

# 1 Geometry aware graph attention networks to 2 explain single-cell chromatin state and gene 3 expression

4

5 Gabriele Malagoli<sup>1,2</sup>, Patrick Hanel<sup>1</sup>, Anna Danese<sup>2</sup>, Guy Wolf<sup>3,4</sup>, Maria  
6 Colomé-Tatché<sup>1,2</sup>

7

8 <sup>1</sup> Helmholtz Zentrum München, Institute of Computational Biology, Neuherberg, Bavaria, Germany

9 <sup>2</sup> LMU Munich, Biomedical Center (BMC), Physiological Chemistry, Faculty of Medicine,

10 Planegg-Martinsried, Bavaria, Germany

11 <sup>3</sup> Université de Montréal, Department of Mathematics and Statistics, Montréal, Québec, Canada

12 <sup>4</sup> Mila-Quebec Artificial Intelligence Institute, Montréal, Québec, Canada

13

## 14 Abstract

15 High-throughput measurements that profile the transcriptome or the epigenome of  
16 single-cells are becoming a common way to study cell identity. These data are high  
17 dimensional, sparse and non linear. Here we present SEAGALL (Single-cell  
18 Explainable Geometry-Aware Graph Attention Learning pipeLine), a hypothesis free  
19 method to extract biologically relevant features from single-cell experiments based  
20 on geometry regularised autoencoders (GRAE) and explainable graph attention  
21 networks (GAT). We use a GRAE to embed the data into a latent space preserving  
22 the data geometry and we construct a cell-to-cell graph computing distances in the  
23 GRAE bottleneck. Exploiting the attention mechanism to dynamically learn the  
24 relevant edges, we use GATs to classify the cells and we explain the predictions of  
25 the model with XAI methods to unravel the features which are driving cell identity  
26 beyond marker genes. We apply our method to data sets from scRNA-seq,  
27 scATAC-seq and scChIP-seq experiments. SEAGALL can extract cell type specific  
28 and stable signatures which not only differ from the ones found in classical linear  
29 approaches but are less biased by coverage and high expression.

30

31

32

33

34

35

36

## 37 Introduction

38 Single-cell sequencing technologies have provided a breakthrough in molecular  
39 biology by allowing the measurement of transcriptomic and epigenomic profiles at  
40 the scale of single-cell with high resolution. Many reproducible and ready-to-apply  
41 kits have become common and affordable, leading to a great increase of interest in  
42 this field. For instance, single-cell RNA sequencing (scRNA-seq), also known as  
43 gene expression (GEX), or single-cell Assay for Transposase-Accessible Chromatin  
44 using sequencing (scATAC-seq)<sup>1,2</sup> can be performed with the readily available kits  
45 commercialized by 10X Genomics by applying their well described protocols. A new  
46 step towards the understanding of molecular biology is the possibility to do  
47 multi-omics single-cell sequencing, which allows the measurement of more than one  
48 omics modality at the same time for the same single-cell. Among others, the 10X  
49 Genomics Multiome Platform, which quantifies the chromatin openness and the  
50 transcriptome of single nuclei, also allows to perform such measurements.

51

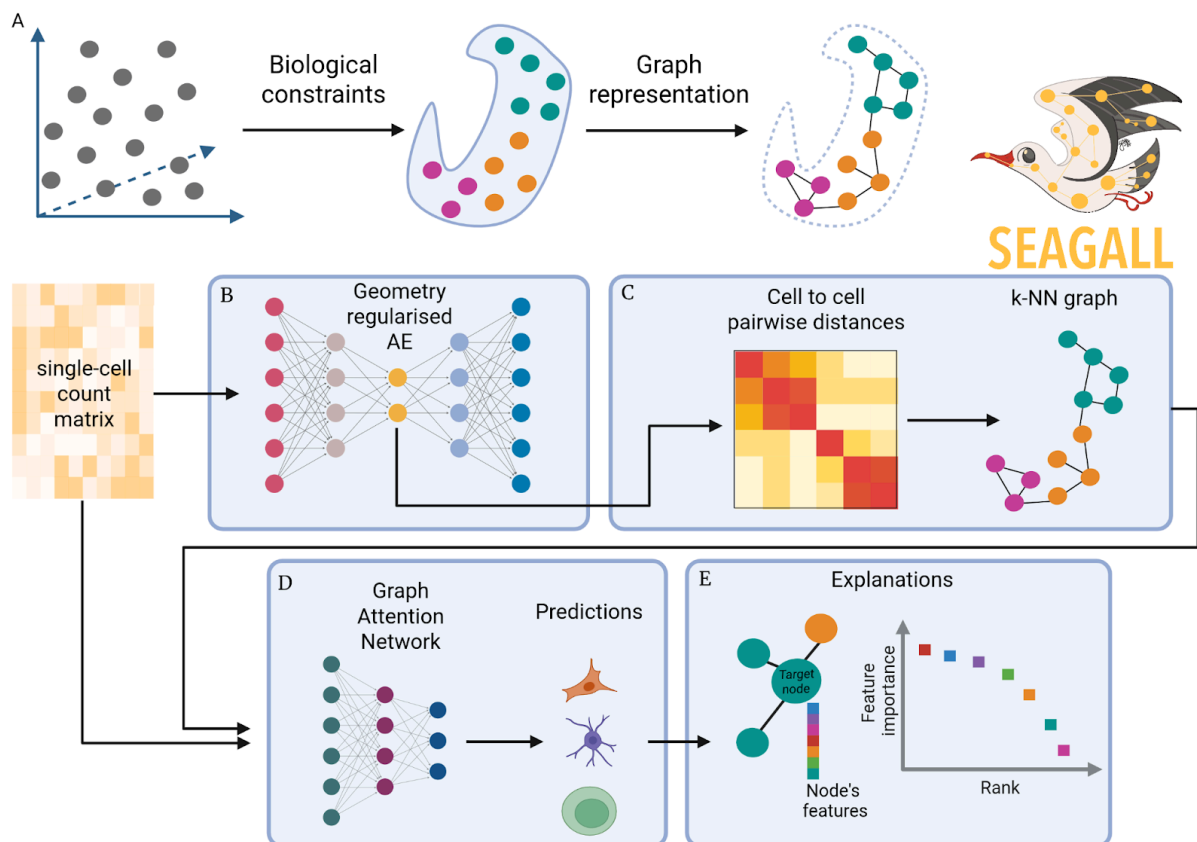
52 The standard analysis of single-cell data involves low-dimensional embedding,  
53 followed by cell clustering and identification of cell types<sup>3</sup>. The common assumption  
54 is known as the “manifold hypothesis”<sup>4</sup>: high dimensional data lie along a latent and  
55 unknown manifold with a smaller dimension than the observed space. In single-cell  
56 biology, we measure tens of thousands of variables, such as genes (for gene  
57 expression measurements) or genomic loci (for epigenomic measurements). These  
58 features cannot take any possible value, but rather they vary within well defined  
59 ranges given by biological constraints, like gene regulatory networks<sup>5</sup>. These  
60 constraints define the underpinning manifold whose exact equations are unknown. A  
61 single-cell experiment can be seen as a method to sample (cells) from this manifold.  
62 From the distances between cells we can create a graph that resembles the manifold  
63 as accurately as possible (Fig.1A).

64 Many tools exist to perform single-cell analysis<sup>3</sup>, yet they show several limitations.  
65 First of all, they are in general omic-specific<sup>6-11</sup>, relying on omic-specific  
66 assumptions, and forcing the users to choose a different tool for each omic.  
67 Moreover, most standard tools compute distances in a low dimensional linear space  
68 such as Principal Components (pca) or Independent Components (ICA)<sup>6,8,9,12-14</sup>.  
69 These linear assumptions lead to the loss of the intrinsic nonlinearity present in  
70 biological data sets, and prevent the discovery of complex insights between features  
71 and cells. Due to these shortcomings, autoencoders (AE)<sup>15</sup> have recently become  
72 very popular because of their ability to learn the input and embed it in a nonlinear  
73 fashion<sup>10,11,16-18</sup>. Indeed, their strongest characteristic is the ability to take into  
74 account nonlinear dependencies within the data sets without making strong data  
75 assumptions. Yet, autoencoders often fail to represent the intrinsic data structure<sup>19</sup>,  
76 such as topology or geometry. To better fit the single-cell data, omic-specific AE have  
77 been developed; they assume a probability distribution from which the data are  
78 sampled and they use variational AE to embed the data sets<sup>10,11,16,18</sup>. As a

79 consequence, a specific AE is needed for each modality, which grows in number day  
80 by day.

81 Finally, another important aspect of single-cell data analysis is the identification of  
82 features defining cell identity. The standard method to investigate what are the  
83 important features, such as genes or peaks, for each group of cells is based on  
84 differential analysis (DA). It consists in computing the distributions of the features in  
85 the different treatments, conditions or cell types to then quantify the difference  
86 between these distributions, giving as a result a list of features ranked by the most  
87 different to least ones. This approach will output features which are different between  
88 two groups of cells, but there is no guarantee that they are also relevant and  
89 important for each group specifically.

90 To address these limitations, we developed SEAGALL (Single-cell ExplAinable  
91 Geometry-Aware Graph Attention Learning piLine), a deep learning method based  
92 on manifold learning and explainable AI for downstream analysis of single-cell data  
93 sets. SEAGALL first learns a low-dimensional embedding of the cells based on a  
94 graph-regularised autoencoder (GRAE)<sup>19</sup> (Fig.1B). This embedding preserves both  
95



96

97 Figure 1: **A** without any constraint, the measurement of molecular features can take any arbitrary value (left). In  
98 reality, the gene regulation network imposes constraints which define the cell type and the possible values of the  
99 variables, defining a manifold on which the data live (center). As a proxy for the manifold it is possible to use a  
100 cell-to-cell graph, defined by pairwise distance between cells (right). **B - E** the SEAGALL model. The initial count  
101 matrix is reduced with a GRAE to preserve data geometry, i.e. both local and global structure of the data (**B**),  
102 within the latent space we compute pairwise distance between cells to build a cell-cell graph, (**C**). The graph and  
103 the count matrix are subsequently the input to a GAT classifier, whose predictions (**D**) will be explained in order to  
104 identify the most relevant features for each cell type (**E**).

105 local and global structure of the data without particular assumptions on the data  
106 generation process. Then the tool computes the cell-to-cell graph on that  
107 low-dimensional space (Fig.1C), which is used as input to a graph attention network  
108 (GAT)<sup>20,21</sup> together with the count matrix defining feature vectors of the nodes. The  
109 GAT classifies the cells into different cell types or states (Fig.1D) and the final output  
110 of SEAGALL are the explanations of the model, i.e. the set of input features which  
111 are the most important for the label prediction<sup>22</sup> (Fig.1E). We applied our new  
112 method to ten different single-cell data sets spanning three omics (sc-RNAseq,  
113 sc-ATACseq, sc-ChIPseq) showing that it is able to reconstruct and embed the data,  
114 explain the cell types beyond common marker genes and extracting stable and  
115 specific features important for the cells which are not seen by standard differential  
116 analysis.

## 117 Results

### 118 The SEAGALL model

119 In a single-cell sequencing experiment, a set of genomic variables are measured for  
120 every cell, for instance gene expression in scRNA-seq, or openness of genomic loci  
121 in scATAC-seq. These measurements can be represented as a count matrix, where  
122 the set of specific variables are quantified within each cell. For example, in the  
123 measurement of the expression of  $F$  genes in  $N$  cells, the output is a  $N \times F$  count  
124 matrix. The same holds for sc-ATACseq, where  $F$  is the number of loci in the  
125 genome for which the openness is quantified, called “peaks”. The count matrix can  
126 be seen as a method to store the position of  $N$  cells in an  $F$ -dimensional space,  
127 commonly called point cloud. Without constraints on the expression of genes or to  
128 openness of chromatin in the cells, the point cloud may span a homogenous volume  
129 in space. Yet constraints do exist, imposed for example by gene regulatory networks;  
130 therefore the data do not occupy a homogenous volume, but rather live on a  
131 manifold<sup>23</sup>, whose equations are unknown. However, the manifold is high  
132 dimensional, making little sense to compute distances on it due to the curse of  
133 dimensionality. Hence, the first step of SEAGALL is learning a low dimensional  
134 representation of the data that conserves the intrinsic geometry of the manifold,  
135 exploiting the recent development of geometry regularised autoencoders (GRAE)<sup>19</sup>  
136 (Methods) (Fig.1B). GRAE first applies a kernel method named PHATE<sup>24</sup> to learn the  
137 geometry of the data and uses it to regularise the structure of its latent space. Within  
138 the latent space of the GRAE, it is now possible to compute reliable pairwise  
139 distances between cells in order to create a cell-to-cell k-NN graph (Fig.1C). In the  
140 next step of SEAGALL, the cell-to-cell graph is used as input to a graph attention  
141 network<sup>21</sup> (GAT) (Fig.1D), a graph neural network<sup>25</sup> (GNNs) with an attention  
142 mechanism on the edges. The GAT is applied to classify the cells into multiple cell  
143 types, which are known based on previous cell type annotation. In this scenario the  
144 classification of a cell depends on its neighbourhood, via the joint embedding of  $k$

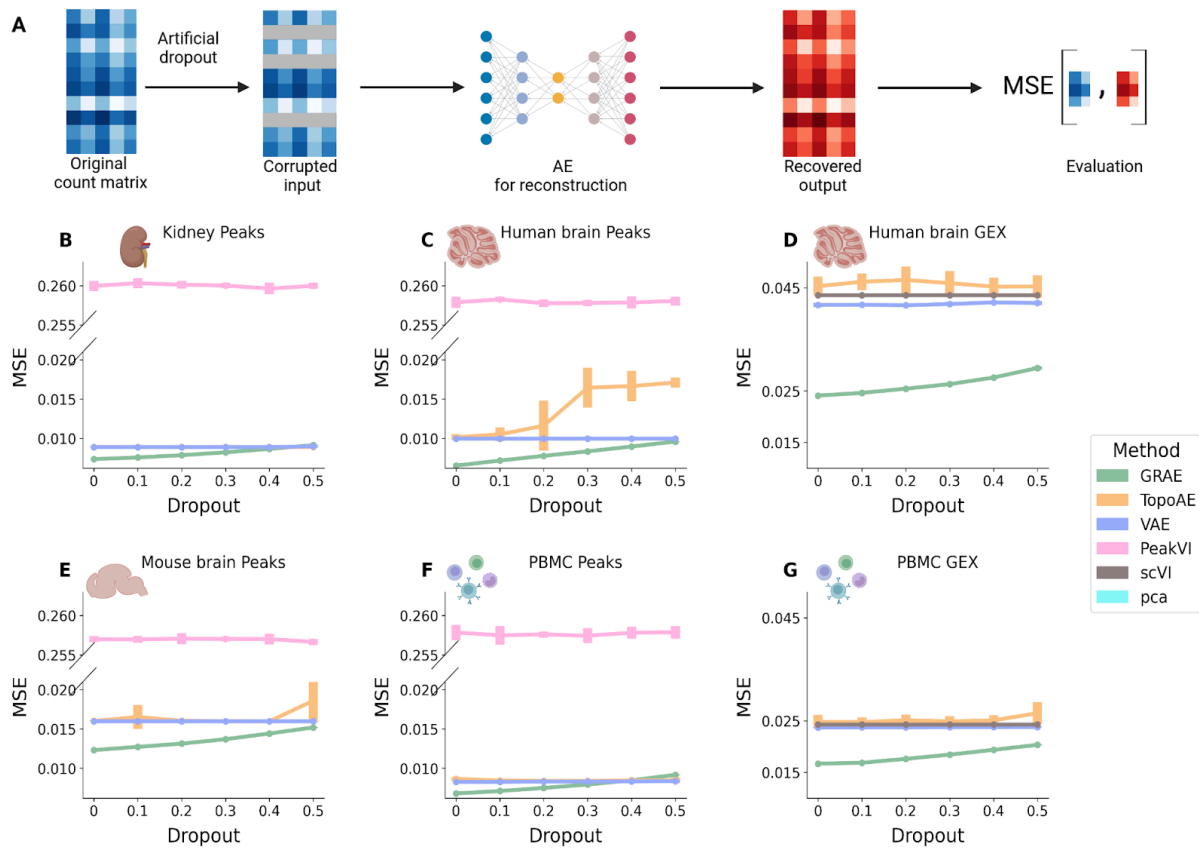
145 features vectors, if the cell has degree  $k$  (see Methods). The attention mechanism is  
146 important to dynamically learn the relevance of each edge: spurious edges will be  
147 ignored so the model can focus on the important ones. Finally, an explainable  
148 artificial intelligence method (GNNExplainer)<sup>22</sup> (Fig.1E) is applied to explain the  
149 predictions: each cell is classified into a cell type and we can extract the most  
150 relevant features that the model has used to predict the label. This last step is critical  
151 to understand what are the most relevant genomic regions or genes that define cell  
152 types beyond the common markers.

## 153 Model benchmarking

154 We carried out a breakdown of SEAGALL, testing each of its main parts: the  
155 embedding method (GRAE), the classifier (GAT) and the explainer (GNNExplainer).  
156 We used six count matrices for benchmarking the embedding method and the  
157 classifier, two from scRNA-seq and four from scATAC-seq (SuppTable1 and  
158 SuppTable2 for the cell type composition and dimensions). They are two multimodal  
159 data sets for which the scRNA-seq and the scATAC-seq were treated separately  
160 (human brain and human PBMC), the scATAC-seq part of a multimodal data set of  
161 mouse embryonic brain, and a scATAC-seq data set of kidney<sup>26</sup> (see Methods for the  
162 count matrix construction and processing). We tested the ability of six embedding  
163 methods (see next paragraph) to recover the original data after adding artificial  
164 dropout and the quality of the cell-to-cell graph computed in the different latent  
165 spaces. To quantify the latter feature, we measure the homogeneity of the cell-to-cell  
166 graph in terms of cell type composition of the neighbourhood and the performance of  
167 a GNN classifier varying the input graph. We then tested the GAT and also a Graph  
168 Convolutional Network (GCN) architecture, computing F1 score accuracy, precision  
169 and recall of the classifiers. Last, we measured the stability and the specificity of the  
170 GNNExplainer. We also measured the classification and explanations performances  
171 of the final model on a scChIP-seq dataset of breast cancer (SuppTable1,  
172 SuppTable2), in which H3K27me3 was measured at the single-cell level<sup>27</sup>.

## 173 Geometrical regularised autoencoders best recover corrupted 174 data and capture biological structure

175 We benchmarked five different AEs architectures. We tested the GRAE together with  
176 a topological autoencoder (TopoAE)<sup>28</sup>, PeakVI<sup>11</sup>, scVI<sup>10</sup>, a standard variational  
177 autoencoder (VAE) and linear pca. The TopoAE was included to compare the GRAE  
178 to an AE with a similar rationale behind: while the GRAE regularises the loss function

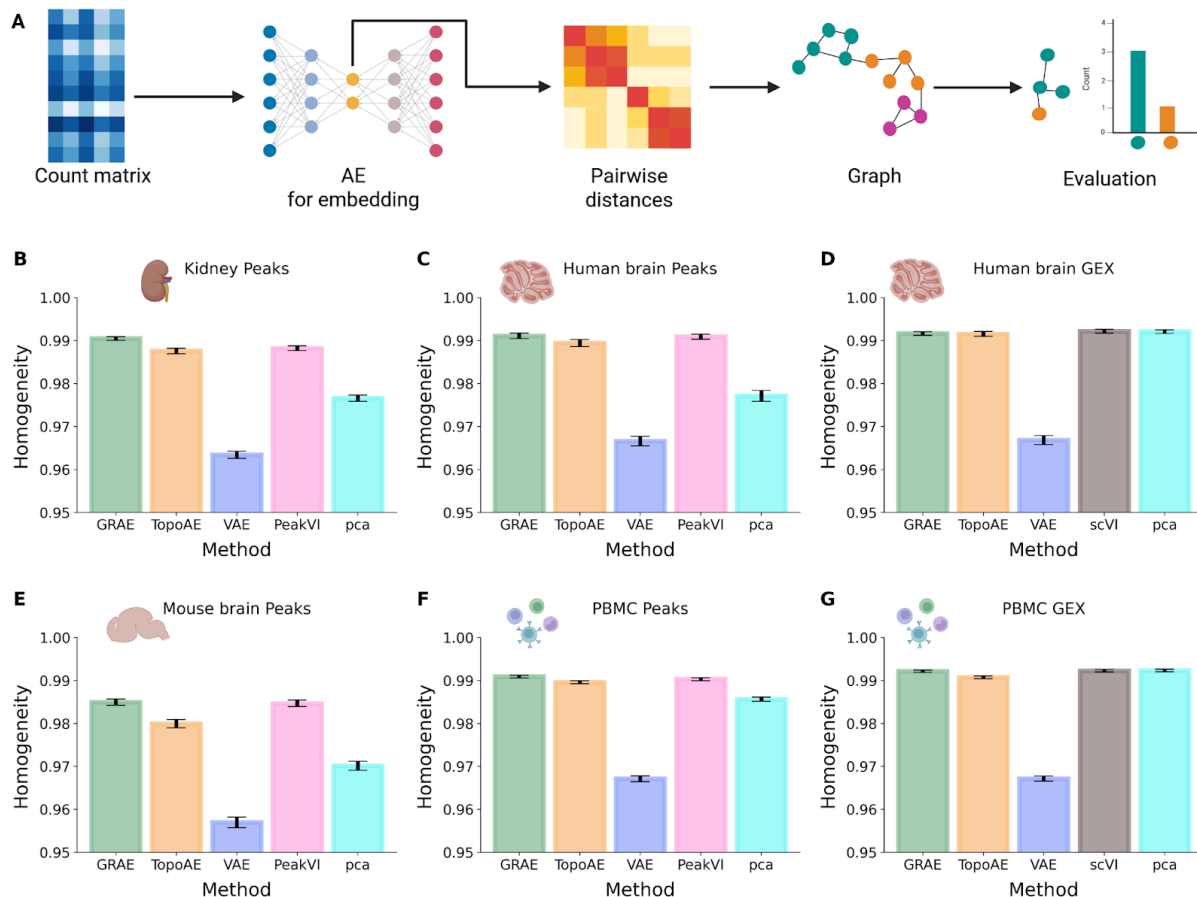


179

180 Figure 2: **A** schematic of the benchmarking process: data is corrupted with artificial dropout and the AEs  
 181 reconstruct them. We evaluated the performance computing the MSE between the recovered output and the  
 182 original count matrix. **B-G** average MSE at different levels of artificial dropout for each AE. Each point is the  
 183 average of ten runs and the height of the error bar represents three times the uncertainty on the mean.

184

185 considering that the geometry of the data should be preserved in the latent space,  
 186 the TopoAE preserves the topology of the input space by applying persistent  
 187 homology. Geometry is a more specific and local property than topology; however  
 188 neither GRAE nor TopoAE make assumptions about the data sets, which makes  
 189 them applicable to, in principle, any kind of biological data. The VAE is taken as a  
 190 baseline model of autoencoder, to compare sophisticated methods to a simpler one.  
 191 PeakVI is a state-of-the-art scATAC-seq specific AE, tailored for this data type so it  
 192 can only be applied to it. scVI is a single-cell AE which represents the state of the art  
 193 to embed scRNA-seq data, it therefore can only be applied to this data modality. pca  
 194 is included in the benchmarking to test how different a linear dimensionality reduction  
 195 method performs. We quantified which AE architecture is better to use for  
 196 dimensionality reduction of single-cell data to then construct the cell-to-cell graph.



197

198 Figure 3: **A** schematic of the benchmarking process: each count matrix is embedded using the different AEs, the  
 199 cell-to-cell graph is computed from the latent space and its homogeneity is used to evaluate the performance of  
 200 the AE. **B-G** homogeneity of the k-NN using the different embedding methods. The height of the bar represents  
 201 the average homogeneity across runs and the error bars spread is three times the uncertainty on the mean.  
 202

203 To measure the AEs ability to retrieve corrupted data, we applied a variable dropout  
 204 between 0% and 50% in regular incremental steps of 10% to the six RNA and ATAC  
 205 count matrices and trained each AE on the faulty data. We repeated the experiment  
 206 ten times to ensure that each time the dropout will affect different features and to  
 207 have a statistically reasonable sample size. After the training we measured the mean  
 208 squared error between the decoded matrices of each AE and the original,  
 209 non-corrupted, matrices (Fig.2A, Methods). GRAE outperforms all the methods  
 210 achieving the minimum MSE at every dropout level (Fig.2B-G), except for the highest  
 211 dropout on the peaks of the PBMC data. In particular, the geometry regularised AE is  
 212 better than the two -omic specific AE scVI and PeakVI.

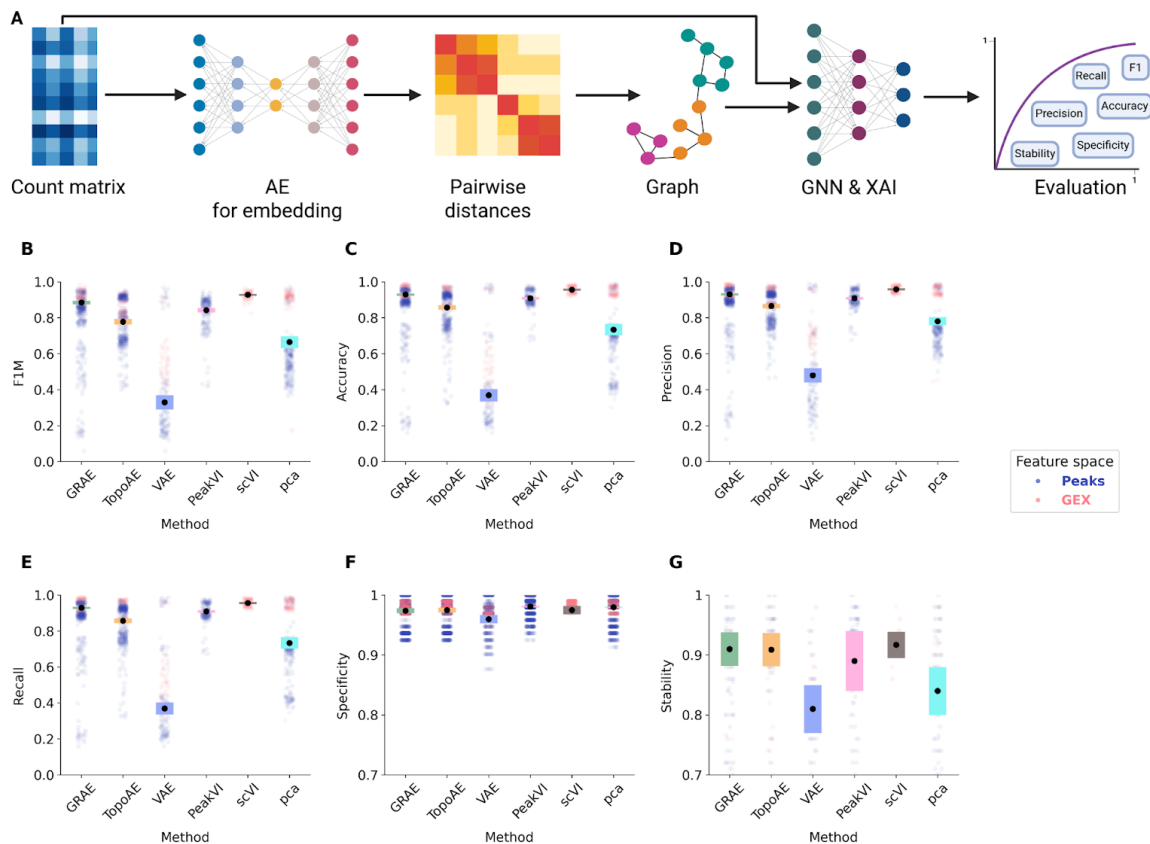
213 Then we measured the homogeneity of the k-NN graphs built from the AEs latent  
 214 spaces. We assume that a good latent space leads to a k-NN graph where  
 215 neighbours of a node belong to the same cell type. The more homogeneous the  
 216 neighbourhood, the more the AE is able to locate close to each other in the latent  
 217 space cells sharing the same biological functions. We trained each AE and we  
 218 computed the k-NN (k=15) graphs from their latent spaces. Then, we measured for  
 219 each cell how many cell types are found in its neighbourhood. We normalised this  
 220 value by the product of the number of the neighbours and the number of cell types

221 within the data set. We call this score “heterogeneity”. The value reported for each  
222 AE and each count matrix is then the average “homogeneity”, which is computed as  
223 one minus the average heterogeneity for each neighbourhood (Fig.3A, Methods).  
224 For the ATAC data sets, GRAE outperforms all the methods on the kidney data set  
225 and on the ATAC part of the PBMC data set (Fig.3B,F); on the other two ATAC count  
226 matrices, PeakVI and GRAE achieve the same homogeneity (Fig.3C,E). The  
227 homogeneity of the graphs computed from GEX count matrices is similar between  
228 GRAE, scVI and pca (Fig.3D,G, SuppTable4).  
229 In conclusion, GRAE, which applies a geometrical regularisation to the loss function,  
230 outperforms all other methods when it comes to reconstructing the initial input, even  
231 with the addition of noise in the data. It also performs either better or equal to the  
232 other methods in recovering the cell type composition in the latent space.

## 233 Geometry aware graph attention networks achieve best 234 classification performances

235

236 Last, we tested together the performances of different embedding strategies and  
237 GNN classifiers. We used the aforementioned AEs and pca as dimensionality  
238 reduction (DR) methods, combined with two types of GNN: GAT<sup>20</sup> and GCN<sup>29</sup>. Each  
239 combination of the DR method and GNN with the explainer was run fifty times,  
240 changing the initial seed, to statistically test the stability of the results. For each  
241 combination and each run we applied six metrics: four metrics for the classification  
242 performance (accuracy, F1 score, precision and recall) and two metrics to measure  
243 the quality of the explanations (specificity and stability) (Fig.4A). Specificity quantifies  
244 how specific the explanations are for each cell type. It is defined as one minus the  
245 average overlap between the most important features for each cell type with the  
246 most important for each other cell type. Stability measures how much the  
247 explanations change by running a new instance of the same classifier and explainer  
248 for the same count matrix. It is computed as the average intersection between the  
249 extracted features of a cell type across different repetitions of the same classification  
250 and explanations (Methods). We did not measure any statistical differences between  
251 the performances of GAT and GCN in terms of F1, accuracy, precision, recall,  
252 specificity and stability (SuppFig1, SuppTable4). Meanwhile, the GRAE graph  
253 outperformed that of PeakVI, TopoAE, VAE and pca in terms of F1 score, accuracy,  
254 precision and recall (Fig.4B,C,D,E, SuppTable5). The scVI graph outperformed the  
255 one from GRAE on most of the metrics; however scVI is sc-RNAseq specific,  
256 preventing the possibility to apply it to every single-cell data The specificity and the  
257 stability of the explanation are very high with all the embedding methods, with GRAE  
258 either leading or being the second in the rank (Fig4F, G, SuppTable5). We tested the  
259 final combination of GRAE and GAT on the scChIP-seq data set, and showed a very  
260 high performance also for that data type in terms of accuracy, F1, precision and  
261 recall, as well as specificity and stability of the discovered features (SuppFig2).



262

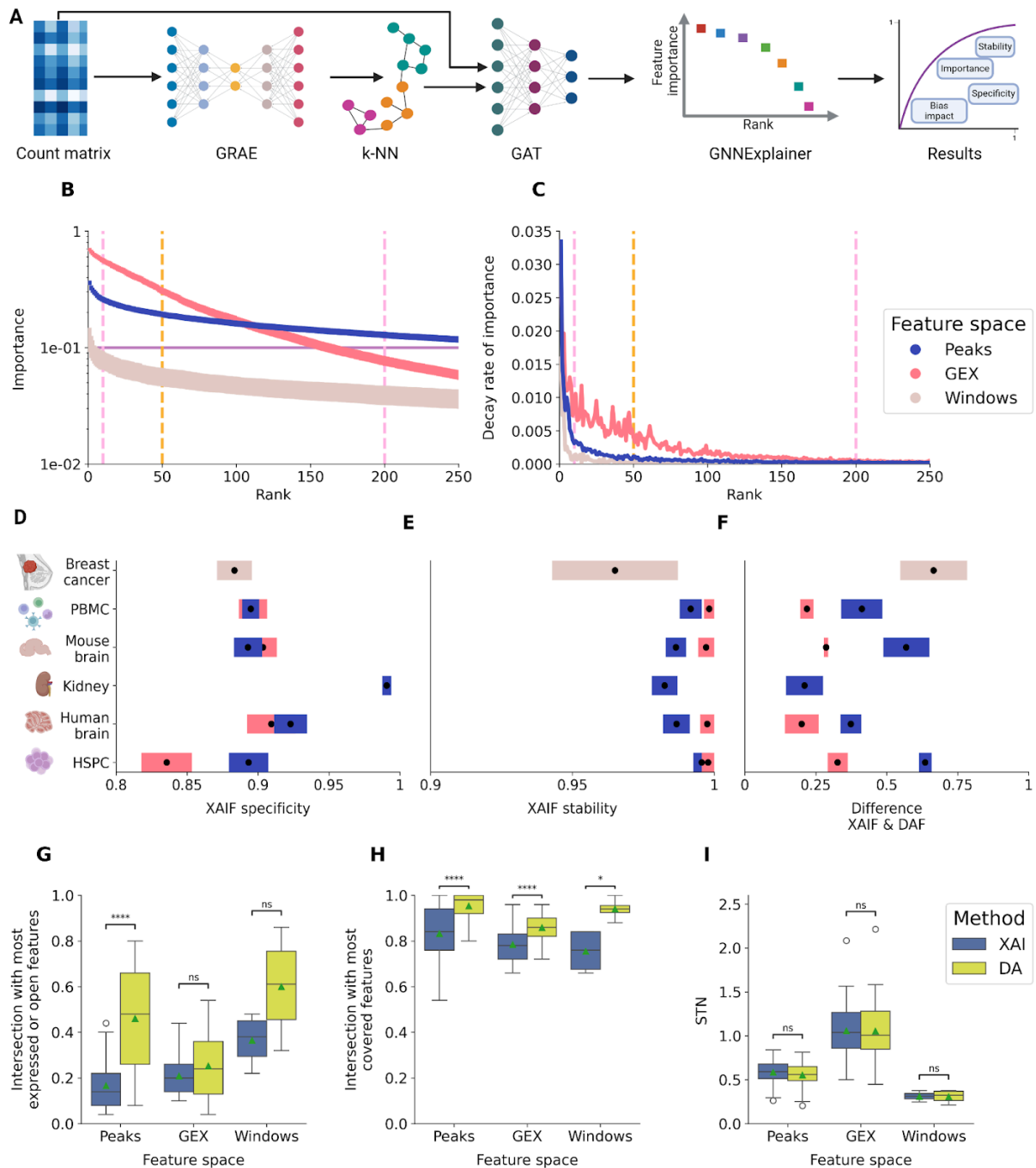
263 Figure 4: **A** schematic of the benchmarking process: After computing the k-NN graph, we trained a GNN classifier  
 264 and applied the GNNExplainer; we finally computed the shown metrics to evaluate the models. **B - E**  
 265 classification performances of the GAT classifier varying the embedding methods. Each black dot represents the  
 266 mean across 50 runs and the height of the bars represent three times the standard deviation of the mean. **F - G**  
 267 specificity and stability of the explanations. The vertical axes start from 0.7 for visualisation purposes. In all  
 268 panels, each pale plot represents the point contributing to the mean, colored by data modality. PeakVI is only run  
 269 on peak data sets and scVI is only run on GEX data sets; all other methods are run on all data sets.

270

271 In conclusion, these results indicate that the GRAE combined with a GNN classifier  
 272 outperforms all other methods except of scVI for classifying cells into cell types and  
 273 have the most stable and specific explanations.

## 274 SEAGALL retrieves stable, specific and unbiased features

275 Given these results, the final SEAGALL model consists of the GRAE to embed the  
 276 data and build the graph, and the GAT to classify the cells (Fig.5A). However, the  
 277 last and crucial step is the explanation of the predictions. This point is crucial since it  
 278 moves the focus from prediction performances to model interpretability, making the  
 279 tool translational and useful for providing new biological insights. Once the GAT is  
 280 trained on the geometry aware graph, SEAGALL investigates what are the features  
 281 which are driving the predictions of the model, assuming that these features are the  
 282 most relevant for the cell type and can define it beyond most common marker genes.  
 283 To address this point, it applies a mask-based graph neural network explainer,  
 284 known as GNNExplainer<sup>22</sup>. Given a node  $v_i$  the explainer finds the subgraph  $G_s$  and  
 285 the subset of features  $X_s = \{x_j | v_j \in G_s\}$  that maximises the probability of having



286

287 Figure 5: **A** final workflow of SEAGALL. **B** rank-importance distribution of the features according to the explainer.

288 Average across data sets and cell type. Vertical dashed lines highlight the interval 10-200 features within the

289 importance is stable. **C** Decay rate of the importance of the features, legend as in A. **D** - **F** specificity (left),

290 stability (center) and difference between xai features (XAIF) and differential features (DAF) (right) of the

291 explanations on the ten count matrices, spanning three different feature spaces and using 50 features for each

292 label, legend as in A. **G** distribution of the overlap between XAIF and most expressed or open features (blue) and

293 overlap between DAF and most expressed or open features (yellow). **H** distribution of the overlap between XAIF

294 and most covered features (blue) and overlap between DAF and most covered features (yellow). **I** distribution of

295 the signal-to-noise (STN) ratio of XAIF and DAF.

296

297 the observed prediction  $\hat{y}_i = \Phi(G_s, X_s)$  where  $\Phi$  is the trained GNN. In other

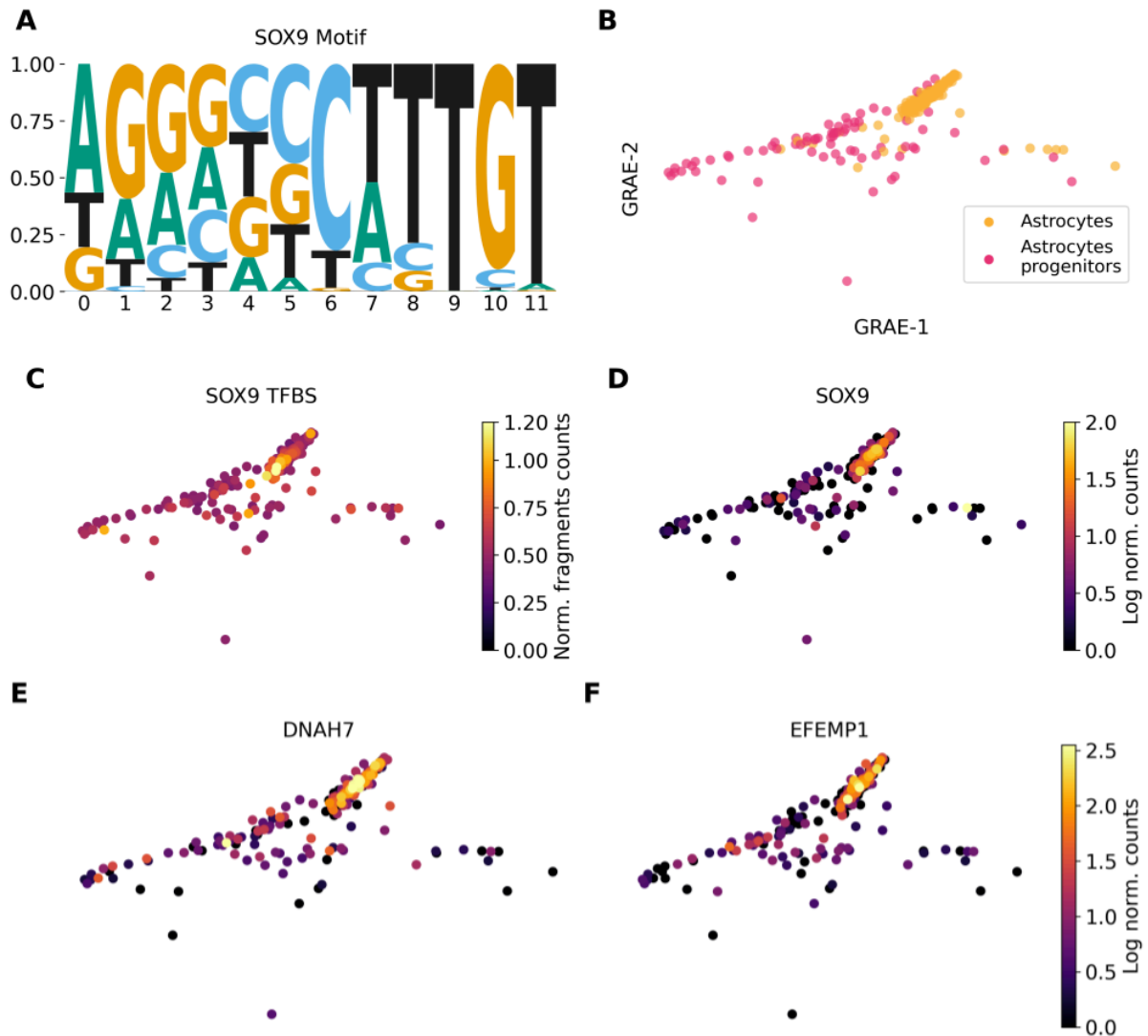
298 words, the explainer finds the subset of nodes features (and a subset of nodes links)

299 that are most important to predict the nodes label. The importance is defined as the

300 mutual information between the feature and the predictions (Methods). The  
301 distribution of the features importance drops fast with the rank, especially in GEX  
302 data. For both genes and peaks, at around the two hundredth feature, the  
303 importance drops one order of magnitude (Fig.5B, SuppFig3). Therefore, we suggest  
304 keeping a lower number of features for downstream analysis. For the single  
305 scChIP-seq data set, window features show a different behaviour: the maximum  
306 importance is significantly smaller than in peaks and GEX, and the importance decay  
307 is slower (Fig.5B). We speculate that because of the more noisy and sparse nature  
308 of the data each individual feature has a lower impact on the final prediction.  
309 However, the rate decay of the importance, computed as the absolute value of the  
310 derivative of the importance by the rank, shows that for all the three feature spaces  
311 the importance of the features does not change anymore after rank two hundred.  
312 Windows have a slower decay, again suggesting that each individual window has a  
313 lower impact on the results (Fig.5C). We picked fifty features for the downstream  
314 analysis, where we quantify the impact of technical biases on the DA and XAI.  
315 However, the influence of this threshold on the following results is very limited, within  
316 the range of 10-200 features (SuppFig4). We measured the stability and specificity of  
317 the XAI features (XAIFs). Stability is defined as the average overlap between the  
318 explanation of the same cell type running a new instance of the model and the  
319 explainer, therefore it measures the robustness of the method. Specificity is defined  
320 as the average overlap of explanations between cell types, quantifying the ability of  
321 the method to retrieve different features for different cell types. Our method is able to  
322 extract cell type specific (Fig.5D) and highly stable (Fig.5E) features in all the ten  
323 count matrices we tested. This means that it can consistently understand and explain  
324 a data set in order to suggest the key degrees of freedom which can be studied in  
325 downstream analysis and wet lab experiments, linking the deep learning method to  
326 real-world information. Notably, the XAIFs consistently differ from the differential  
327 features (DAFs) (Fig.5F), providing a potential for the discovery of novel data  
328 characteristics. This is due to the non-linearity and awareness of geometry of the  
329 model we propose, as replacing either GRAE with pca or GAT with standard NN  
330 leads to different explanations (SuppFig5). A direct comparison between XAIF and  
331 DAF shows that the former are less biased by high openness or expression and  
332 coverage (Fig.5G,H). This is particularly strong with peaks. Nevertheless, the lower  
333 biases of XAIFs are not traded off with a higher noise: the signal to noise (STN) ratio  
334 of XAIFs and DAFs is, on average, the same in all the feature spaces (Fig.5I).

## 335 SEAGALL identifies chromatin priming states and known cell 336 type predictors

337 To study the biological significance of our results, we studied the features which were  
338 identified only by SEAGALL and not by differential analysis. For the scATAC-seq  
339 feature spaces, to link genomic loci to transcription factors (TFs), we run a motif  
340 analysis on the top important features, using HOMER<sup>30</sup>.



341

342 Figure 6: **A** SOX9 motif. **B** GRAE embedding showing the differentiation from astrocytes progenitors to  
343 astrocytes. **C** openness of SOX9 TFBSs. **D-E** SOX9, DNAH7 and EFEMP1 expression.

344 HOMER takes as input a set of genomic intervals and identifies enriched motifs, i.e.  
345 recurrent patterns of bases, and it checks whether these patterns match known  
346 motifs of TF binding.

347 In the human brain data set we explored both the GEX and ATAC modalities. Taking  
348 the scATAC-seq XAIF and running motif analysis, we identified several brain specific  
349 motifs, which were not retrieved by motif analysis on the DA specific features  
350 (ExtendedTable1 for the complete motif results). In the astrocytes progenitors,  
351 SEAGALL could identify motifs belonging to the well known family of TFs SOX, such  
352 as SOX9 (Fig.6A), SOX17 (SuppFig6A) and SOX1 (SuppFig6B). SOX9 is known to  
353 be essential for the correct development of astrocytes<sup>31</sup> and its promoter is activated  
354 to determine astrocyte differentiation<sup>32</sup>. Notably, the two dimensional embedding  
355 obtained with GRAE can well capture the differentiation process from astrocytes  
356 progenitors to astrocytes along its horizontal axis (Fig.6B). We measured the  
357 openness of all SOX9 TFBS (obtained from HOMER) and it turns out that they are

358 already open in the astrocyte progenitors with a maximal openness in astrocytes  
359 (Fig.6C). On the other hand, the scRNA-seq modality shows that the expression of  
360 SOX9 is very limited in the progenitors but very high in the mature cells (Fig.6D). The  
361 other motif we retrieved is SOX17 (SuppFig6A), which is a TF known to be  
362 upregulated in astrocytes<sup>33</sup>. We discovered the openness of its TFBS as relevant for  
363 the identity of astrocyte progenitors; hence, we found a relevant TFBS openness in a  
364 progenitor population which is related to the expression of the TF in the direct next  
365 cellular state. For both SOX9 and SOX17 we therefore see the relevance of  
366 chromatin state priming the gene expression in progenitor cells, as suggested in<sup>34</sup>.  
367 Standard differential analysis could not highlight this dynamic behaviour. Focusing  
368 on GEX, SEAGALL ranked in the top fifty features of astrocytes the genes DNAH7  
369 and EFEMP1, which were not identified with standard differential analysis. The  
370 former is known to be expressed in intermediate astrocytes<sup>35</sup> and the latter is known  
371 to be expressed during synaptic development of astrocytes from iPSCs<sup>36</sup>. We  
372 correctly identified these genes during their positive gradient expression from the  
373 astrocytes progenitors to the astrocytes (Fig.6E,F). In addition, reprogrammed  
374 astrocytes have been shown to express a SOX1 positive state with neuronal stem  
375 cells characteristics<sup>37</sup> and we identified its motifs (SuppFig6B) within the XAI  
376 features. Also in the human brain data set, only SEAGALL was able to obtain the  
377 motifs of JUNb, FOSL2 and FOS (SuppFig6C,D,E) as enriched among the  
378 discovered XAIF for microglia in the ATAC modality. These TFs are known lineage  
379 determining for microglia<sup>38</sup>. For GEX in microglia only our method highlighted two  
380 important genes, TLR2 and RIPK2 (SuppFig7B,C): the former modulates microglial  
381 activity<sup>39</sup> and the latter plays an essential role in the microglia inflammatory  
382 response<sup>40</sup>. Last, in the brain cells annotated as inhibitory neurons, we found the  
383 motif of ASCL1 (SuppFig6F), which is known to specify and promote differentiation  
384 of GABAergic interneurons (i.e. inhibitory neurons)<sup>41</sup>.

385 For the PBMC data set, we also analysed the ATAC and GEX modalities. Focusing  
386 on the plasmacytoid dendritic cells (pDCs), which are cells responsible for presenting  
387 antigens to other immune cells, SEAGALL uniquely identified the RUNX1 and  
388 RUNX2 binding motifs (SuppFig6G,H), which are known TFs determining the pDCs  
389 lineage<sup>42</sup>. In the same cell type but from GEX data, the TF SOX4 is ranked in the  
390 most important features but not in the most differential ones (SuppFig7E), and it is  
391 known to be involved in pDCs ontogeny<sup>43</sup>. In addition, we exclusively found CR1  
392 (SuppFig7F) in the explanation of memory B cells, which is known to be necessary  
393 for the correct development of this cell type<sup>44</sup>. In the natural killers (NK) and T MAIT  
394 cells we uniquely retrieved, respectively, LAIR2 and CD8 (SuppFig7G,H), which are  
395 their cell type markers<sup>45,46</sup>.

396 Combining the features discovered by SEAGALL with motif analysis and manual  
397 inspection, we show how SEAGALL can identify several relevant TFBS and genes  
398 which are known to be determinants of cell types and their lineages. These TF motifs  
399 and genes were not discovered by the classical differential analysis pipeline,  
400 showing that our method is able to extract meaningful biological insights which can  
401 contribute to the discovery of new determinants of cell identity. In particular, we

402 identified several features, which were not identified using differential analysis, which  
403 related to the development and differentiation of cells, suggesting the ability of  
404 SEAGALL to capture features which are important in a dynamical state rather than  
405 only differences between populations.

## 406 Discussion

407 In this study we present SEAGALL (Single-cell ExplAInable Geometry-Aware Graph  
408 Attention Learning pipLine), a deep learning method based on manifold learning and  
409 explainable AI to analyse different modalities of single-cell data. SEAGALL combines  
410 a graph-regularised autoencoder (GRAE) and a graph attention network (GAT)  
411 together with an explainable artificial intelligence (XAI) method to classify the cells  
412 into cell types or states and extract the most important input features for the label  
413 prediction (Fig.1). We applied SEAGALL to ten single-cell data sets from three  
414 different omics (sc-RNAseq, sc-ATACseq, sc-ChIPseq) and have shown that  
415 SEAGALL can consistently understand and explain the cell identity from a different  
416 perspective than the classical differential analysis. The combination of a manifold  
417 learning method and an autoencoder (GRAE) to reduce the dimensionality of the  
418 data has been for the first time extensively applied to the single-cell field. This part of  
419 the workflow has been able to reconstruct corrupted data with the highest success  
420 (Fig.2), and we also showed that the biological information about cell type was  
421 robustly preserved while building the cell-to-cell graph (Fig.3). Therefore, this  
422 strategy has revealed an effective and reliable method to build a cell-to-cell graph,  
423 which is the final representation of the input data. Moreover, using the geometry  
424 aware graph as input to a classifier which applies an attention mechanism to  
425 increase the flexibility of the model, the classification performances reach a  
426 maximum (Fig.4). The main innovation of SEAGALL is the use of explainable AI to  
427 explore the cell type phenotype, making our method highly translational. Often, deep  
428 learning has focused on prediction performances, keeping the black box closed and  
429 preventing the direct gain of new biological knowledge. Here, we exploit a novel  
430 method of graph neural network explainer (GNNEExplainer) to open the black box and  
431 extract specific, stable and novel features (Fig.5, 6) which drive the cell type  
432 classification predictions. Thanks to its user-friendly code and tutorial, we have made  
433 our method suitable and useful for real world applications, since it can be directly  
434 applied to any count matrix from single-cell data. The deep learning method to learn  
435 the data sets ensures that the nonlinearity of the manifold, determined by the  
436 complicated gene regulatory networks, is taken into account, while standard  
437 approaches based on pca and DA do not. We have applied SEAGALL to several  
438 single-cell datasets and have shown that we are able to retrieve TFBSs which are  
439 driving factors of cell identity, but that would not have been identified using standard  
440 differential analysis pipeline (Fig.6). Finally, SEAGALL can be applied to different  
441 single-cell data modalities, such as scATAC-seq, scRNA-seq and scChIP-seq data,  
442 reflecting the omic-free hypothesis framework we proposed.

## 443 Methods

### 444 Single-cell RNAseq data processing

445 Single-cell RNAseq quantifies the abundance of RNA, mainly mRNA, molecules  
446 within a cell. For each single-cell the sequencer reads the transcripts that belong to  
447 it; hence the output is a raw set of reads which need to be aligned and quantified.  
448 For the two human multi-ome data sets (PBMC and brain), raw reads were  
449 processed using Cell Ranger Arc 2.0.2 aligning the reads onto the complete human  
450 genome (T2T)<sup>47</sup>. The GEX count matrix of the HSPC data set has been downloaded  
451 from<sup>48</sup>. The GEX count matrix of the mouse brain data set has been taken from 10X  
452 website<sup>48</sup>. We computed the probability distribution across cells of the number of non  
453 zero genes and the number of mitochondrial reads; we filtered out all the cells having  
454 a value of one of these variables lying outside the 5% or 95% quantile of their  
455 distribution. Similarly, genes present in less than 5% or more than 95% quantile of  
456 the cells were moved. Data were library-size normalised. We kept the top 10% highly  
457 variable genes. Last, data were log transformed. Differential expressed genes (DEG)  
458 between cell types were calculated using the Wilcoxon test. We kept for our analysis  
459 the fifty most differentially expressed genes.

### 460 Single-cell ATACseq data processing

461 Single-cell ATAC-seq is a popular technique to profile chromatin openness at the  
462 single-cell level. Typically, When analysing scATAC-seq data, typically the  
463 measurements are summarized in a count matrix using the positions of signal  
464 enrichment on the genome, called peaks<sup>49</sup>. To construct a count matrix, peaks are  
465 called on the pseudo-bulk signal and for each cell and each peak the number of  
466 reads that fall into the peak are counted. The structure of the matrix is identical to  
467 scRNA-seq, but in the latter case the features are the transcripts.

468 The reads of the kidney data set<sup>26</sup> have been processed using Cell Ranger ATAC  
469 2.1.0<sup>2</sup> and the ones of the PBMC and human brain data sets have been aligned with  
470 Cell Ranger Arc 2.0.2; in both cases the reference genome is the T2T human  
471 genome<sup>47</sup>. Count matrices were built using episcanpy<sup>12</sup>, given the fragments file and  
472 peak file obtained with MACS2<sup>50</sup>. We computed the distribution of the number of  
473 features per cell and we filtered out cells having a number of features lower than the  
474 5% quantile or higher than the 95% quantile of this distribution. Cells with lower than  
475 2 for the transcription start site (TSS) enrichment score, and higher than 2 for the  
476 nucleosome signal, have been filtered out. Features (peaks) present in lower than  
477 5% or more than the 95% quantile of the cells have been removed. Data were  
478 library-size normalised. Only peaks with a variance higher than the one defining the  
479 80% quantile of the variance distribution were kept, with a maximum of 30000  
480 features. Last, data were log transformed.

481 For the mouse data set, we downloaded the fragments file from 10X database. The  
482 fragments file of the HSPC data set was downloaded from the original publication<sup>48</sup>.  
483 Before building the count matrix and filtering, following the procedure described  
484 above, we called peaks using MACS2<sup>50</sup>.

485 The choice to use quantile as thresholds is motivated by the fact that this an  
486 automatic and fully reproducible method to apply quality controls on the data: instead  
487 of inspecting matrix by matrix and apply to every case a different threshold without  
488 being able to fully motivate the choice, the quantile computation takes into account  
489 the specific properties of the data set, ensuring the same rigidity on the filtering.

490 sc-ChIPseq experiment count matrices were downloaded from the original  
491 publication<sup>27</sup>. In this case the features are windows, that are constant size (50kb)  
492 intervals spanning the whole genome. We processed those data as peaks since the  
493 processing does not rely on any peaks-specific assumption. Differential open peaks  
494 or windows between cell types are calculated with the Wilcoxon test. We kept for our  
495 analysis the fifty most differentially open peaks or windows.

## 496 Cell type annotation

497 For the HSPC<sup>48</sup>, kidney<sup>26</sup> and breast cancer<sup>27</sup> data the cell type annotation is  
498 provided from the authors. The cell type annotation of the mouse brain is taken  
499 from<sup>48</sup> and it is based on marker genes. Human PBMC has been manually annotated  
500 following the muon tutorial<sup>51</sup>. Mouse brain has been manually annotated with marker  
501 genes and the procedure is shown in our github. Each data set consists of a different  
502 number of cell types (SuppTable1,2).

## 503 Embedding and graph construction

504 Once the count matrices are cleaned we use GRAE to build the cell-to-cell graph.  
505 First, PHATE is applied as a manifold learning method; it is able to capture both  
506 global and local structure of the data and embed it into a smaller representation with  
507 arbitrary dimension. The loss function of the autoencoder, which is the mean  
508 squared error (MSE) between original and reconstructed space, is then regularised  
509 by adding a term which increases if the AEs latent space differs more from the  
510 PHATE embedding. In other words, the total loss function  $L$  is composed by two  
511 terms: a reconstruction term  $L_r$  and and a regularisation term  $L_g$

512

$$(1) L(X, E) = L_r(X, f^{-1}(f(X))) + \lambda L_g(f(X), \Xi) = MSE(X, f^{-1}(f(X))) + \lambda \sum_{i=1}^N \|\xi_i - f(x_i)\|^2$$

513

514

515 where  $X$  is a set of  $N$  data points such that  $x_i \in \mathbb{R}^d$ ,  $\Xi$  is the PHATE embedding  
516 of  $X$  such that  $\xi_i \in \mathbb{R}^p$  with  $p \ll d$ ,  $f$  and  $f^{-1}$  are, respectively, the encoding  
517 and decoding function. The dimension of the latent representation varies for each

518 count matrix to fit the data set complexity and it is set as the cubic root of the number  
519 of features.

520 Within the latent space, pairwise euclidean distance between cells is computed and  
521 then a k-NN graph is built, with  $k = 15$ . k-NN graph had already been used in the  
522 literature as input graph for GNNs<sup>52</sup> but there is also a technical motivation that led  
523 us to a constant degree network: building a correlation or distance based graph is  
524 intrinsically problematic; after computing pairwise distances or correlations a cut-off  
525 is applied to the maximum distance or minimum correlation. Each node may have  
526 any number of neighbours in the interval  $[0, N - 1]$ . We tested this possibility and it  
527 turns out the resulting graph is extremely dense (SuppFig8), which may lead to  
528 nonsensical connections and makes the training of the GNN extremely time and  
529 energy demanding.

530 The graph is the final representation of the data set, which contains the connectivity  
531 pattern and the geometry of the input manifold.

## 532 Cell type classification with GNN

533 Graph neural networks are a particular type of neural network able to process data  
534 with a graph structure. GNNs take as input a graph  $G = (V, E)$ , where  $V \in \mathbb{N}$  is the  
535 set of nodes and  $E \subseteq V \times V$  is a set of edges, also known as links, between  
536 nodes. Each node can have a feature vector that defines the properties of the nodes.  
537 In our context, the feature vector is the gene expression or the chromatin openness  
538 vector. From a point cloud perspective, the embedding of each point is a function of  
539 the point itself and the points close to it. Let  $G = (V, E)$  be an undirected graph  
540 containing  $N$  vertices,  $x_i \in \mathbb{R}^d$  is the initial representation of node  $i$ ,  
541  $\mathcal{N}_i = \{j \in V | (j, i) \in E\}$  the neighbours of node  $i$ , then the first layer of the GNN  
542 will create a new representation of the node  $x'_i$  according to (2), known as the  
543 message passing equation<sup>53</sup>

544

$$(2) \mathbf{x}'_i = \gamma_{\Theta} \left( \mathbf{x}_i, \bigoplus_{j \in \mathcal{N}_i} \phi_{\Theta}(\mathbf{x}_i, \mathbf{x}_j, \mathbf{e}_{j,i}) \right).$$

545

546

547 GNN is a broad class of neural networks which rely on equation (2). We decided to  
548 apply a more refined version of the base message-passing layer called “GAT”<sup>20</sup>  
549 which applies an attention mechanism to the embedding function, meaning that the  
550 model learns the importance of each feature and link, following

551

$$(3) \mathbf{x}'_i = \sum_{j \in \mathcal{N}_i \cup \{i\}} \alpha_{i,j} \Theta_t \mathbf{x}_j$$

552

553

$$(4) \alpha_{i,j} = \frac{\exp(\mathbf{a}^\top \text{LeakyReLU}(\Theta_s \mathbf{x}_i + \Theta_t \mathbf{x}_j))}{\sum_{k \in \mathcal{N}(i) \cup \{i\}} \exp(\mathbf{a}^\top \text{LeakyReLU}(\Theta_s \mathbf{x}_i + \Theta_t \mathbf{x}_k))},$$

554

555

556 where  $\Theta \in \mathbb{R}^{d \times d'}$ ,  $\mathbf{a} \in \mathbb{R}^{2d'}$  are learned parameters,  $\bigoplus$  is any differentiable and  
557 permutation invariant function such as sum or mean, and  $\gamma_\Theta$  and  $\phi_\Theta$  are  
558 differentiable functions such as MLPs.

559 Thus, the key property of GNNs is the ability to create latent representations of a  
560 local neighbourhood rather than a single point. The rationale behind the choice of  
561 GNN relies on this property: we want to have a local analysis of each cell aiming for  
562 a local ensemble study, rather than treat them totally independently. Each count  
563 matrix with its own graph is given as input to the GNN classifier; the target output  
564 is the cell type of each node, which is defined as explained in the “cell type  
565 annotation” paragraph.

566 Our specific model consists of a graph neural network with two layers, first one to  
567 create a latent representation of the input and second one which performs the  
568 classification task. The dimension of each layer is defined with hyperparameter  
569 optimization (HPO)<sup>54</sup> case by case. The model is trained with Adam<sup>55</sup> optimizer with  
570 learning rate and weight decay estimated with HPO.

## 571 GAT explanation

572 Once the model is trained an XAI method is applied to it. We choose to apply  
573 “GNExplainer”<sup>22</sup>, which is a model agnostic method. It creates a graph and a  
574 feature mask to spot which is the minimum set of features and edges of each node  
575 sufficient to predict the class. We assume that the nature of our data defines a real  
576 function  $f$  that labels objects, nodes in the case of GNN, representing cells in our  
577 context. The GNN model  $\Phi$  receives as input a graph  $G$  and a feature vector  $X$  as  
578 explained in the previous paragraph. In practice,  $\Phi$  learns a probability  $P_\Phi(Y|G, X)$   
579 with  $Y$  random variable for the classes  $\{c_i\}_{i=1, \dots, C}$  representing the probability  
580 of nodes to belong to each of the classes. After the training, the model is fixed and it  
581 will be used to make predictions. The crucial point of the explainer is the fact that  
582 each node has a computation graph  $G$  and certain node features  $X$  that completely  
583 determine all the information that are necessary to predict  $\hat{y}$  at certain node  $v$ .  
584 Given a node  $v_i$  the explainer finds the subgraph  $G_s \subseteq G$  and the associated  
585 features  $X_s = \{x_j | v_j \in G_s\}$  that maximise the probability of having seen the  
586 prediction  $\hat{y} = \Phi(G_s, X_s)$  where  $\Phi$  is the trained GNN. Indicating as  $MI$  the mutual  
587 information function and  $H$  the entropy function, the GNExplainer solves the  
588 following problem

589

$$590 \max_{\{G_s\}} MI(Y, (G_s, X_s)) = H(Y) - H(Y|G = G_s, X = X_s).$$

591

592 MI quantifies the variation in the prediction probability when the graph and the  
593 features are  $G_s$  and  $X_s$  instead of  $G$  and  $X$ , with the feature vector constrained to  
594 be much smaller than the original one. In practice, for each node we obtain the  
595 features ranked by their importance. Since we are interested in the cell types  
596 explanations, we average the feature importance of all the nodes belonging to the  
597 same class to obtain the most relevant features for each label.

## 598 Topological and variational autoencoder models

599 Whereas GRAE<sup>19</sup>, PeakVI<sup>11</sup> and scVI<sup>10</sup> are released as packages, we had to  
600 implement the models for topological and variational autoencoder. Both the methods  
601 are based on the same architecture, which consists of one input layer, one hidden  
602 layer, with dimension equal to the square root of the input length, and a latent space,  
603 with dimension equal to the cubic root of the input layers. The varying size of the  
604 layers are important to account for the data set complexity. The best values of  
605 dropout, learning rate, weight decay, weight of topological regularisation and  
606 signature of the p-norm (for TopoAE) and weight of Kullback-Leibler divergence for  
607 VAE, have been estimated using HPO implemented with the *optuna* package<sup>56</sup>. Each  
608 HPO consists of 25 runs to explore the parameters within defined intervals  
609 (SuppTable8). We used a subset of count matrices to explore the HPO and we then  
610 applied the same parameters for each matrix.

## 611 Input data reconstruction

612 To test the robustness of the AEs we measured their ability to reconstruct the input  
613 data after corruption. To corrupt the data, we applied an increasing dropout from  
614 10% to 50% of the features in linear steps of 10% to each count matrix and we  
615 trained each model with the corrupted data. All the models have been trained with  
616 the same patience (20) and maximum number of epochs (300). We used 85% of the  
617 data for training and 15% for validation. When applying dropout the choice of the  
618 removed features is random, therefore it may happen that we remove some features  
619 particularly important for a specific model but not for another. To make sure our  
620 results are not biased by this factor, we repeat this experiment ten times varying the  
621 features to drop out at each level. For each run, for each level of dropout and for  
622 each model we measured the MSE between the original data (the not corrupted one)  
623 and the model reconstructed data. We then computed the average MSE for each  
624 level and model and the uncertainty of the mean.

## 625 Graph homogeneity

626 After applying each dimensionality reduction method as described in the previous  
627 paragraph (GRAE, TopoAE, VAE, PeakVI, scVI and pca), but without dropout, to  
628 each one of the count matrices, we computed the k-NN graph (k=15) from their latent  
629 spaces. For each cell we computed how many different cell types are found in its

630 neighbourhood. We divided this value by both the number of neighbours (15) and the  
631 number of cell types (different for each data set) to create a score called  
632 heterogeneity. Last we computed the homogeneity as 1-heterogeneity.

## 633 Classification and XAI experiments

634 To test the quality of each embedding method we used their latent space to build the  
635 cell-to-cell k-NN graphs (k=15) and we gave the graphs as input to a graph neural  
636 network node classifier. We tested the combination of the six dimensionality  
637 reduction methods (GRAE, TopoAE, VAE, PeakVI, scVI and pca) and two different  
638 GNN: GAT and GCN (see *Cell type classification with GNN* paragraph for the details  
639 of the models). We run each combination fifty times. Each training started with a  
640 different random seed to make sure that the models do not always start from the  
641 same point in the parameter space. Both GAT and GCN are trained for 250 epochs  
642 with a patience of 20 epochs. The data sets have been splitted into train, validation  
643 and test sets with a ratio of, respectively, 70%, 10% and 20%. Before training the  
644 classifier we run a 25 steps HPO study to select the best values (SuppFig9) of each  
645 hyperparameter of the GNN, within defined ranges (SuppTable9). After each training  
646 we applied the explainer for 200 epochs and saved the fifty most important features  
647 for each label, i.e. for each cell type. Accuracy, precision, recall and F1 score have  
648 been computed in the standard way using *sklearn*<sup>57</sup>; the specificity of the explainer is  
649 defined as one minus the average intersection of the top fifty most relevant features  
650 across cell types. Stability is defined as the average intersection of the explanation  
651 for the same cell type across different runs of the classifier and the explainer.

## 652 Data and code availability

653 Code and tutorial for SEAGALL are available at  
654 <https://github.com/gmalagol10/seagall.git>  
655 Data and notebook to reproduce all the results are available at  
656 <https://github.com/gmalagol10/seagall/tree/main/reproducibility>.  
657

## 658 References

- 659 1. Buenrostro, J. D., Giresi, P. G., Zaba, L. C., Chang, H. Y. & Greenleaf, W. J.  
660 Transposition of native chromatin for fast and sensitive epigenomic profiling of open  
661 chromatin, DNA-binding proteins and nucleosome position. *Nat Methods* **10**, 1213–1218  
662 (2013).
- 663 2. Satpathy, A. T. *et al.* Massively parallel single-cell chromatin landscapes of human  
664 immune cell development and intratumoral T cell exhaustion. *Nat. Biotechnol.* **37**,  
665 925–936 (2019).
- 666 3. Heumos, L. *et al.* Best practices for single-cell analysis across modalities. *Nature*  
667 *Reviews Genetics* **24**, 550–572 (2023).

- 668 4. Cayton, L. Algorithms for manifold learning. (2008).
- 669 5. Gene Regulatory Network. <http://dx.doi.org/10.1016/B978-0-12-809633-8.20294-X>  
670 doi:10.1016/B978-0-12-809633-8.20294-X.
- 671 6. Wolf, F. A., Angerer, P. & Theis, F. J. SCANPY: large-scale single-cell gene expression  
672 data analysis. *Genome Biol.* **19**, 1–5 (2018).
- 673 7. Bravo González-Blas, C. *et al.* cisTopic: cis-regulatory topic modeling on single-cell  
674 ATAC-seq data. *Nat Methods* **16**, 397–400 (2019).
- 675 8. Stuart, T., Srivastava, A., Madad, S., Lareau, C. A. & Satija, R. Single-cell chromatin  
676 state analysis with Signac. *Nat Methods* **18**, 1333–1341 (2021).
- 677 9. Fang, R. *et al.* Comprehensive analysis of single cell ATAC-seq data with SnapATAC.  
678 *Nat Commun* **12**, 1337 (2021).
- 679 10. Lopez, R., Regier, J., Cole, M. B., Jordan, M. I. & Yosef, N. Deep generative modeling  
680 for single-cell transcriptomics. *Nature Methods* **15**, 1053–1058 (2018).
- 681 11. PeakVI: A deep generative model for single-cell chromatin accessibility analysis. *Cell*  
682 *Reports Methods* **2**, 100182 (2022).
- 683 12. Danese, A. *et al.* EpiScanpy: integrated single-cell epigenomic analysis. *Nat. Commun.*  
684 **12**, 1–8 (2021).
- 685 13. Granja, J. M. *et al.* ArchR is a scalable software package for integrative single-cell  
686 chromatin accessibility analysis. *Nature Genetics* **53**, 403–411 (2021).
- 687 14. Integrated analysis of multimodal single-cell data. *Cell* **184**, 3573–3587.e29 (2021).
- 688 15. Rumelhart, D. E., Hinton, G. E. & Williams, R. J. Learning representations by  
689 back-propagating errors. *Nature* **323**, 533–536 (1986).
- 690 16. Gayoso, A. *et al.* Joint probabilistic modeling of single-cell multi-omic data with totalVI.  
691 *Nature Methods* **18**, 272–282 (2021).
- 692 17. Tangherloni, A., Ricciuti, F., Besozzi, D., Liò, P. & Cvejic, A. Analysis of single-cell RNA  
693 sequencing data based on autoencoders. *BMC Bioinformatics* **22**, 1–27 (2021).
- 694 18. Drost, F. *et al.* Multi-modal generative modeling for joint analysis of single-cell T cell  
695 receptor and gene expression data. *Nature Communications* **15**, 1–15 (2024).
- 696 19. Duque, A. F., Morin, S., Wolf, G. & Moon, K. R. Geometry Regularized Autoencoders.  
697 <https://ieeexplore.ieee.org/document/9950332>.
- 698 20. Brody, S., Alon, U. & Yahav, E. How Attentive are Graph Attention Networks? in  
699 *International Conference on Learning Representations* (2021).
- 700 21. Veličković, P. *et al.* Graph Attention Networks. (2017).
- 701 22. Ying, Z., Bourgeois, D., You, J., Zitnik, M. & Leskovec, J. GNNExplainer: Generating  
702 Explanations for Graph Neural Networks. *Adv. Neural Inf. Process. Syst.* **32**, (2019).
- 703 23. van Dijk Guy Wolf Smita Krishnaswamy, K. R. M. J. S. S. D. B. D. Manifold  
704 learning-based methods for analyzing single-cell RNA-sequencing data. *Current*  
705 *Opinion in Systems Biology* **7**, 36–46 (2018).
- 706 24. Moon, K. R. *et al.* Visualizing structure and transitions in high-dimensional biological  
707 data. *Nat. Biotechnol.* **37**, 1482–1492 (2019).
- 708 25. Scarselli, F., Gori, M., Tsoi, A. C., Hagenbuchner, M. & Monfardini, G. The Graph Neural  
709 Network Model. <https://ieeexplore.ieee.org/document/4700287>.
- 710 26. Sheng, X. *et al.* Mapping the genetic architecture of human traits to cell types in the  
711 kidney identifies mechanisms of disease and potential treatments. *Nat. Genet.* **53**,  
712 1322–1333 (2021).
- 713 27. Grosselin, K. *et al.* High-throughput single-cell ChIP-seq identifies heterogeneity of  
714 chromatin states in breast cancer. *Nat. Genet.* **51**, 1060–1066 (2019).
- 715 28. Moor, M., Horn, M., Rieck, B. & Borgwardt, K. Topological Autoencoders. in *International*

- 716 *Conference on Machine Learning* 7045–7054 (PMLR, 2020).
- 717 29. Kipf, T. N. & Welling, M. Semi-Supervised Classification with Graph Convolutional  
718 Networks. in *International Conference on Learning Representations* (2022).
- 719 30. Heinz, S. *et al.* Simple combinations of lineage-determining transcription factors prime  
720 cis-regulatory elements required for macrophage and B cell identities. *Mol. Cell* **38**,  
721 576–589 (2010).
- 722 31. Claus Stolt, C. *et al.* The Sox9 transcription factor determines glial fate choice in the  
723 developing spinal cord. *Genes Dev.* **17**, 1677–1689 (2003).
- 724 32. The transcription factor PITX1 drives astrocyte differentiation by regulating the SOX9  
725 gene. *Journal of Biological Chemistry* **295**, 13677–13690 (2020).
- 726 33. Leonard, J. *et al.* Transcriptomic alterations in cortical astrocytes following the  
727 development of post-traumatic epilepsy. *Scientific Reports* **14**, 1–12 (2024).
- 728 34. Ma, S. *et al.* Chromatin Potential Identified by Shared Single-Cell Profiling of RNA and  
729 Chromatin. *Cell* **183**, 1103–1116.e20 (2020).
- 730 35. Serrano-Pozo, A. *et al.* Astrocyte transcriptomic changes along the spatiotemporal  
731 progression of Alzheimer’s disease. *Nature Neuroscience* **27**, 2384–2400 (2024).
- 732 36. Supakul, S. *et al.* Mutual interaction of neurons and astrocytes derived from iPSCs with  
733 APP V717L mutation developed the astrocytic phenotypes of Alzheimer’s disease.  
734 *Inflammation and Regeneration* **44**, 1–21 (2024).
- 735 37. Nakajima-Koyama, M., Lee, J., Ohta, S., Yamamoto, T. & Nishida, E. Induction of  
736 Pluripotency in Astrocytes through a Neural Stem Cell-like State. *The Journal of*  
737 *Biological Chemistry* **290**, 31173 (2015).
- 738 38. Holtman, I. R., Skola, D. & Glass, C. K. Transcriptional control of microglia phenotypes  
739 in health and disease. *J Clin Invest* **127**, 3220–3229 (2017).
- 740 39. Laflamme, N., Soucy, G. & Rivest, S. Circulating cell wall components derived from  
741 gram-negative, not gram-positive, bacteria cause a profound induction of the  
742 gene-encoding Toll-like receptor 2 in the CNS. *Journal of neurochemistry* **79**, (2001).
- 743 40. Yang, C. *et al.* RIPK2 Is Crucial for the Microglial Inflammatory Response to Bacterial  
744 Muramyl Dipeptide but Not to Lipopolysaccharide. *International Journal of Molecular*  
745 *Sciences* **25**, 11754 (2024).
- 746 41. Liu, Y.-H. *et al.* Ascl1 promotes tangential migration and confines migratory routes by  
747 induction of Ephb2 in the telencephalon. *Scientific Reports* **7**, 1–17 (2017).
- 748 42. Sawai, C. M. *et al.* Transcription factor Runx2 controls the development and migration of  
749 plasmacytoid dendritic cells. *J Exp Med* **210**, 2151–2159 (2013).
- 750 43. Transcriptomic and genomic heterogeneity in blastic plasmacytoid dendritic cell  
751 neoplasms: from ontogeny to oncogenesis. *Blood Advances* **5**, 1540–1551 (2021).
- 752 44. Fischer, M. B. *et al.* Dependence of Germinal Center B Cells on Expression of  
753 CD21/CD35 for Survival. *Science* (1998) doi:10.1126/science.280.5363.582.
- 754 45. Rebuffet, L. *et al.* High-dimensional single-cell analysis of human natural killer cell  
755 heterogeneity. *Nature Immunology* **25**, 1474–1488 (2024).
- 756 46. Dias, J. *et al.* The CD4–CD8– MAIT cell subpopulation is a functionally distinct subset  
757 developmentally related to the main CD8+ MAIT cell pool. *Proceedings of the National*  
758 *Academy of Sciences* **115**, E11513–E11522 (2018).
- 759 47. Nurk, S. *et al.* The complete sequence of a human genome. *Science* (2022)  
760 doi:10.1126/science.abj6987.
- 761 48. Li, C., Virgilio, M. C., Collins, K. L. & Welch, J. D. Multi-omic single-cell velocity models  
762 epigenome–transcriptome interactions and improves cell fate prediction. *Nat.*  
763 *Biotechnol.* **41**, 387–398 (2022).

- 764 49. Yan, F., Powell, D. R., Curtis, D. J. & Wong, N. C. From reads to insight: a hitchhiker's  
765 guide to ATAC-seq data analysis. *Genome Biol.* **21**, 1–16 (2020).
- 766 50. Zhang, Y. *et al.* Model-based Analysis of ChIP-Seq (MACS). *Genome Biol.* **9**, 1–9  
767 (2008).
- 768 51. Processing gene expression of 10k PBMCs — muon-tutorials documentation.  
769 [https://muon-tutorials.readthedocs.io/en/latest/single-cell-rna-atac/pbmc10k/1-Gene-Exp](https://muon-tutorials.readthedocs.io/en/latest/single-cell-rna-atac/pbmc10k/1-Gene-Expression-Processing.html)  
770 [ression-Processing.html](https://muon-tutorials.readthedocs.io/en/latest/single-cell-rna-atac/pbmc10k/1-Gene-Expression-Processing.html).
- 771 52. Wang, J. *et al.* scGNN is a novel graph neural network framework for single-cell  
772 RNA-Seq analyses. *Nat. Commun.* **12**, 1–11 (2021).
- 773 53. Gilmer, J., Schoenholz, S. S., Riley, P. F., Vinyals, O. & Dahl, G. E. Neural Message  
774 Passing for Quantum Chemistry. (2017).
- 775 54. Akiba, T., Sano, S., Yanase, T., Ohta, T. & Koyama, M. Optuna: A Next-generation  
776 Hyperparameter Optimization Framework. (2019).
- 777 55. Kingma, D. P. & Ba, J. Adam: A Method for Stochastic Optimization. (2014).
- 778 56. Optuna. <https://dl.acm.org/doi/10.1145/3292500.3330701>  
779 [doi:10.1145/3292500.3330701](https://doi.org/10.1145/3292500.3330701).
- 780 57. Pedregosa, F. *et al.* Scikit-learn: Machine Learning in Python. *Journal of Machine*  
781 *Learning Research* **12**, 2825–2830 (2011).
- 782 58. Scientific Image and Illustration Software. <https://BioRender.com>.

## 783 Acknowledgments

784 We thank Samuele Firmani for the insightful discussion about the evaluation of the models.  
785 We thank Vera Manelli for the important help in the interpretation of motif analysis. We thank  
786 Federica Tosato for the suggestions about visualisation and graphics. We thank Gaia  
787 Fontana for drawing the logo of the tool.

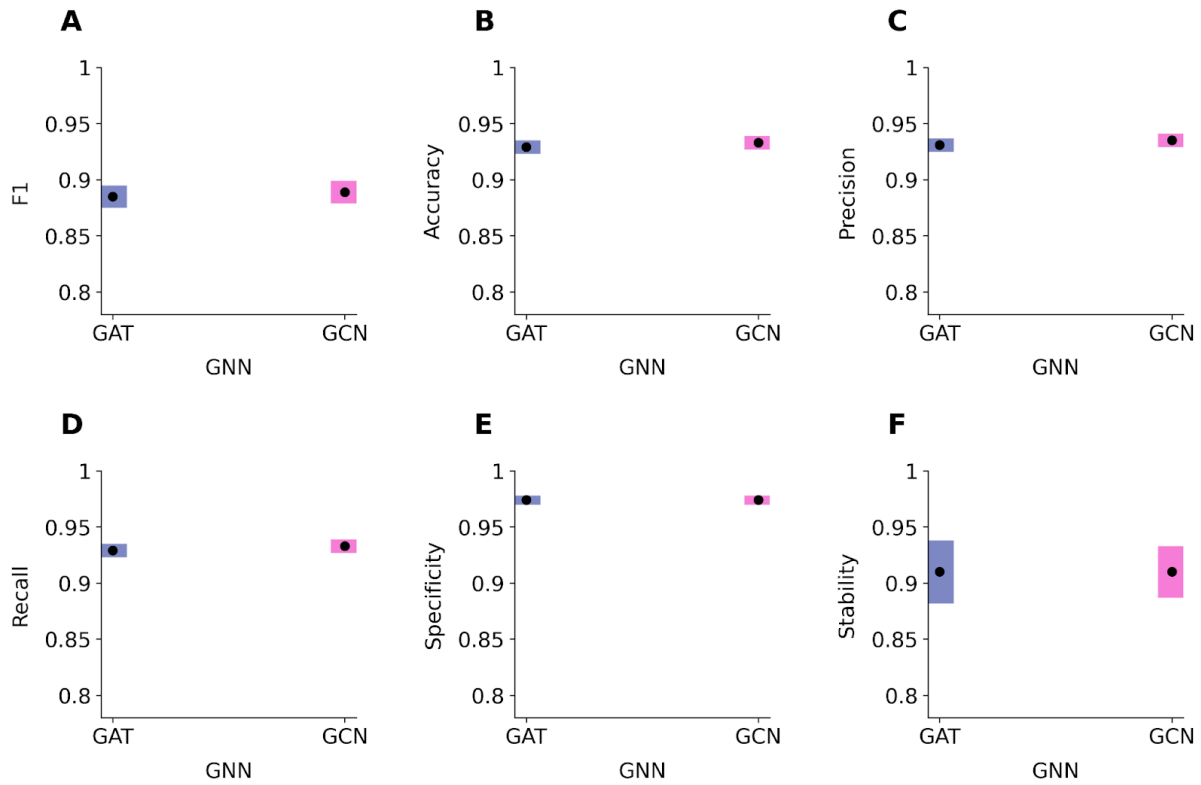
788 G.M. is supported by the Helmholtz International Lab Causal Cell Dynamics  
789 (InterLabs-0029) - Grant support from the Initiative and Networking Fund of the Hermann  
790 von Helmholtz-Association Deutscher Forschungszentren e.V.-. G.M is also supported by the  
791 Helmholtz Association under the joint research school “Munich School for Data Science —  
792 MUDS and by German Research Foundation project ID 213249687–SFB 1064.

793 We thank the BMC Bioinformatics Core Facility for providing access to their HPC cluster.  
794 Figures have been made with the help of Biorender<sup>58</sup>.

## 795 Contributions

796 G.M., M.C.T. and G.W. designed the study and conceived the algorithm. G.M. implemented  
797 the algorithm. P.H. provided code and helped the implementation of it. A.D. annotated the  
798 human brain data set. G.M. and M.C.T. wrote the manuscript with help from G.W. and  
799 additional inputs from all co-authors. All authors reviewed and approved the manuscript.

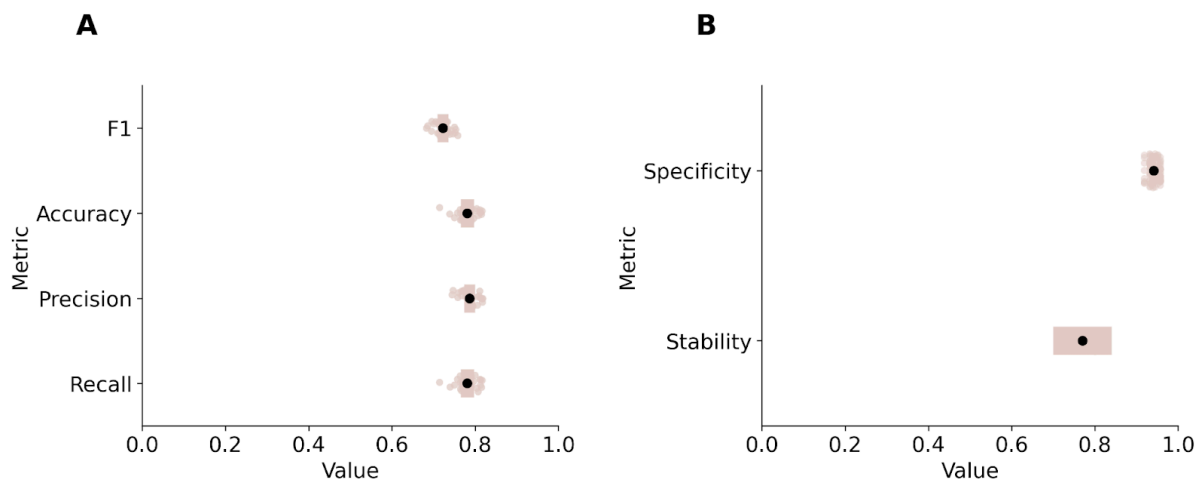
## 800 Supplementary Figures



801

802 SuppFig1: **A-D** classification of the two GNN architectures. Black dots indicate the mean over all the  
803 runs and the different embedding methods and the bar represent three times the uncertainty on the  
804 mean. **E-F** specificity (left) and stability (right) of the two GNN architectures. In none of the six metrics  
805 here represented we can appreciate a significant difference between the two models.

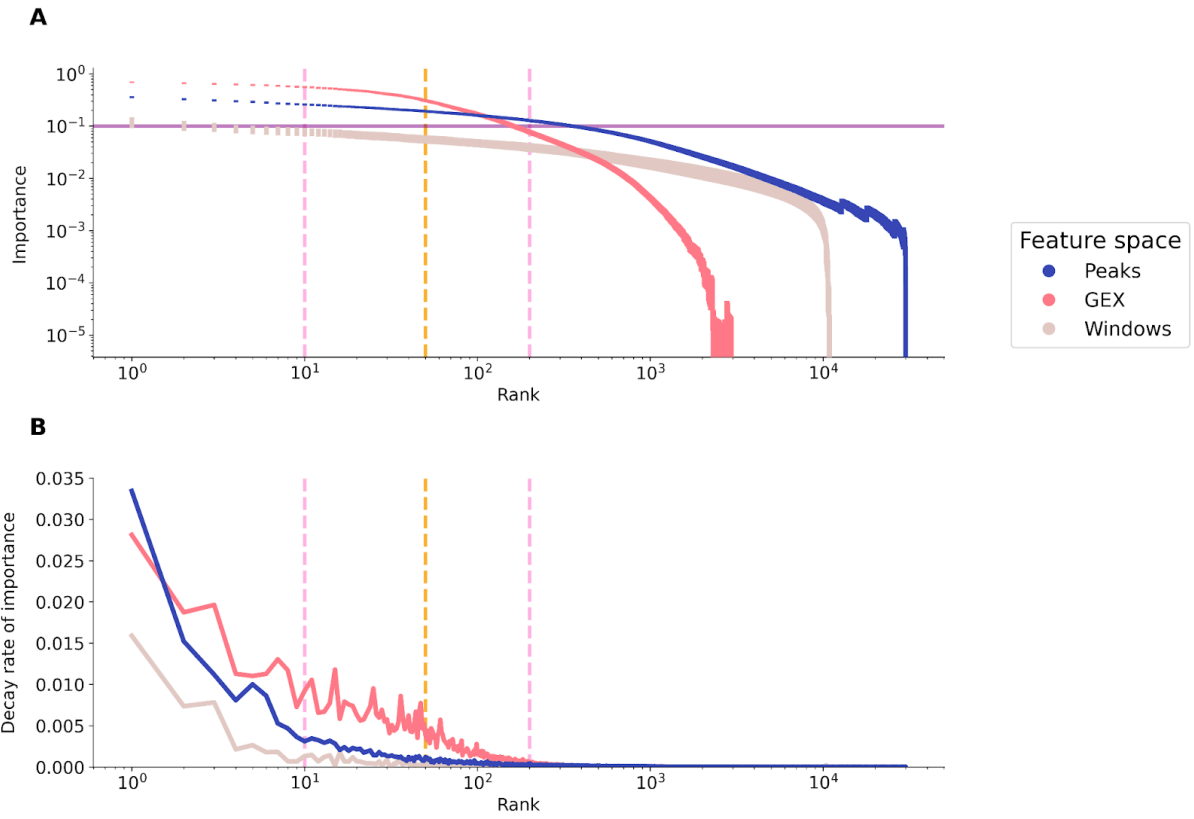
806



807

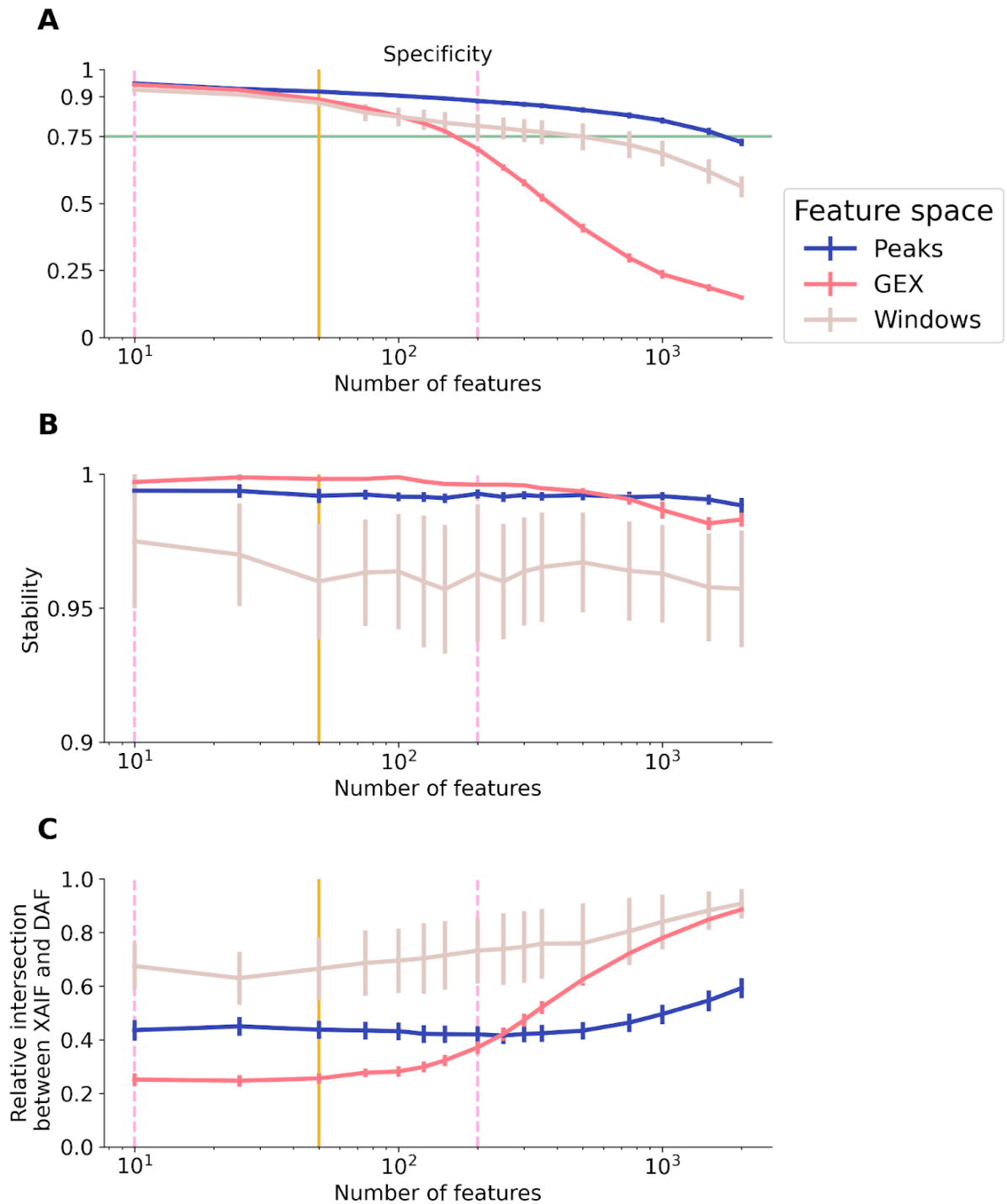
808 SuppFig2: **A** classification performances of the final model on the scChIP-seq. **B** stability and  
809 specificity of the explainer on the scChIP-seq data set.

810



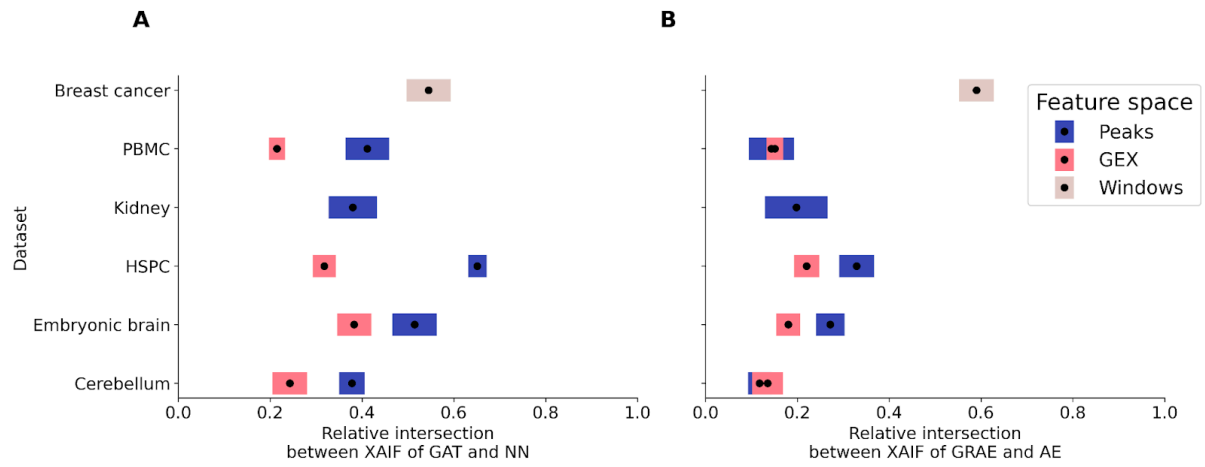
811

812 SuppFig3: **A** full range rank-importance distribution of the features for each feature space. **B** full range  
813 derivative of importance distribution of the features for each feature space.



814

815 SuppFig4: **A** specificity of the explanations varying the number of features we asked the model to  
816 keep. Each color represents one feature space; the mean value at each step is the average of the  
817 specificity for each cell type in each data set for each after running the explainer fifty times. The error  
818 bar height represents the standard deviation of the mean. **B** stability of the explanation computed as  
819 described in A. **C** Similarity of the XAIF to the DAF measured in terms of relative overlapping.



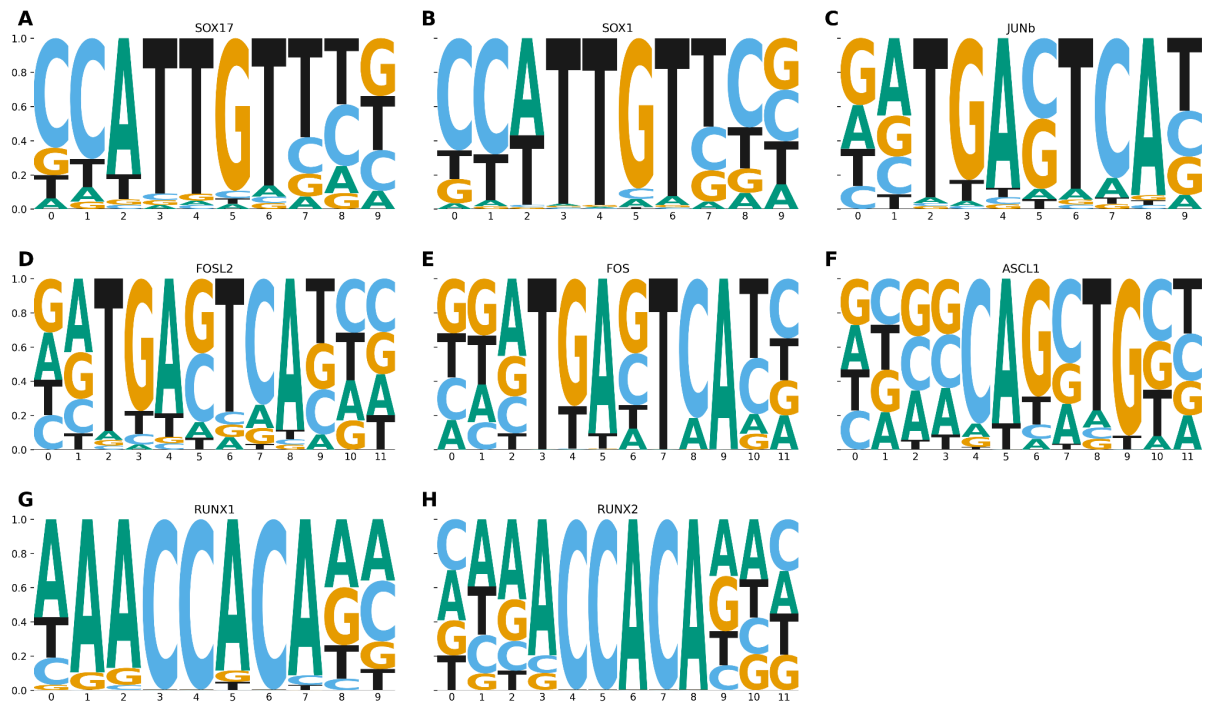
820

821 SuppFig5: **A** Difference between the XAIFs of GAT and a normal NN without any graph information. **B**

822 Difference between the XAIFs of GRAE versus AE for graph construction, applying the same GAT.

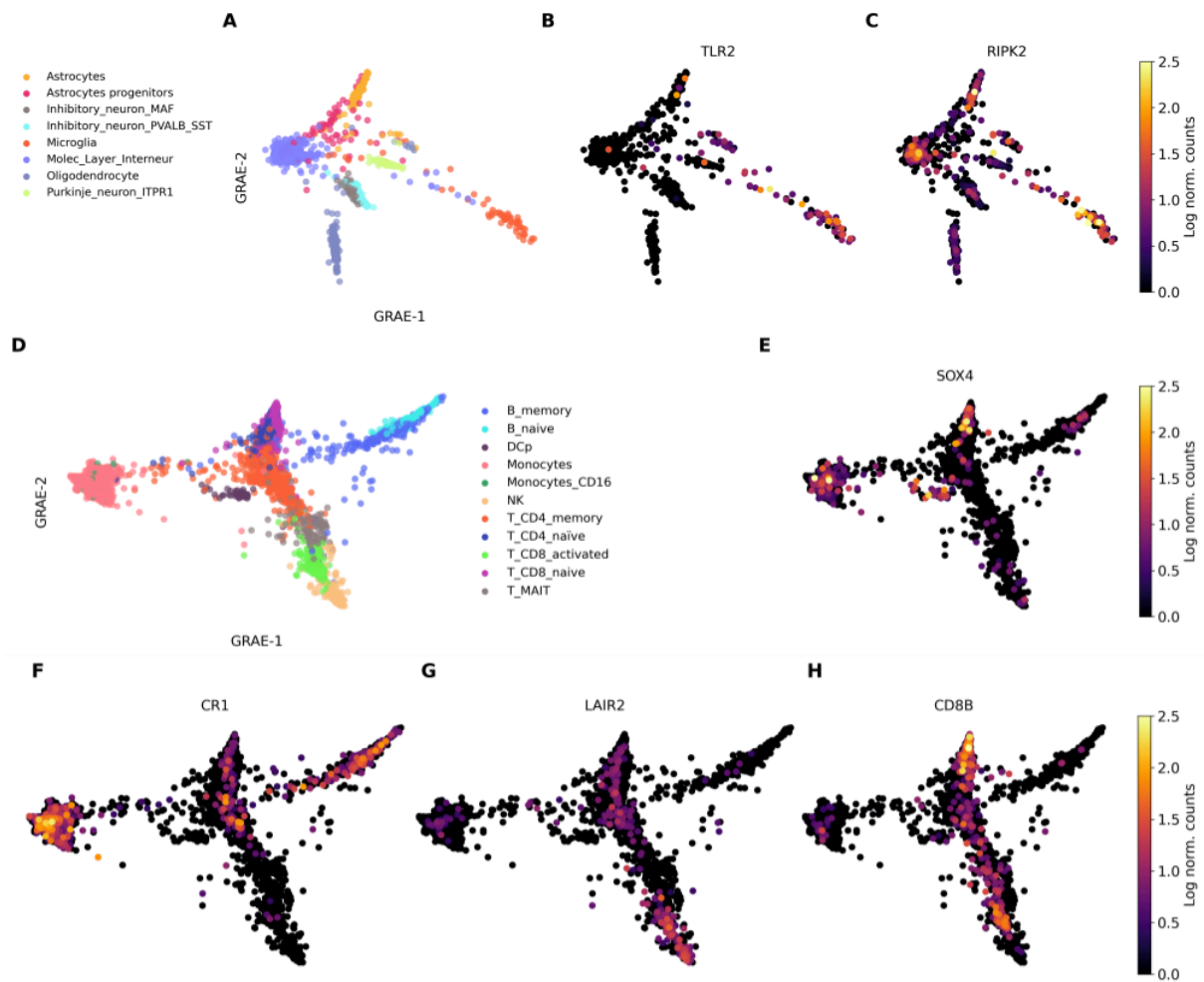
823 Removing either geometry or attention leads to different explanations.

824



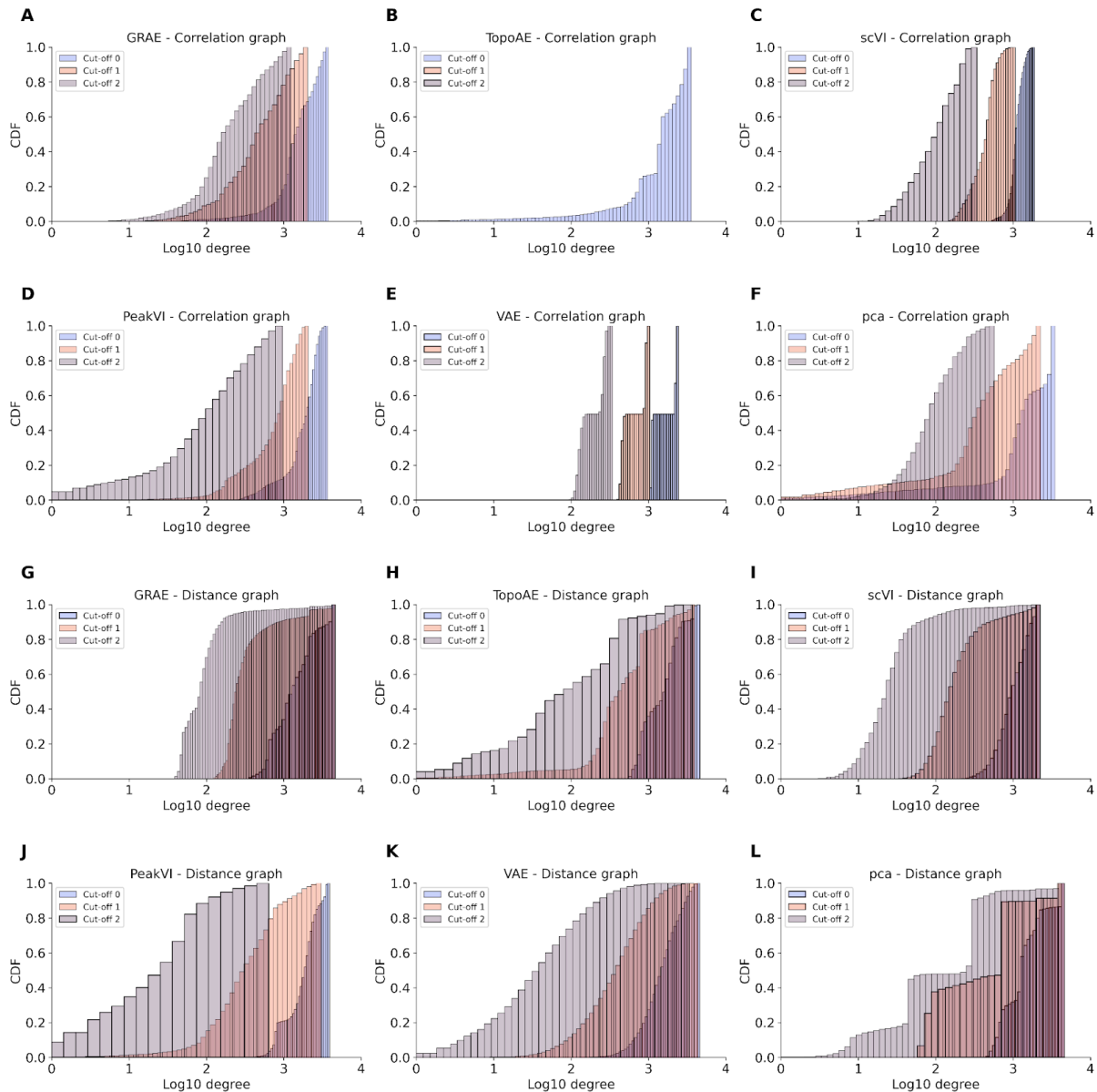
825

826 SuppFig6 **A-B** motif visualisation



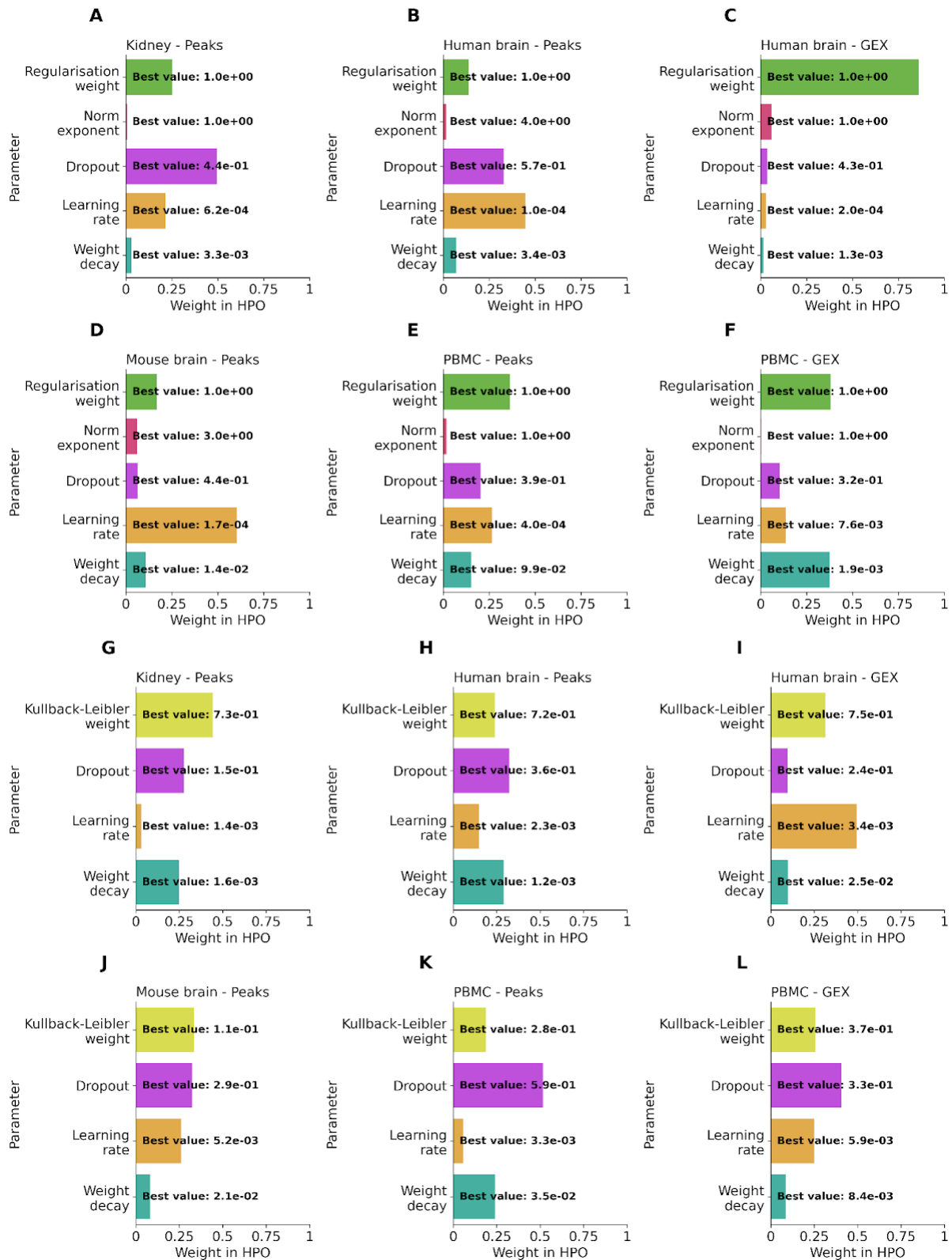
827

828 SuppFig7: **A** GRAE embedding of the sc-ATAC of the human brain data set. **B-C** expression of TLR2 and RIPK2  
829 in normalised log counts. **D** GRAE embedding of the GEX of the PBMC data set. **E-H** expression of SOX4, CR1,  
830 LAIR2 and CD8 (sub unit B) in normalised log counts.



831

832 SuppFig8 A-F: cumulative distribution function of nodes degree of the correlation graph varying the  
833 latent space (across panel) and applying different cut-off on the minimum correlation between nodes  
834 to keep the link. G-L show the same information but based on euclidean distance computed in the  
835 latent space of the embedding method. In all the cases, the graphs are very dense making the edges  
836 meaningless and the computation very time and energy demanding.



837

838 SuppFig9: impact (weight) of each hyperparameter in the optimization of either TopoAE or VAE. The  
 839 bold annotations indicate the best values of the hyperparameter, i.e. the ones used in all the  
 840 experiments involving TopoAE (A-F) and VAE (G-L). Each panel corresponds to one data set. The  
 841 final value of each hyperparameter is the average across the six experiments.

## 842 Supplementary Tables

843

	Human brain	Breast cancer	PBMC	HSPC	Mouse brain	Kidney
0	Molec_Layer_Interneur	HBCx-95	Monocytes_CD16	HSC	Subplate	Endo
1	Microglia	HBCx-95-Ca paR	T_CD8_naive	LMPP	Deeper_Layer	PT
2	Oligodendrocyte	HBCx-22	B_naive	MPP	Upper_Layer	CDPC
3	Purkinje_neuron_layer	HBCx-22-Ta mR	B_memory	MEP	IPC	Immune
4	Inhibitory_neuron_PVALB_SST		NK	Prog_B	V_SVZ	LOH
5	Purkinje_neuron_ITPR1		T_MAIT	Prog_DC	RG_Astro_OPC	CDIC
6	Inhibitory_neuron		T_CD4_naive	Erythrocyte	Ependymal_cells	Podo
7	Purkinje_neuron_FOXP2		DCm	GMP		DCT
8	Astrocyte		Monocytes	Granulocyte		
9	Inhibitory_neuron_MAF		T_CD4_memory	Prog_MK		
10	Astrocyte_progenitor		T_CD8_activated	Platelet		
11			DCp			

844 Supplementary table 1: cell type composition of each data set.

845  
846  
847

Data set	Feature space	Number of features	Number of cells
HSPC	GEX	2689	10389
HSPC	Peaks	30000	10350
Kidney	Peaks	25901	4624
Human brain	GEX	3003	2441
Human brain	Peaks	17886	2726
PBMC	GEX	2761	5700
PBMC	Peaks	12781	6372
Mouse brain	GEX	2237	2721
Mouse brain	Peaks	29114	3027
Breast cancer	Windows	10908	4636

848 Supplementary table 2: number of features and cells of each count matrix.

849  
850

	Kidney Peaks	Human brain Peaks	Human brain GEX	Mouse brain Peaks	PBMC Peaks	PBMC GEX
GRAE	<b>0.9905±0.0003</b>	<b>0.9911±0.0004</b>	<b>0.9917±0.0003</b>	<b>0.9849±0.0004</b>	<b>0.9909±0.0002</b>	<b>0.9921±0.0002</b>
TopoAE	0.9875±0.0004	0.9894±0.0005	<b>0.9915±0.0004</b>	0.9799±0.0006	0.9896±0.0002	0.9907±0.0002
VAE	0.9634±0.0005	0.9666±0.0007	0.9668±0.0006	0.9570±0.0008	0.9671±0.0004	0.9671±0.0004
PeakVI	0.9882±0.0003	<b>0.9909±0.0004</b>		<b>0.9846±0.0005</b>	0.9903±0.0002	
scVI			<b>0.9921±0.0003</b>			<b>0.9922±0.0002</b>
pca	0.9766±0.0004	0.9771±0.0008	<b>0.9921±0.0003</b>	0.9701±0.0006	0.9856±0.0003	<b>0.9923±0.0007</b>

851 Supplementary table 3: average homogeneity of each k-NN graph computed from the latent space of  
852 the AEs. The uncertainty is three times the standard deviation of the mean. Bold numbers highlight  
853 the method achieving the best performances; multiple bold numbers indicate that the methods  
854 performances fall within the uncertainty. The empty spaces represent the forbidden combinations of  
855 AE and feature spaces, which are GEX with PeakVI and peaks with scVI.

856  
857  
858  
859  
860  
861  
862

	Accuracy	F1	Precision	Recall	Specificity
GAT	0.929±0.006	0.89±0.01	0.93±0.01	0.93±0.01	0.974±0.004
GCN	0.933±0.006	0.89±0.01	0.94±0.01	0.93±0.01	0.974±0.004

863 Supplementary table 4: classification and explanations performances of the two GNN architectures.  
 864 GCN and GAT achieve the same performances. Values are the mean over all the runs and the using  
 865 GRAE as embedding method. The uncertainty is three times the standard deviation of the mean.  
 866  
 867

	Accuracy	F1	Precision	Recall	Specificity	Stability
GRAE	<u>0.929±0.006</u>	<u>0.885±0.010</u>	<u>0.931±0.006</u>	<u>0.929±0.006</u>	<u>0.974±0.004</u>	<b>0.91±0.03</b>
TopoAE	0.857±0.013	0.778±0.015	0.865±0.012	0.857±0.013	<u>0.975±0.004</u>	<b>0.91±0.03</b>
VAE	0.369±0.034	0.33±0.04	0.48±0.04	0.369±0.034	0.960±0.007	0.81±0.04
PeakVI	0.909±0.007	0.842±0.013	0.909±0.007	0.909±0.007	<b>0.981±0.001</b>	<b>0.89±0.05</b>
scVI	<b>0.956±0.006</b>	<b>0.928±0.005</b>	<b>0.959±0.006</b>	<b>0.956±0.006</b>	<b>0.975±0.007</b>	<b>0.92±0.02</b>
pca	0.734±0.032	0.665±0.033	0.781±0.023	0.734±0.032	<b>0.980±0.001</b>	0.84±0.04

868 Supplementary table 5: performances of the GAT classifier varying the embedding methods. For each  
 869 method are reported the mean across 50 runs and three times the uncertainty on the mean. GRAE  
 870 outperforms all the methods except scVI in terms of accuracy, F1 score, precision and recall. In bold  
 871 is highlighted the best performing method, the second best one is underlined. See SuppTable4 and  
 872 SuppTable5 for the performances divided by feature space.  
 873  
 874

	Accuracy	F1	Precision	Recall	Specificity	Stability
GRAE	<b>0.919±0.008</b>	<b>0.86±0.01</b>	<b>0.922±0.008</b>	<b>0.919±0.008</b>	<b>0.969±0.006</b>	<b>0.92±0.04</b>
TopoAE	0.83±0.01	0.73±0.04	0.84±0.04	0.83±0.04	<b>0.977±0.002</b>	<b>0.91±0.04</b>
VAE	0.30±0.03	0.23±0.03	0.38±0.04	0.30±0.03	<b>0.96±0.01</b>	0.78±0.04
PeakVI	<b>0.909±0.007</b>	<b>0.84±0.01</b>	<b>0.909±0.007</b>	<b>0.91±0.01</b>	<b>0.981±0.001</b>	<b>0.89±0.05</b>
pca	0.63±0.03	0.54±0.02	0.70±0.01	0.63±0.03	<b>0.979±0.002</b>	0.79±0.05

875 Supplementary table 6: classification and explanation performances of the embedding methods as  
 876 listed in Table but only for peaks count matrices.  
 877  
 878  
 879  
 880  
 881  
 882  
 883

884  
885  
886  
887  
888

	Accuracy	F1	Precision	Recall	Specificity	Stability
GRAE	<b>0.95±0.01</b>	<b>0.929±0.007</b>	<b>0.950±0.009</b>	<b>0.95±0.01</b>	<b>0.9819±0.0007</b>	<b>0.89±0.04</b>
TopoAE	0.91±0.02	0.87±0.02	0.92±0.02	<b>0.92±0.02</b>	<b>0.97±0.01</b>	<b>0.90±0.03</b>
VAE	0.50±0.05	0.53±0.03	0.67±0.03	0.50±0.05	<b>0.9668±0.0009</b>	<b>0.87±0.04</b>
scVI	<b>0.956±0.006</b>	<b>0.928±0.005</b>	<b>0.959±0.006</b>	<b>0.956±0.006</b>	<b>0.975±0.007</b>	<b>0.92±0.02</b>
pca	<b>0.950±0.007</b>	<b>0.921±0.006</b>	<b>0.953±0.007</b>	<b>0.950±0.007</b>	<b>0.9822±0.0006</b>	<b>0.92±0.02</b>

889 Supplementary table 7: classification and explanation performances of the embedding methods as  
890 listed in Table but only for GEX count matrices.

891  
892

	Regularisation weight	P for the P-norm	Weight decay	Learning rate	KL weight	Dropout
Min value	1	1	0.0001	0.0001	0.1	0.1
Max value	30	5	0.1	0.1	0.9	0.7
AE	TopoAE	TopoAE	TopoAE VAE	TopoAE VAE	VAE	TopoAE VAE

893 Supplementary table 8: minimum and maximum explored values during HPO of the  
894 parameters of AEs.

895  
896

	Hidden dimension	Number of heads	Weight decay	Learning rate
Min value	32	4	0.0001	0.0001
Max value	256	12	0.1	0.1

897 Supplementary table 9: minimum and maximum explored values during HPO for the GNNs  
898 parameters.

899  
900  
901  
902  
903  
904  
905  
906  
907

## 908 Extended Tables

909

Data set	Cell type	DA Motif Name	DA Consensus	DA q-value (Benjamini)	SEAGALL Motif Name	SEAGALL Consensus	SEAGALL q-value (Benjamini)
Human brain	Astrocyte_progenitor				Sox10(HMG)/SciaticNerve-Sox3-ChIP-Seq(GSE35132)/Homer	CCWTTGTY YB	0.0048
Human brain	Astrocyte_progenitor				Sox4(HMG)/proB-Sox4-ChIP-Seq(GSE50066)/Homer	YCTTTGTT CC	0.0048
Human brain	Astrocyte_progenitor				Sox9(HMG)/Limb-SOX9-ChIP-Seq(GSE73225)/Homer	AGGVNCCT TTGT	0.0152
Human brain	Astrocyte_progenitor				Sox15(HMG)/CPA-Sox15-ChIP-Seq(GSE62909)/Homer	RAACAATG GN	0.0152
Human brain	Astrocyte_progenitor				Sox17(HMG)/Endoderm-Sox17-ChIP-Seq(GSE61475)/Homer	CCATTGTT YB	0.0152
Human brain	Astrocyte_progenitor				Sox21(HMG)/ESC-SOX21-ChIP-Seq(GSE110505)/Homer	BCCWTTGT BYKV	0.0197
Human brain	Astrocyte_progenitor				SOX1(HMG)/NPC-SOX1-ChIP-Seq(GSE138215)/Homer	CCATTGTT CB	0.0197

Human brain	Astrocyte_progenitor					hINR(CPE)	SBCABW	0.0374
Human brain	Astrocyte_progenitor					Sox3(HMG)/NPC-Sox3-ChIP-Seq(GSE33059)/Homer	CCWTTGTY	0.0374
Human brain	Astrocyte_progenitor					Sox6(HMG)/Myotubes-Sox6-ChIP-Seq(GSE32627)/Homer	CCATTGTTY	0.0413
Human Brain	Microglia	ETS1(ETS)/Jurkat-ETS1-ChIP-Seq(GSE17954)/Homer	ACAGGAAGTG	0				
Human Brain	Microglia	SpiB(ETS)/OCILY3-SPIB-ChIP-Seq(GSE56857)/Homer	AAAGRGGAGTG	0				
Human Brain	Microglia	Fli1(ETS)/CD8-FLI-ChIP-Seq(GSE20898)/Homer	NRYTTCCGH	0				
Human Brain	Microglia	PU.1(ETS)/ThioMac-PU.1-ChIP-Seq(GSE21512)/Homer	AGAGGAAGTG	0				
Human Brain	Microglia	Elf4(ETS)/BMDM-Elf4-ChIP-Seq(GSE88699)/Homer	ACTTCCKGKT	0				
Human Brain	Microglia	GABPA(ETS)/Jurkat-GABPa-ChIP-Seq(GS	RACCGGAAGT	0				

		E17954)/Homer					
Human Brain	Microglia	ETV1(ETS)/GIST48-ETV1-ChIP-Seq(GSE22441)/Homer	AACCGGAA GT	0			
Human Brain	Microglia	ERG(ETS)/CaP-ER G-ChIP-Seq(GSE14097)/Homer	ACAGGAAG TG	0			
Human Brain	Microglia	Etv2(ETS)/ES-ER71-ChIP-Seq(GSE59402)/Homer	NNAYTTCCT GHN	0			
Human Brain	Microglia	ELF5(ETS)/T47D-ELF5-ChIP-Seq(GSE30407)/Homer	ACVAGGAA GT	0			
Human Brain	Microglia	ELF3(ETS)/PDAC-ELF3-ChIP-Seq(GSE64557)/Homer	ANCAGGAA GT	0			
Human Brain	Microglia	EHF(ETS)/LoVo-EHF-ChIP-Seq(GSE49402)/Homer	AVCAGGAA GT	0			
Human Brain	Microglia	Elk1(ETS)/Hela-Elk1-ChIP-Seq(GSE31477)/Homer	HACTTCCG GY	0			
Human Brain	Microglia	ETV4(ETS)/HepG2-ETV4-ChIP-Seq(ENCODE)/Homer	ACCGGAAG TG	0			

Human Brain	Microglia	Elk4(ETS)/Hela-Elk4-ChIP-Seq(GSE31477)/Homer	NRYTTCCGY	0			
Human Brain	Microglia	IRF2(IRF)/Erythroblasts-IRF2-ChIP-Seq(GSE36985)/Homer	GAAASYGA AASY	0			
Human Brain	Microglia	EWS:FLI1-fusion(ETS)/SK_N_MC-EWS:FLI1-ChIP-Seq(SRA014231)/Homer	VACAGGAA AT	0			
Human Brain	Microglia	IRF8(IRF)/BMDM-IRF8-ChIP-Seq(GSE77884)/Homer	GRAASTGA AAST	0.0001			
Human Brain	Microglia	IRF3(IRF)/BMDM-Irf3-ChIP-Seq(GSE67343)/Homer	AGTTTCAKT TTC	0.0001			
Human Brain	Microglia	ETS(ETS)/Promoter/Homer	AACCGGAA GT	0.0001			
Human Brain	Microglia	PU.1-IRF(ETS:IRF)/Bcell-PU.1-ChIP-Seq(GSE21512)/Homer	MGGAAGTG AAAC	0.0002			
Human Brain	Microglia	ELF1(ETS)/Jurkat-ELF1-ChIP-Seq(SRA014231)/Homer	AVCCGGAA GT	0.0002			

Human Brain	Microglia	bZIP50(bZIP)/colamp-bZIP50-DAP-Seq(GSE60143)/Homer	GATGACGTCA	0.0002			
Human Brain	Microglia	PU.1:IRF8(ETS:IRF)/pDC-Irf8-ChIP-Seq(GSE66899)/Homer	GGAAGTGA AAST	0.0004			
Human Brain	Microglia	IRF1(IRF)/PBMC-IRF1-ChIP-Seq(GSE43036)/Homer	GAAAGTGA AAGT	0.0005			
Human Brain	Microglia	EWS:ERG-fusion(ETS)/CADO_ES1-EWS:ERG-ChIP-Seq(SRA014231)/Homer	ATTCCTGT N	0.0013			
Human Brain	Microglia	SPDEF(ETS)/VCaP-S PDEF-ChIP-Seq(SRA014231)/Homer	ASWTCCTG BT	0.0013			
Human Brain	Microglia	Fra2(bZIP)/Striatum-Fra2-ChIP-Seq(GSE43429)/Homer	GGATGACT CATC	0.004			
Human Brain	Microglia	ETS:RUNX(ETS,Runt)/Jurkat-RUNX1-ChIP-Seq(GSE17954)/Homer	RCAGGATG TGGT	0.0055			

Human Brain	Microglia	Foxo1(For khead)/RA W-Foxo1-ChIP-Seq(Fan_et_al.)/Homer	CTGTTTAC	0.0085			
Human Brain	Microglia	Nrf2(bZIP)/Lymphoblast-Nrf2-ChIP-Seq(GSE37589)/Homer	HTGCTGAGTCAT	0.0102			
Human Brain	Microglia	Fra1(bZIP)/BT549-Fra1-ChIP-Seq(GSE46166)/Homer	NNATGASTCATH	0.0102			
Human Brain	Microglia	IRF4(IRF)/GM12878-IRF4-ChIP-Seq(GSE32465)/Homer	ACTGAAACCA	0.0103			
Human Brain	Microglia	TGA6(bZIP)/colamp-TGA6-DAP-Seq(GSE60143)/Homer	TGACGTCA BC	0.0119			
Human Brain	Microglia	MafA(bZIP)/Islet-MafA-ChIP-Seq(GSE30298)/Homer	TGCTGACTCA	0.0172			
Human Brain	Microglia	Ets1-distal(ETS)/CD4+-PolII-ChIP-Seq(Bar ski_et_al.)/Homer	MACAGGAA GT	0.0206			
Human Brain	Microglia	RAP211(A P2EREBP)/colamp-RAP211-DA	RGCCGGCY WW	0.0373			

		P-Seq(GS E60143)/Homer					
Human Brain	Microglia	ISRE(IRF)/ThioMac-L PS-Expression(GSE23622)/Homer	AGTTTCATTC	0.0429			
Human Brain	Microglia				EWS:ERG-fusion(ETS)/CADO_ES1-EWS:ERG-ChIP-Seq(SRA014231)/Homer	ATTCCTGTN	0.0001
Human Brain	Microglia				ELF5(ETS)/T47D-ELF5-ChIP-Seq(GSE30407)/Homer	ACVAGGAA GT	0.0001
Human Brain	Microglia				ETV1(ETS)/GIST48-ETV1-ChIP-Seq(GSE22441)/Homer	AACCGGAA GT	0.0001
Human Brain	Microglia				PU.1(ETS)/ThioMac-PU.1-ChIP-Seq(GSE21512)/Homer	AGAGGAA GTG	0.0001
Human Brain	Microglia				Fli1(ETS)/C8-D8-FLI-ChIP-Seq(GSE20898)/Homer	NRYTTCCGH	0.0001
Human Brain	Microglia				ERG(ETS)/VCaP-ERG-ChIP-Seq(GSE14097)/Homer	ACAGGAAG TG	0.0001
Human Brain	Microglia				Etv2(ETS)/ES-ER71-ChIP-Seq(GS	NNAYTTCC TGHN	0.0001

					E59402)/Homer		
Human Brain	Microglia				GABPA(ETS)/Jurkat-ABPa-ChIP-Seq(GSE17954)/Homer	RACCGGAA GT	0.0003
Human Brain	Microglia				ETS1(ETS)/Jurkat-ETS1-ChIP-Seq(GSE17954)/Homer	ACAGGAAG TG	0.0003
Human Brain	Microglia				EHF(ETS)/LoVo-EHF-ChIP-Seq(GSE49402)/Homer	AVCAGGAA GT	0.0003
Human Brain	Microglia				ETV4(ETS)/HepG2-ETV4-ChIP-Seq(ENCODE)/Homer	ACCGGAA GTG	0.0004
Human Brain	Microglia				JunB(bZIP)/DendriticCells-Junb-ChIP-Seq(GSE36099)/Homer	RATGASTC AT	0.0005
Human Brain	Microglia				Elf4(ETS)/MDM-Elf4-ChIP-Seq(GSE88699)/Homer	ACTTCKKG KT	0.0007
Human Brain	Microglia				ELF3(ETS)/PDAC-ELF3-ChIP-Seq(GSE64557)/Homer	ANCAGGAA GT	0.0008
Human Brain	Microglia				Fosl2(bZIP)/3T3L1-Fosl2-ChIP-Seq(GSE56872)/Homer	NATGASTC ABNN	0.0072

Human Brain	Microglia				Fos(bZIP)/TSC-Fos-ChIP-Seq(GSE110950)/Homer	NDATGASTCAYN	0.0082
Human Brain	Microglia				EWS:FLI1-fusion(ETS)/SK_N_MC-EWS:FLI1-ChIP-Seq(SRA014231)/Homer	VACAGGAAAT	0.0095
Human Brain	Microglia				Fra1(bZIP)/BT549-Fra1-ChIP-Seq(GSE46166)/Homer	NNATGASTCATH	0.0152
Human Brain	Microglia				ELF1(ETS)/Jurkat-ELF1-ChIP-Seq(SRA014231)/Homer	AVCCGGAA GT	0.0152
Human Brain	Microglia				AP-1(bZIP)/ThioMac-P.U.1-ChIP-Seq(GSE21512)/Homer	VTGACTCATC	0.0153
Human Brain	Microglia				Fra2(bZIP)/Striatum-Fra2-ChIP-Seq(GSE43429)/Homer	GGATGACTCATC	0.0215
Human Brain	Microglia				BATF(bZIP)/Th17-BATF-ChIP-Seq(GSE39756)/Homer	DATGASTCAT	0.0255
Human Brain	Microglia				Ets1-distal(ETS)/CD4+-PolII-ChIP-Seq(Barski_et_al.)/Homer	MACAGGAAGT	0.0319
Human Brain	Microglia				Atf3(bZIP)/GBM-ATF3-	DATGASTCATHN	0.0323

					ChIP-Seq( GSE33912) /Homer		
Human Brain	Microglia				bZIP69(bZIP)/col-bZIP69-DAP-Seq(GSE60143)/Homer	GACAGCTG KCAW	0.0366
Human brain	Inhibitory_neuron_MAF	Fra1(bZIP)/BT549-Fra1-ChIP-Seq(GSE46166)/Homer	NNATGAST CATH	0.0443			
Human brain	Inhibitory_neuron_MAF	Fos(bZIP)/TSC-Fos-ChIP-Seq(GSE110950)/Homer	NDATGAST CAYN	0.0443			
Human brain	Inhibitory_neuron_MAF				Ap4(bHLH)/AML-Tfap4-ChIP-Seq(GSE45738)/Homer	NAHCAGCT GD	0.0014
Human brain	Inhibitory_neuron_MAF				MyoG(bHLH)/C2C12-MyoG-ChIP-Seq(GSE36024)/Homer	AACAGCTG	0.0014
Human brain	Inhibitory_neuron_MAF				Myf5(bHLH)/GM-Myf5-ChIP-Seq(GSE24852)/Homer	BAACAGCT GT	0.0014
Human brain	Inhibitory_neuron_MAF				Tcf21(bHLH)/ArterySmoothMuscle-Tcf21-ChIP-Seq(GSE61369)/Homer	NAACAGCT GG	0.0018
Human brain	Inhibitory_neuron_MAF				CTCF(Zf)/C4D4+-CTCF-ChIP-Seq(B	AYAGTGCC MYCTRGTGCCA	0.0018

					arski_et_al.)/Homer		
Human brain	Inhibitory_neuron_MAF				Fra1(bZIP)/BT549-Fra1-ChIP-Seq(GSE46166)/Homer	NNATGASTCATH	0.0018
Human brain	Inhibitory_neuron_MAF				MyoD(bHLH)/Myotube-MyoD-ChIP-Seq(GSE21614)/Homer	RRCAGCTGYTSY	0.0018
Human brain	Inhibitory_neuron_MAF				SCL(bHLH)/HPC7-Scl-ChIP-Seq(GSE13511)/Homer	AVCAGCTG	0.0018
Human brain	Inhibitory_neuron_MAF				BATF(bZIP)/Th17-BATF-ChIP-Seq(GSE39756)/Homer	DATGASTCAT	0.0084
Human brain	Inhibitory_neuron_MAF				BORIS(Zf)/K562-CTCF-L-ChIP-Seq(GSE32465)/Homer	CNNBRGC GCCCCCT GSTGGC	0.0084
Human brain	Inhibitory_neuron_MAF				Atf3(bZIP)/GBM-ATF3-ChIP-Seq(GSE33912)/Homer	DATGASTCATHN	0.009
Human brain	Inhibitory_neuron_MAF				Twist2(bHLH)/Myoblast-Twist2.Ty1-ChIP-Seq(GSE127998)/Homer	MCAGCTG BYH	0.0147
Human brain	Inhibitory_neuron_MAF				Tcf12(bHLH)/GM12878-Tcf12-ChIP-Seq(GSE32465)/Homer	VCAGCTGYTG	0.0151

Human brain	Inhibitory_neuron_MAF				Ascl1(bHLH)/NeuralTubes-Ascl1-ChIP-Seq(GSE55840)/Homer	NNVVCAGC TGBN	0.0168
Human brain	Inhibitory_neuron_MAF				Fos(bZIP)/TSC-Fos-ChIP-Seq(GSE110950)/Homer	NDATGAST CAYN	0.0169
Human brain	Inhibitory_neuron_MAF				Barx1(Homobox)/Stomach-Barx1.3xFlag-ChIP-Seq(GSE69483)/Homer	AAACMATT AN	0.0225
Human brain	Inhibitory_neuron_MAF				LHX9(Homobox)/Hct116-LHX9.V5-ChIP-Seq(GSE116822)/Homer	NGCTAATT AG	0.0305
Human brain	Inhibitory_neuron_MAF				JunB(bZIP)/DendriticCells-Junb-ChIP-Seq(GSE36099)/Homer	RATGASTC AT	0.0305
Human brain	Inhibitory_neuron_MAF				Nkx6.1(Homobox)/Islet-Nkx6.1-ChIP-Seq(GSE40975)/Homer	GKTAATGR	0.0409
PBMC	DCp				RUNX1(Runt)/Jurkat-RUNX1-ChIP-Seq(GSE29180)/Homer	AAACCACAR M	0
PBMC	DCp				RUNX-AML(Runt)/CD4+-PolII-ChIP	GCTGTGGT TW	0

					P-Seq(Barski_et_al.)/Homer		
PBMC	DCp				RUNX(Runt)/HPC7-Runx1-ChIP-Seq(GSE22178)/Homer	SAAACCACAG	0.0005
PBMC	DCp				Fli1(ETS)/CD8-FLI-ChIP-Seq(GSE20898)/Homer	NRYTTCCGGH	0.0018
PBMC	DCp				RUNX2(Runt)/PCa-RUNX2-ChIP-Seq(GSE33889)/Homer	NWAACCA CADNN	0.002
PBMC	DCp				ERG(ETS)/VCaP-ERG-ChIP-Seq(GSE14097)/Homer	ACAGGAAG TG	0.0046
PBMC	DCp				ETS1(ETS)/Jurkat-ETS1-ChIP-Seq(GSE17954)/Homer	ACAGGAAG TG	0.0073
PBMC	DCp				EWS:ERG-fusion(ETS)/CADO_ES1-EWS:ERG-ChIP-Seq(SRA014231)/Homer	ATTTCTGTN	0.028
PBMC	DCp				Etv2(ETS)/ES-ER71-ChIP-Seq(GSE59402)/Homer	NNAYTTCC TGHN	0.028

910 Extended Table 1: complete list of motifs identified only by either XAI or differential analysis.