></script>
472 <link href='static/jspanel/dist/jspanel.css' rel='stylesheet'>
473 <script src='static/jspanel/dist/jspanel.js'></script>
474 <script src='static/jspanel/dist/extensions/modal/jspanel.modal.js'></script>
475 <script src='static/jspanel/dist/extensions/tooltip/jspanel.tooltip.js'></script>
476 <script src='static/jspanel/dist/extensions/hint/jspanel.hint.js'></script>
477 <script src='static/jspanel/dist/extensions/layout/jspanel.layout.js'></script>
478 <script src='static/jspanel/dist/extensions/contextmenu/jspanel.contextmenu.js'></script>
479 <script src='static/jspanel/dist/extensions/dock/jspanel.dock.js'></script>
480 <script>
481     // execute JavaScript code in panel content
482     var setInnerHTML = function(elm, html) {
483         elm.innerHTML = html;
484         Array.from(elm.querySelectorAll('script')).forEach( oldScript => {
485             const newScript = document.createElement('script');
486             Array.from(oldScript.attributes)
487                 .forEach( attr => newScript.setAttribute(attr.name, attr.value) );
488             newScript.appendChild(document.createTextNode(oldScript.innerHTML));
489             oldScript.parentNode.replaceChild(newScript, oldScript);
490         });
491     }
492     var plotPanel = jsPanel.create({
493         panelSize: '190 0',
494         position: 'left-top 160 6',
495         dragit: { containment: [-10, -2000, -4000, -2000] }, // set dragging range of VIP window
496         boxShadow: 1,
497         border: "solid #D4DBDE thin",
498         contentOverflow: 'scroll scroll', // adding scrolling bars
499         headerControls: {
500             close: 'remove',
501             minimize: 'remove',
502             maximize: 'remove'
503         },
504         headerTitle: function () {return '<strong>Visualization in Plugin</strong>'},
505         contentAjax: {
506             url: 'static/interface.html',
507             done: function (panel) {
508                 setInnerHTML(panel.content, this.responseText);
509             }
510         },
511         onwindowresize: function(event, panel) {
512             var jptop = parseInt(this.currentData.top);
513             var jpleft = parseInt(this.currentData.left);
514             if (jptop<-10 || window.innerHeight-jptop<10 || window.innerWidth-jpleft<10 ||
515             jpleft+parseInt(this.currentData.width)<10) {
516                 this.reposition("left-top 160 6");
517             }
518         },
519         onunsmallified: function (panel, status) {
520             this.reposition('center-top -370 180');
521             this.resize({ width: 740, height: function() { return Math.min(480, window.innerHeight*0.6); } });
522         },
523         onsmallified: function (panel, status) {
524             this.reposition('left-top 160 6');
525             this.style.width = '190px';
526         }
527     }).smallify();
528     plotPanel.headerbar.style.background = "#D4DBDE";
529 </script>
530 EOF
531 insertL=$(sed -e 's/[&\\/]/{&\\}/g; s/|||/|/g; s/$/\\// -e '$s/\\$//` <<"$insertL"
532 sed -i "s|<div id=\"root\"></div>|$insertL\n&|" "cellxgene/client/index_template.html"
```

533 All functional VIP HTML and JavaScript code will be in “interface.html” that is independent of cellxgene
534 code base.

536 **Sup. Method 2 – Server-side Integration**

537

538 Following section in config.sh file.

539

```
540 echo '
541 from server.app.biogenInterface import route
542 @webbp.route("/biogen", methods=["POST"])
543 def biogen():
544     return route(request.data,current_app.app_config) >> cellxgene/server/app/app.py
545 .
546 .
547 .
548
549 strPath=$(python -c 'import site; print(site.getsitepackages())')"
550 strPath=${strPath//\"/}
551 strPath=${strPath//\'/}
552 strweb="${strPath}/server/common/web/static/"
553 echo $strweb
554 cp interface.html $strweb
555 cp jquery.min.js $strweb
556 cp color_map.png $strweb
557
558 cp -R DataTables $strweb
559 cp -R jspanel $strweb
560
561 cp cellxgene/server/test/decode_fbs.py $strPath/server/app/
562 cp VIPInterface.py $strPath/server/app/.
```

563

564 **Sup. Method 3 – Communication between VIP and cellxgene web GUI**

565

566 Cellxgene client utilizes [React Redux](#) that is the official React binding for [Redux](#). It lets your React
567 components read data from a Redux store, and dispatch actions to the store to update data.

568 So, this line of code is appended to the end of client/src/reducers/index.js of cellxgene source code to
569 expose the store to the browser.

570 *window.store = store;*

571

572 By doing this, Redux store holding client data and user selections are visible to VIP to access variables
573 and dispatch actions to control cellxgene user interface. For example,

574

575 1. Unselect / select a feature. GUI is refreshed automatically after dispatching.

576

```
577       window.store.dispatch({type: "categorical metadata filter deselect", metadataField: "louvain",  
578       categoryIndex: 5})  
579       window.store.dispatch({type: "categorical metadata filter select", metadataField: "louvain",  
580       categoryIndex: 5})
```

581

582

583

584 2. Get state of just finished action and synchronize gene input and cell selections from main
585 window to VIP if corresponding action was performed.

586

```
587       const unsubscribe = window.store.subscribe(() => {  
588       if (window.store.getState()["@@undoable/filterState"].prevAction) {  
589           actionType = window.store.getState()["@@undoable/filterState"].prevAction.type;  
590           if (actionType.includes("user defined gene success")) ||  
591               actionType.includes("store current cell selection as differential set")) {  
592               sync();  
593           }  
594       }  
595       });
```

596 **Sup. Method 4 – Diffxpy Integration**

597

598 This is the sample pseudocode, please see VIPInterface.py for actual implementation.

599

```
600     import scanpy as sc
601     import pandas as pd
602     import diffxpy.api as de
603     import server.app.app as app
604     # set 1 of cells as cell1; set 2 of cells as cell2
605
606     with app.get_data_adaptor() as data_adaptor:
607         X1 = data_adaptor.data.X[cell1]
608         X2 = data_adaptor.data.X[cell2]
609
610         adata = sc.AnnData(pd.concat([X1,X2]),pd.DataFrame(['grp1' for i in range(X1.shape[0])]+['grp2' for i
611             in range(X2.shape[0])],columns=['comGrp']))
612         deg = de.test.two_sample(adata,'comGrp').summary()
613         #deg is a dataframe contains the following columns ['gene','log2fc','pval','qval']
```

614 **Sup. Method 5 – Create h5ad file from Seurat object**

615 First, export the following from Seurat object in R: **expression matrix (assume normalized), metadata**
616 **and coordinates (pca, tsne, umap) as separate txt files.**

617 Next in Python, create an AnnData object from 10x (*scanpy.read_h5ad function*) as a starting point.
618 Then replace the expression matrix, meta data and coordinates as following, a h5ad file would be
619 generated.

620 *Python code:*

```
621 import sys
622 import scanpy as sc
623 import pandas as pd
624 import numpy as np
625 import seaborn as sns
626 from numpy import ndarray, unique
627 from scipy.sparse.csc import csc_matrix
628
629 adata= sc.read_h5ad("previous generated .h5ad")
630
631 # read clustering res
632 xPCA = pd.read_csv("./data/harmony_clustered.h5ad.pca_coordinates.txt", sep='\t', encoding='utf-8')
633 xTSNE = pd.read_csv("./data/harmony_clustered.h5ad.tsne_coordinates.txt", sep='\t', encoding='utf-8')
634 xUMAP = pd.read_csv("./data/harmony_clustered.h5ad.umap_coordinates.txt", sep='\t', encoding='utf-8')
635
636 xObs = pd.read_csv("./data/harmony_clustered.h5ad.meta_data.txt", sep='\t', encoding='utf-8')
637
638 xPCA.set_index('index', inplace=True)
639 xTSNE.set_index('index', inplace=True)
640 xUMAP.set_index('index', inplace=True)
641 xObs.set_index('index', inplace=True)
642
643 adata.obsm['X_pca'] = np.array(xPCA.loc[adataRaw.obs.index])
644
645 adata.obsm['X_tsne'] = np.array(xTSNE.loc[adataRaw.obs.index])
646 adata.obsm['X_umap'] = np.array(xUMAP.loc[adataRaw.obs.index])
647 adata.obs = xObs.loc[adataRaw.obs.index] # this is a pandas dataframe
648
649 # read in expression matrix as numpy.ndarray as following:
650 exp_mat = np.loadtxt(fname ="expression matrix .txt")
651 adata.X = exp_mat
652 # convert dense matrix into sparse matrix to save storage space and memory usage
653 adata.X = csc_matrix(adata.X)_matrix
654
655
656
657
```

```
658  
659  
660 # add short description and initial graph settings. “|” and “by” are delimiters for VIP to parse the  
661 initial settings. Please follow the same rule for your own h5ad files.  
662  
663 adata.obs['>Description'] = ['Human brain snRNAseq 46k cells (MS Nature 2019 Schirmer et al.); data -  
664 normalized, log transformed and scaled UMI; platform - 10X v2 chemistry | embedding by umap; color  
665 by cell_type']*adata.n_obs  
666  
667 # Then last step to save h5ad:  
668 adata.write_h5ad("final output.h5ad")  
669  
670  
671
```

672 Helpful tips

673 *Handle nulls in categorical annotation*
674 Such nan's in categorical annotation would cause trouble in VIP because it cannot be converted to
675 string. Here is how to handle it, let's call the annotation X_annotation:

```
676 # Cast to str from categorical  
677 adata.obs = adata.obs.astype({'X_annotation':'str'})  
678 # replace all of nan by 'nan'  
679 adata.obs["X_annotation "][adata.obs["X_annotation "].isnull()] = 'nan'
```

680
681 *Display full traceback stack for debugging in VIP*
682 It follows the global setting. Please set “—verbose” to launch cellxgene server.

683 *Pitfall of using special characters*
684 In the mode which allows user to create manual annotation in cellxgene, user should try to avoid using
685 hyphen (“-”) in name label. It would cause client-side issue. Please try to use underscores.

686 *Potential use for bulk or pseudo bulk sample dataset*
687 Once the data matrix is replaced by sample x gene matrix, cellxgene VIP framework can handle regular
688 bulk / pseudobulk RNAseq datasets. Simply replace “cells” by “samples”. All plotting functions would still
689 work.