**1** **A Novel Variable Neighborhood Search Approach for Cell Clustering for Spatial**
**2** **Transcriptomics**

**3**

**4** **Authors:**

**5** Aleksandra Djordjevic[1,*], aleksandradjordjevic@genomics.cn
**6** Junhua Li[1,3], lijunhua@genomics.cn
**7** Shuangsang Fang[1,2], fangshuangsang@genomics.cn
**8** Lei Cao[1,2], caolei2@genomics.cn
**9** Marija Ivanovic[1,*], marijaivanovic@genomics.cn

**10**

**11** 1 BGI Research, Shenzhen 518083, China
**12** 2 BGI Research, Beijing 102601, China
**13** 3 BGI Research, Lidostas Parks, Riga 49276, Latvia

**14**

**15** * contributed equally, corresponding authors.

**16** ORCID IDs:
**17** Marija Ivanovic [0000-0002-3372-8271], Aleksandra Djordjevic [0009-0002-9450-4121], Junhua Li
**18** [0000-0001-6784-1873], Shuangsang Fang [0000-0002-4126-0074], Lei Cao [0000-0002-7170-9602];

**19**

**20** **Abstract**

**21** This paper introduces a new approach to cell clustering using the Variable Neighborhood Search (VNS)
**22** metaheuristic. The purpose of this method is to cluster cells based on both gene expression and spatial
**23** coordinates. Initially, we confronted this clustering challenge as an Integer Linear Programming
**24** minimization problem. Our approach introduced a novel model based on the VNS technique,
**25** demonstrating the efficacy in navigating the complexities of cell clustering. Notably, our method
**26** extends beyond conventional cell-type clustering to spatial domain clustering. This adaptability
**27** enables our algorithm to orchestrate clusters based on information gleaned from gene expression
**28** matrices and spatial coordinates. Our validation showed the superior performance of our method
**29** when compared to existing techniques. Our approach advances current clustering methodologies and
**30** can potentially be applied to several fields, from biomedical research to spatial data analysis.

**31** **Subject areas:** Software and Workflows, Bioinformatics, Transcriptomics.

**32**

## Statement of Need

**34**

**35** In high-throughput omics, deciphering the intricate cellular dynamics within tissues is pivotal
**36** [1,2]. Cell clustering is essential for dissecting the mosaic of cellular diversity [3,4]. This analytical
**37** approach seeks to categorize individual cells based on shared molecular signatures, allowing the
**38** identification of discrete subpopulations within heterogeneous tissues. In exploring cellular behavior
**39** and function, cell clustering emerges as an indispensable tool, providing insights into the subtle
**40** nuances of gene expression profiles. The ability to stratify cells into meaningful clusters not only refines

41 our understanding of tissue composition but also lays the groundwork for precise insights into disease
42 etiology and potential therapeutic interventions.
43     In tandem with cell clustering, spatial transcriptomics [5,6] constitutes a revolutionary frontier
44 for understanding cellular dynamics with their native microenvironments. Beyond the traditional
45 scope of genomics, spatial transcriptomics integrates the spatial context of cells into the analysis,
46 allowing researchers to explore how gene expression patterns unfold across complex tissue structures.
47 This multidimensional approach surpasses the limitations of conventional transcriptomic studies,
48 providing a spatially resolved perspective that is indispensable for decoding the orchestration of
49 cellular interactions and the emergence of tissue-specific functions.
50     In order to contribute to this dynamic landscape, we introduce a novel methodology rooted in
51 the Variable Neighborhood Search approach [7]. Our innovation seeks to elevate the precision and
52 efficacy of cell clustering in spatial transcriptomic analyses, promising to reveal hidden facets of cellular
53 organization and functionality. In this work, we introduce a novel Variable Neighborhood Search (VNS)
54 approach tailored for cell clustering in spatial transcriptomics. Although our initial investigations
55 focused on datasets designed for cell-type clustering, it is essential to emphasize that our method's
56 design accommodates spatial domain clustering as well. Here, we present a synthesis of computational
57 skills and biological insights aimed at pushing the boundaries of our understanding of the complex cell
58 interactions within tissues.
59

60 ## Background
61 ## Clustering methods from the literature
62
63 Many methods in the literature can be used to partition an $N$-dimensional population into $K$ sets
64 based on specific rules. In this paper, we focus on some of the most popular clustering methods used
65 in the field of data analysis, such as $k$-Means [8], Louvain [9], Leiden [10], and MClust [11]. While these
66 methods share the goal of grouping data points, they differ in the types of data they are designed for,
67 the principle they optimize, and the algorithms they are well-suited for. $k$-Means is a general-purpose
68 clustering algorithm, Louvain and Leiden are tailored for community detection in networks, while
69 MClust is a model-based clustering method. In the following subsections, we briefly describe each of
70 these methods.
71

72 ### *$k$-Means algorithm*
73 The $k$-means algorithm [8] is a partitioning algorithm that divides a dataset into $k$-clusters based on
74 the similarity of data points. It starts by establishing $k$ groups, each comprising a singular randomly
75 chosen point. Points are then added to these groups according to the principle that new points are
76 assigned to the group whose mean point is the most similar by some rule. After point allocation, the
77 means of all groups are adjusted to incorporate the influence of newly added points. Consequently, at
78 each stage, the $k$-means are reflective of the means of the groups they represent.
79     While this method is computationally efficient and adeptly handles extensive datasets, it does
80 not guarantee convergence to an optimal solution. Notably, issues arise from the random initialization
81 of centroids, leading to unexpected convergence patterns. Moreover, the algorithm requires users to
82 choose the cluster number beforehand, influencing cluster shapes and susceptibility to outlier effects.
83 However, it is known that certain special cases of the $k$-means algorithm exist in the literature where
84 convergence to an optimal solution is assured.
85

##### Louvain algorithm

The Louvain algorithm, developed by V. D. Vondel et al. [9], is designed for detecting communities in network or graph data. This algorithm aims to optimize modularity, a measure of the quality of network division into communities, using two phases: (1) local moving of nodes and (2) aggregation of the network. In the first phase, individual nodes are moved to the community that yields the largest increase in the quality function. In the second phase, an aggregation network is obtained based on partitions, with each community in a partition becoming a node in the aggregate network. These two phases are repeated until the quality function cannot be increased further. However, the Louvain algorithm can potentially produce communities with arbitrarily poor connectivity. In the most adverse scenarios, these communities may become entirely disconnected, particularly during iterative executions of the algorithm.

##### Leiden algorithm

To address the connectivity issues of the Louvain algorithm, V. A. Traag et al. introduced the Leiden algorithm [10]. The Leiden algorithm guarantees that communities are well connected and, when applied iteratively, the algorithm converges to a partition where all subsets of all communities are locally optimally assigned. The Leiden algorithm is partly based on the smart local move algorithm, which itself can be seen as an improvement of the Louvain algorithm and takes advantage of the idea of speeding up the local moving of nodes and the idea of moving nodes to random neighbors, the Leiden algorithm considers these ideas to represent the most promising directions in which the Louvain algorithm can be improved. The Leiden algorithm consists of three phases: (1) local moving of nodes, (2) refinement of the partition, and (3) aggregation of the network based on the refined partition, using the non-refined partition to create an initial partition for the aggregate network. Thus, this algorithm optimizes a quality function to identify communities by considering the density of connections within the communities.

##### MClust

MClust [11], applied in cell clustering, identifies distinct cell groups based on observed features using Gaussian mixture models [12]. Unlike other clustering algorithms, MClust accommodates various cluster shapes, making it suitable for complex situations. It utilizes the Expectation-Maximization [13] algorithm for parameter estimation, offering robust handling of missing data and complex distributions. This model-based clustering tool is powerful in uncovering patterns within complex biological datasets, such as those from single-cell omics technologies. Initially designed for single-cell RNA sequencing data, it can also be applied to spatial transcriptomic data, its effectiveness depending on data characteristics and analysis goals.

## Embedding methods from the literature

In spatial transcriptomics, where data is organized as a matrix with cells and genes, the high dimensionality (often exceeding 30,000 genes) and sparsity pose analytical challenges. Dimensionality reduction methods play key roles in addressing these issues. These techniques help distill meaningful patterns from the data, facilitating more efficient analyses.

The generation of embeddings, achieved through established literature methods, aims to transform the high-dimensional gene space into a more manageable form. This process enables a clearer exploration of spatial relationships, cell heterogeneity, and underlying biological processes. By leveraging validated methods from existing literature, we ensure a scientifically rigorous approach,

132 condensing rich gene expression profiles into interpretable embeddings while addressing
133 computational complexities.
134      As mentioned previously, we performed dimensionality reduction using five different
135 embedding methods: STAGATE [14], Principal Component Analysis (PCA) [15], GraphST [16], Cell
136 Clustering for Spatial Transcriptomics (CCST) data [17], and STAligner [18].
137

138 *STAGATE*
139 The STAGATE method [14] has been designed for spatial clustering and denoising in spatially resolved
140 transcriptomics data. This method generates low-dimensional latent embeddings with both spatial
141 information and gene expressions via a graph attention auto-encoder. Notably, the method adopts an
142 attention mechanism in the middle layer of the encoder and decoder, which learns the edge weights
143 of spatial neighbor networks and uses them to update spot representations by collectively aggregating
144 information from their neighbors.
145

146 *Principal Component Analysis*
147 PCA [15] is a statistical method for dimensionality reduction and data visualization. It is a mathematical
148 procedure that transforms a set of correlated variables into a new set of uncorrelated variables known
149 as principal components. The principal components are linear combinations of the original variables
150 and are sorted based on how much they account for the variance within the data; i.e., the first principal
151 component accounts for the highest variance. PCA finds widespread application across domains,
152 including data analysis, machine learning, and image processing, aiming to streamline intricate
153 datasets and uncover patterns or associations between variables.
154

155 *GraphST*
156 GraphST [16] is an advanced self-supervised contrastive learning technique designed to maximize the
157 potential of spatial transcriptomics data. Integrating graph neural networks with self-supervised
158 contrastive learning, this method acquires spot representations that are both informative and
159 distinctive. This is achieved by minimizing the embedding distance between spatially neighboring spots
160 reciprocally.
161

162 *Cell Clustering for Spatial Transcriptomics data*
163 CCST [17] leverages graph convolutional networks (GCNs) to integrate gene expression data and
164 comprehensive spatial information from individual cells in spatial gene expression data. The
165 relationships between variables are captured as a graph, with the adjacency matrix representing
166 connections among variables and the node feature matrix reflecting variable observations. The GCN
167 layer is strategically designed to fuse graph (in our case, spatial structure) and node features (gene
168 expression). Initially, the data is transformed into a graph, where nodes represent cells with gene
169 expression profiles as attributes, and edges represent neighborhood relationships between cells.
170 Subsequently, a sequence of GCN layers is used to incorporate graph and gene expression details into
171 cell node embedding vectors. Concurrently, the graph is perturbed to generate negative embeddings.
172 By learning the discrimination task, the neural network model is trained to encode cell embeddings
173 derived from spatial gene expression data, subsequently used for cell clustering.
174

175 *STAligner*
176 STAligner [18] is a specialized tool for aligning and integrating spatially-resolved transcriptomics data.
177 It begins by normalizing expression profiles for all spots and creating a spatial neighbor network based
178 on spatial coordinates. Employing a graph attention auto-encoder neural network, STAligner extracts

179 spatially-aware embeddings and uses spot triplets to guide the alignment process, fostering similarity
180 among related spots and distinction among dissimilar ones across slices. The introduction of triplet
181 loss refines spot embeddings by minimizing the distance from the anchor to positive spots and
182 increasing the distance to negative spots. This iterative process optimizes triplet construction and auto-
183 encoder training until batch-corrected embeddings are obtained. Furthermore, STAligner's versatility
184 extends to integrating spatial transcriptomics datasets, facilitating alignment and concurrent
185 identification of spatial domains across diverse biological samples, technological platforms,
186 developmental stages, disease conditions, and consecutive tissue slices for 3D alignment.
187

## Implementation

### Mathematical model

190

191 Let $C = [c_i]$ represent the set of cells $c_i$, $i = 1, \dots, n$, and the total number of cells equal $n$. For each
192 cell $c_i$, $i = 1, \dots, n$, let $c_i^x$ and $c_i^y$ represent its $x$ and $y$ coordinates, and let vector $c_i^{emb} =$
193 $[c_i^{emb_1}, \dots, c_i^{emb_M}]$ represent embedding values ($M$ is the total number of embedding values).
194 Furthermore, let the distance function $D: C \times C \to \mathcal{R}^+$ be defined as a measure of the similarity
195 between the cells. In our model, for two cells $c_i$ and $c_j$, the distance $D$ was calculated as follows: $D =$
196 $\alpha \, D_{gene} + (1 - \alpha)D_{coord}$, where $\alpha$ is the input parameter, $D_{gene}$ is the cosine similarity between cell
197 embeddings, and $D_{coord}$ is the Euclidian distance between cell coordinates:

198 $$D_{gene}(c_i, c_j) = cosine(c_i, c_j),$$

199 $$D_{coord}(c_i, c_j) = \sqrt{\left(c_i^x - c_j^x\right)^2 + \left(c_i^y - c_j^y\right)^2}.$$

200 In our model, we chose $K$ different cells from the set of cells $C$ to represent clusters and called these
201 cells *centroids*. Therefore, let the binary variables $x_{ij}$ ($i, j = 1, \dots, n$) and $y_i$ be defined in the following
202 way:

203 $$x_{ij} = \begin{cases} 1, & \text{if cell } c_i \text{ belongs to the cluster represented by centroid } c_j \\ 0, & \text{otherwise} \end{cases}$$

204 $$y_i = \begin{cases} 1, & \text{if cell } c_i \text{ represents the centorid} \\ 0, & \text{otherwise} \end{cases}$$

205 The Integer Linear Programming formulation of the clustering problem can be described as follows:

206 $$\min \sum_{i=1}^{n} \sum_{j=1}^{n} x_{ij} D(c_i, c_j) \qquad (1)$$

207 subject to these constraints:

208 $$\sum_{i=1}^{n} x_{ij} = 1, \ 1 \le j \le n, \qquad (2)$$

209 $$x_{ij} \le q_j, \ 1 \le i \le n, \qquad 1 \le j \le n, \quad (3)$$

210 $$\sum_{i=1}^{n} y_i = K, \qquad (4)$$

211 $$x_{ij}, y_j \in \{0,1\}, \ 1 \le i \le n, \ 1 \le j \le n. \quad (5)$$

212 The objective function (1) represents the sum of distances from each cell to its most similar cluster
213 representative. This function should be minimized. Equation (2) indicates that each cell is assigned to
214 only one cluster. Before assigning a cell to a cluster, the cluster needs to be defined (3). The total
215 number of clusters is equal to $K$ (4). All variables are constrained to be binary (5).
216 The model described with equations (1)-(5) is based on the $p$-median classification and is
217 presented in a similar form by Davidović et al. [19].
218

## Variable Neighborhood Search Method

The VNS method is a well-known metaheuristic method. It starts from one point in the search space, explores its neighborhoods, and repeats the process until a better solution or stopping criteria are reached. This method was proposed for the first time by Mladenović [20] and later elaborated by Mladenović and Hansen [21] and Hansen and Mladenović [22].

Before we introduce the VNS method, let us define the set $N_k(X)$, $k = k_{\{min\}}, \dots, k_{max}$ as the set of all vectors $X'$ that have a difference of the $k^{th}$ order from the solution $X$, and call that set $k^{th}$ Neighborhood to the solution X.

The VNS-based heuristic can be defined in a way that it starts from the initial feasible solution $X$, shakes it by creating another solution $X' \in N_k(X)$, and then applies a local search method to create a better feasible solution $X''$. If the feasible solution $X''$ obtained by the local search procedure is not better than the current incumbent $X$ ($F(X'') \geq F^*$), the VNS algorithm repeats the procedure of shaking in the neighborhood $N_{k+k_{step}}$ (i.e., $k$ is incremented by $k_{step}$) and local searches within it. It repeats this passage until $k$ reaches its maximum $k_{max}$. Otherwise, if $F(X'') < F^*$, $F^*$ becomes $F(X'')$ and $k$ becomes $k_{min}$. The procedure of changing the neighborhood enables the VNS algorithm to get out from the local minima. The process is repeated until a certain number of iterations or other stop criteria are reached.

Pseudo-code for the basic VNS algorithm is presented as Algorithm 1. Implementations of the functions *InitialSolution(), Shake(), LocalSearch(),* and *StoppingCondition()* defined for our clustering problem are described in the following subsection.

---

**Algorithm 1** (Basic) Variable Neighborhood Search Method

---

1: $X^* \leftarrow InitialSolution()$
2: $F^* \leftarrow F(X^*)$
3: **while** *StoppingCondition()* **do**
4: $\quad k \leftarrow k_{min}$
5: $\quad$ **while** $k \leq k_{max}$ **do**
6: $\quad\quad X \leftarrow X^*$
7: $\quad\quad X' \leftarrow Shake(X, k)$
8: $\quad\quad X'' \leftarrow LocalSearch(X')$
9: $\quad\quad$ **if** $F(X'') < F^*$ **then**
10: $\quad\quad\quad X^* \leftarrow Y''$
11: $\quad\quad\quad F^* \leftarrow F(X^*)$
12: $\quad\quad\quad k \leftarrow k_{min}$
13: $\quad\quad$ **else**
14: $\quad\quad\quad k \leftarrow k + k_{step}$
15: $\quad\quad$ **end if**
16: $\quad$ **end while**
17: **end while**

---

242

## VNS for the cell clustering problem

244

245 With respect to the problem's definition, let us assume that all cells can be represented by numbers
246 from 1 to $n$. Specifically, cells can be represented by the set $C = [c_i]$, $n = |C|$, and that for each cell
247 $c_i$ there are two types of data: the $x$ and $y$ coordinates of the cell ($c_i^x$ and $c_i^y$) and the embedding
248 values (vector $emb_i$). In our representation, the solution vector $Y = [y_1, ..., y_K]$ contains indexes of $K$
249 cells chosen as cluster representatives. Also, cell $y_i$ is a centroid of the $i$-th cluster. From the centroid
250 solution vector $Y$ we obtain vector $X = [x_i]$ of size $n$ in the following way: $x_i$, $i = 1, ..., n$, represents
251 the closest centroid from the $Y$ vector to the $i$-th cell. Our representation satisfies all conditions
252 described by equations (2) - (5). Using this representation, our goal was to minimize the value of the
253 function $F: C \times C \to \mathcal{R}^+$, where $F$ is defined as $F(X) = \sum_{i=1}^{n} \left( \alpha D_{gene}(i, x_i) + (1 -$
254 $\alpha) D_{coord}(i, x_i) \right)$.

255 The function *InitialSolution()* randomly chooses $K$ mutually different numbers from the set of
256 numbers $\{1, ..., n\}$ and returns them as a $K$-dimensional vector $Y$. For every solution vector $Y$, vector
257 $X$ is obtained in the following way: for each cell $i$, the distance $D$ between the cell $i$ and all centroids
258 $y_j$ from the vector $Y$ is calculated; next, $x_i$ is set equal to the $y_j$ for which the distance $D$ is minimal.
259 That is, whenever the vector $Y$ is changed, vector $X$ is also updated. Also, to avoid repeated
260 calculations, the distance $D$ between all cells is calculated and saved as a *distance* matrix.

261 The *Shake()* function takes two inputs: the incumbent $Y$ and the size $k$ of the neighborhood
262 that needs to be explored. As a result, the *Shake()* function randomly chooses $k$ elements from the
263 vector $Y$ and replaces them with $k$ randomly chosen elements from the set $\{1, ..., n\}$ that are different
264 from all elements from the current $Y$. This means that when some elements are changed, all elements
265 in vector $Y$ will still be mutually different. In other words, the $Shake()$ function chooses a vector $Y'$
266 from $N_k(Y)$.

267 The *LocalSearch()* function takes vector $Y'$, the distance matrix $distance$, and the parameters
268 $m$ and $p$ as inputs. In our implementation, we used *the first improvement strategy*. Based on the value
269 of the parameter $m$, for each element of the vector $Y'$, the *LocalSearch()* function first chooses a
270 random integer number $ind \in [0, m]$; next, based on the $ind$ value, keeps the observed element of
271 the vector $Y'$ as it is ($ind == 0$) or replace it with the new one ($ind > 0$). For $ind \geq 2$, the observed
272 element is replaced with one of the candidates from the set of candidates that are created within the
273 *LocalSearch()* function (the *LocalSearch()* function searches for $ind$ candidates for which the $distance$
274 value from the observed candidate is the smallest, sorts the list, excludes all candidates that are
275 already present in the vector $Y'$, and then chooses one candidate for the replacement). Please note
276 that the smallest $distance$ value between the observed candidate and itself will be zero, so the
277 condition $ind > 1$ is necessary. In case $ind == 1$, $ind$ will be chosen again until its value is not equal
278 to 1. Additionally, if the candidate list is empty after excluding all elements that already exist in the
279 vector $Y'$, a random candidate will be chosen from the set $\{1, ..., n\} \setminus \{y_1, ..., y_K\}$.

280 Finally, after the procedure of replacing or keeping elements from the vector $Y'$ is finished, i.e., a new
281 vector $Y''$ is obtained, the $LocalSearch()$ function calculates $F(Y'')$ and, if $F(Y'') < F^*$, the first
282 improvement has been made, and the function returns the vector $Y''$ as the output or repeats the
283 whole process. The process of examining elements of the vector $Y'$ and replacing them with new
284 values is repeated only if no improvement is made, but not more than $p$ times. In case no improvement
285 is made and the process has been repeated p times, the vector $Y'' = Y'$ will be returned as the output
286 of this function.

287         In other words, the *LocalSearch()* function examines elements in the close neighborhood of
288 the observed vector $Y'$ by creating a new vector $Y''$, calculates the function value $F(Y'')$ and, if the
289 function value is less than the currently best value $F^*$, returns that vector. Otherwise, it will continue
290 the process of examining elements of the vector $Y'$ but not more than $p$ times.

291         Usually, the *StoppingCondition()* function checks if the maximal number of iterations
292 ($max_{iter}$) or the maximal running time ($t_{max}$) have been reached. In our code, the *StoppingCondition()*
293 function checks only if the maximal number of iterations has been reached and, if the answer is $true$,
294 returns the best solution found as the result of the VNS procedure. If the maximal number of iterations
295 has not been reached, the VNS procedure continues its search.

296

# Data Description

298

299 We assessed the performance of the clustering methods through quantitative evaluation, employing
300 datasets sourced from two distinct spatially resolved transcriptomic technologies: Stereo-seq [23] and
301 10x Visium [24].

302         From Stereo-seq technology, two datasets were used for testing: a large dataset of a field
303 mouse brain hemisphere (**SS200000128TR E2 benchmark**) and another from the dorsal midbrain
304 (**Forebrain**). The large field mouse brain contains more than 38,000 cells and more than 20,000 genes
305 and can be downloaded from [25], while Forebrain contains more than 18,000 cells and more than
306 23,000 genes and can be downloaded from [26]. Please note that Forebrain contains the whole dorsal
307 midbrain. In our study, we used manual lasso to separate a part of this dataset and called that part
308 **Forebrain**. Both datasets are composed of only one slice.

309         In order to evaluate the performance of the presented VNS method on multi-slice datasets,
310 we used a 10x Visium dataset containing spatial expressions of 12 human-layered dorsolateral
311 prefrontal cortex (DLPFC) sections. Since these 12 sections are from three different human donors,
312 they were used as multi-section (4-layers) datasets in our study. All layers of the DLPFC sections were
313 manually annotated by Maynard et al. [24] and can be downloaded from [27]. Viewing them as the
314 ground truth, we compared the clustering accuracy of the VNS method with other clustering methods
315 using only embedding obtained by the vertical spatial transcriptomic integration provided by
316 STAGATE.

# Analysis

## Input parameters

319

320 Testing was conducted on the AWS instance **m6a.48xlarge** under the Linux operative system.

321         Input parameters for our algorithm are the number of clusters ($K$), the percentage of the
322 influence of the embedding values ($\alpha$), the maximal number of neighborhoods that should be
323 searched ($k_{max}$,), the maximal number of iterations ($max_{iter}$), and the *local search* parameters $m$ and
324 $p$. The minimal ($k_{min}$) number of neighborhoods and step ($k_{step}$) are set to 1 by default.

325         The input parameters used for testing are $\alpha \in \{1, 0.95\}$ ($\alpha = 1$ means that no additional
326 spatial information is included, while $\alpha = 0.95$ means that 5% of spatial information is used for
327 calculating the distance between the cells), $k_{max} \in \{10,15,20,25,30\}$, $m \in \{10,12,15,20,30\}$, and
328 $p \in \{10,12,15,20\}$.

329

## Evaluation method

We used the Adjusted Rand Index (ARI) [28] to evaluate the results and compare them with each other. ARI is a measure used to evaluate the performance and similarity between two clustering algorithms. It quantifies the agreement between the true and predicted clustering, adjusting for the amount of agreement that could occur by chance. ARI values range from -1 to 1: where 1 indicates the perfect agreement, 0 indicates agreement expected by chance, and negative values suggest less agreement than expected by chance.

## Results of the VNS method across various scenarios with single-slice datasets

Due to the sparsity of the gene expression matrix and to ensure a fair comparison, embeddings were obtained using various methods from the literature (PCA, STAGATE, GraphST, and CCST) for both Stereo-seq datasets. Moreover, all methods create embedding that significantly reduces the number of genes to a much smaller set of features. For instance, the CCST method reduced the number of genes from the Forebrain dataset to 128 features, STAGATE to 64 features, PCA to 50 features, and GraphST to 20 features. For the E2 dataset, all parameters were the same except for STAGATE, where the number of features was lowered to 30. Hence, the input data depend on the number of cells and the number of obtained features (embeddings). The standard clustering methods from the literature ($k$-Means, MClust, Louvain, and Leiden) and the proposed VNS method for cell clustering were applied to the generated embeddings. The results of the testing are presented in Tables 1 and 2.

The goal of the VNS method was to find the solution with the smallest cost function, and we show these results in Table 1. Table 1 shows results obtained by the VNS method only and is organized as follows: the first column presents the name of the embeddings used as the input to the VNS method, while the following four columns ($f_{VNS}$, $t_{VNS}$, $err$, and $\sigma$) show the smallest cost function value, the corresponding running time, and the statistical analysis of all solutions obtained by VNS when comparing to the presented cost function value in that order. In other words, due to the stochastic nature of the metaheuristic, the VNS algorithm was run 20 times (for 20 different seeds) for each embedding, and information regarding the best solution value obtained in these 20 runs is provided in these four columns ($f_{VNS}$, $t_{VNS}$, $err$, and $\sigma$). More precisely, $f_{VNS}$ presents the minimal cost function value obtained after these 20 runs; $t_{VNS}$ is the corresponding running time for the presented solution value; $err$ and $\sigma$ contain additional information on the quality of the solution: $err$ is the average relative error of found solution from the presented one and is calculated as $err = \frac{1}{20}\sum_{i=1}^{20} err_i$, where $err\_i = |VNS_i - f_{VNS}|/|VNS_i|$, where $VNS_i$ is the VNS solution obtained in the $i^{th}$ run (seed). The value $\sigma$ is the standard deviation of $err$ and is calculated by $\sigma = \sqrt{\frac{1}{20}\sum_{i=1}^{20}(err_i - err)^2}$ . For each embedding method, the results obtained by VNS are presented in separate rows.

The results presented in Table 2 are organized into three groups. Similar to Table 1, the first column (first group) presents the name of the method used for creating the embedding. The next ten rows present the results for each clustering method separately; for each method, we provide the ARI score ($ARI$) and the running time ($t$) in seconds. The $ARI$ and $t$ values under the VNS columns stand for the best found $ARI$ score obtained for all testing combinations and the corresponding running time. The highest $ARI$ score achieved for some datasets among all clustering methods is highlighted in bold, while the second-best $ARI$ score is highlighted by an asterisk (*).

374 In both tables, the first set of results corresponds to the E2 dataset, and the next corresponds
375 to the Forebrain dataset. The E2 dataset results are visualized in Figure 1, while the Forebrain dataset
376 results are visualized in Figure 2.
377
378 **Table 1.** VNS solution for single-slice datasets. Values in columns $f_{VNS}, t_{VNS}, err$ and $\sigma$ are
379 obtained as explained in the Analysis section.
380

| Embedding | $f_{VNS}$ | $t_{VNS}$ (s) | $err$ | $\sigma$ |
|---|---|---|---|---|
| E2 | | | | |
| CCST | 1,019.7419 | 48.8355 | 0.1626 | 0.0476 |
| STAGATE | 2,706.7446 | 110.258 | 0.1196 | 0.0415 |
| PCA | 9,550.0142 | 79.1977 | 0.0320 | 0.0118 |
| GraphST | 10,083.5379 | 64.95 | 0.0197 | 0.0059 |
| Forebrain | | | | |
| CCST | 427.8511 | 47.8054 | 0.1579 | 0.0439 |
| STAGATE | 543.0947 | 52.7096 | 0.0925 | 0.0347 |
| PCA | 3,541.7886 | 50.1935 | 0.0214 | 0.0073 |
| GraphST | 2,209.235 | 92.0103 | 0.0473 | 0.0140 |

381
382 **Table 2**. Clustering method comparison for single-slice datasets. The highest *ARI* score
383 achieved for some datasets among all clustering methods is highlighted in bold, while the second-best
384 *ARI* score is highlighted by an asterisk (*).
385

| Embeddings | Leiden | | Louvain | | $k$-Means | | MClust | | VNS | |
|---|---|---|---|---|---|---|---|---|---|---|
| | *ARI* | *t* (s) | *ARI* | *t* (s) | *ARI* | *t* (s) | *ARI* | *t* (s) | *ARI* | *t* (s) |
| E2 | | | | | | | | | | |
| CCST | 0.1553 | 29.1638 | 0.1518 | 5.7702 | 0.1962* | 15.3243 | 0.1401 | 4,799.5287 | **0.2224** | 47.5667 |
| STAGATE | 0.1951 | 7.5198 | 0.2176 | 6.3803 | **0.2907** | 2.62854 | 0.2052 | 516.8929 | 0.2890* | 59.7737 |
| PCA | 0.0001 | 6.8347 | 0.1316 | 9.9780 | 0.2072* | 12.0037 | 0.2024 | 1,128.1911 | **0.2907** | 235.465 |
| GraphST | **0.0841** | 14.8255 | 0.0697* | 13.0344 | 0.0492 | 4.2599 | 0.0635 | 533.1441 | 0.0636 | 47.5184 |
| Forebrain | | | | | | | | | | |
| CCST | 0.0925 | 25.7164 | 0.0961* | 2.5659 | 0.1093 | 8.7788 | 0.0821 | 1,330.3455 | **0.1263** | 18.6987 |
| STAGATE | 0.1753 | 3.6952 | 0.1676 | 3.6263 | 0.1775* | 6.0085 | 0.1718 | 269.9742 | **0.2342** | 24.6907 |
| PCA | 0.1659 | 4.4805 | 0.1674* | 3.7720 | **0.1717** | 6.4302 | 0.1025 | 147.4443 | 0.1568 | 45.2866 |
| GraphST | 0.1738 | 3.8813 | 0.1847* | 4.6558 | 0.1833 | 1.8972 | 0.1709 | 73.0143 | **0.2104** | 9.2064 |

386
387

## VNS clustering achieves better results than other tested methods using the E2 dataset

389
390 From the first part of the results shown in Table 1, we can conclude that, using PCA embedding in all
391 20 runs, the values of the cost function are very close to the lowest cost function value ($err < 3.5$,
392 $\sigma < 1.5\%$). Using STAGATE, we have some differences, although $\sigma$ is still below 5% implying that the

393 VNS method is stable with both embeddings. The results of VNS clustering when the smallest cost
394 function values are reached are visualized in Figure 1a, while the results with the best ARI score
395 achieved by all clustering methods are shown in Figure 1b.

396

## VNS methods outperform other methods when clustering cells from the Forebrain dataset

397
398

399 By examining values from the $err$ and $\sigma$ columns in Table 1 for the Forebrain dataset, it can be easily
400 seen that differences between the results obtained in 20 runs are very small. In fact, the difference
401 between the best-found solution (the solution with the minimal cost function value) and the other 19
402 solutions is less than 5% (the average relative error $\sigma$ is less than 5%). This result means that the
403 solutions found in all 20 runs were very close to the smallest one. Also, from the results in the column
404 $t_{VNS}$, we can observe a running was less than 1 minute for three different embedding types and less
405 than 2 minutes for one embedding type.
406 Moreover, from the results presented in Table 2 for the Forebrain dataset, we can see that, in
407 the majority of cases, VNS had the highest $ARI$ score compared to the other methods (for three types
408 of embedding, the $VNS\ ARI$ score was the highest). Also, the running time was less than 1 minute for
409 each type of embedding. The only embedding for which the VNS did not find a solution with the best
410 $ARI$ score was the PCA one, and for this embedding, the best $ARI$ score was obtained by the $k$-Means
411 method.
412 By analyzing the results in Tables 1 and 2, we conclude that the VNS method achieves the best
413 $ARI$ score with the STAGATE embedding, and that in all 20 runs all solutions were close to the one
414 with the lowest cost function ($err < 1\%$). The results obtained with the minimal cost function and
415 the maximal $ARI$ score are visualized in Figure 2.

416
417
418

## VNS demonstrates a superior performance on multi-slice datasets

419
420

421 Next, we compared the clustering accuracy of the VNS method with other clustering methods by using
422 embeddings obtained by the STAligner method only. Compared to other embedding methods used for
423 single-slice datasets, it is worth mentioning that STAligner reduces the number of genes to 30 features.
424 The results of this comparison are presented in Tables 3 and 4. Table 3 is organized similarly to Table
425 1. The only difference is in the first column, which, in this case, is called Slice name. Since DLPFC
426 datasets are 4-layered slices, this column contains the names of the first and the last slices in this
427 particular dataset. Other slices imply. Thus, each row represents the results for one separate DLPFC
428 dataset.
429 Table 4 is organized similarly to Table 2; however, the column Embeddings is replaced by the
430 column Slice name, and the names of the first and the last slices from particular multi-slice datasets
431 are presented. Other slices imply. The results for each dataset are presented in separate rows, as in
432 Table 3. The results from Table 3 are visualized in Figure 3.
433 As we see from the columns $err$ and $\sigma$ in Table 3, in all 20 runs, the VNS method obtained
434 results similar to the ones with the smallest cost function ($err < 5.8\%$, $\sigma < 2.5\%$). Again, these results
435 imply that the method is stable even for multi-slice datasets. The fact that results from the columns
436 $t_{VNS}$ are smaller than 5 implies that this method can obtain results for four slices of these types of
437 datasets in less than 5 seconds.

438        From the results presented in Table 4, it can be concluded that the method proposed in this
439 paper outperforms other clustering methods in all aspects. Specifically, for each of the datasets we
440 tested, $ARI$ score was the highest and the running time was the lowest when the VNS method was
441 used.
442

443 **Table 3.** VNS solution for multi-slice datasets.

| Slice name | $f_{VNS}$ | $t_{VNS}$ | $err$ | $\sigma$ |
|---|---|---|---|---|
| 151507_151510 | 890.7088 | 4.2262 | 0.0884 | 0.0390 |
| 151669_151672 | 755.7133 | 2.8674 | 0.0866 | 0.0273 |
| 151673_151676 | 513.8781 | 1.1983 | 0.0923 | 0.0396 |

444

445 **Table 4**. Clustering method comparison for multi-slice datasets. The highest *ARI* score achieved for
446 some datasets among all clustering methods is highlighted in bold, while the second-best *ARI* score is
447 highlighted by an asterisk (*).

| Slice name | Leiden | | Louvain | | $k$-Means | | MClust | | VNS | |
|---|---|---|---|---|---|---|---|---|---|---|
| | *ARI* | *t* (s) | *ARI* | *t* (s) | *ARI* | *t* (s) | *ARI* | *t* (s) | *ARI* | *t* (s) |
| 151507_151510 | 0.3440 | 27.3778 | 0.4293* | 4.0119 | 0.3061 | 2.1001 | 0.3489 | 62.5176 | **0.4887** | 2.1094 |
| 151669_151672 | 0.4084 | 26.9197 | 0.4985* | 2.9611 | 0.2213 | 1.6839 | 0.4633 | 39.1007 | **0.6156** | 1.3014 |
| 151673_151676 | 0.4370 | 25.1056 | 0.4754* | 2.6766 | 0.3299 | 1.4413 | 0.4316 | 49.1890 | **0.5016** | 0.8573 |

448

## Discussion and Conclusion

450

451 Here, we introduced a novel approach suitable for clustering both single- and multi-slice spatial
452 transcriptomics datasets. This is the first application of a metaheuristic method, called the VNS, to the
453 clustering of spatial transcriptomic data. The essence of the VNS implementation presented in this
454 study is the utilization of a combinatorial/mathematical optimization algorithm; in this instance, a
455 metaheuristic approach. These methods are strategically designed to deliver sufficiently optimal
456 solutions to optimization and machine learning challenges while minimizing computational resources.
457 This approach is intended to offer a robust and computationally efficient solution for cell clustering in
458 spatial transcriptomics.

459 Our analysis demonstrated that the performance of clustering methods is significantly influenced by
460 the choice of embeddings and the way they were generated. Notably, the VNS approach combined
461 with PCA embeddings yields results that closely align with the ground truth, as illustrated in Figure 2b.
462 When benchmarked against existing techniques, our method consistently outperforms in terms of
463 efficiency and ARI scores. The algorithm's speed and stability are commendable, and its flexibility is
464 evidenced by a comprehensive set of parameters that can be tailored to meet diverse user
465 requirements. Future research will extend the method's application to time-series datasets and
466 explore additional VNS modifications and embedding techniques to enhance its utility.

467

## Availability of source code and requirements:

469
470 • Project name: VNS

471 • Project home page: https://github.com/STOmics/VNS/tree/main
472 • Operating system(s): Linux
473 • Programming language: Python
474 • License: MIT
475 • RRID:
476

# Data availability

479 From Stereo-seq technology, two datasets were used:
480 (1) a large dataset of a field mouse brain hemisphere (**SS200000128TR E2 benchmark**), which can be
481 downloaded from Zenodo [25]
482 (2) **Forebrain**, which can be downloaded from the CNGB MOSTA database
483 https://db.cngb.org/stomics/mosta/download/.
484 Additional data is also available in GigaDB [29]. We used only one part of Forebrain, which was
485 extracted using a manual lasso.

# Declarations

## Abbreviations

490 ARI, Adjusted Rand Index; CCST, Clustering for Spatial Transcriptomics; DLPFC, dorsolateral prefrontal
491 cortex; GCN, graph convolutional network; PCA, Principal Component Analysis; VNS, Variable
492 Neighborhood Search.

## Consent for publication

495 Not applicable.

## Competing Interests

498 The author(s) declare that they have no competing interests.

## Ethics approval and consent to participate

501 The authors declare that ethical approval was not required for this type of research.

## Funding

## Author's Contributions

507 AD and MI provided the idea of the solution, implementation, testing, and manuscript. JL and SF
508 supervised the whole process. CL created embeddings for both datasets for testing.
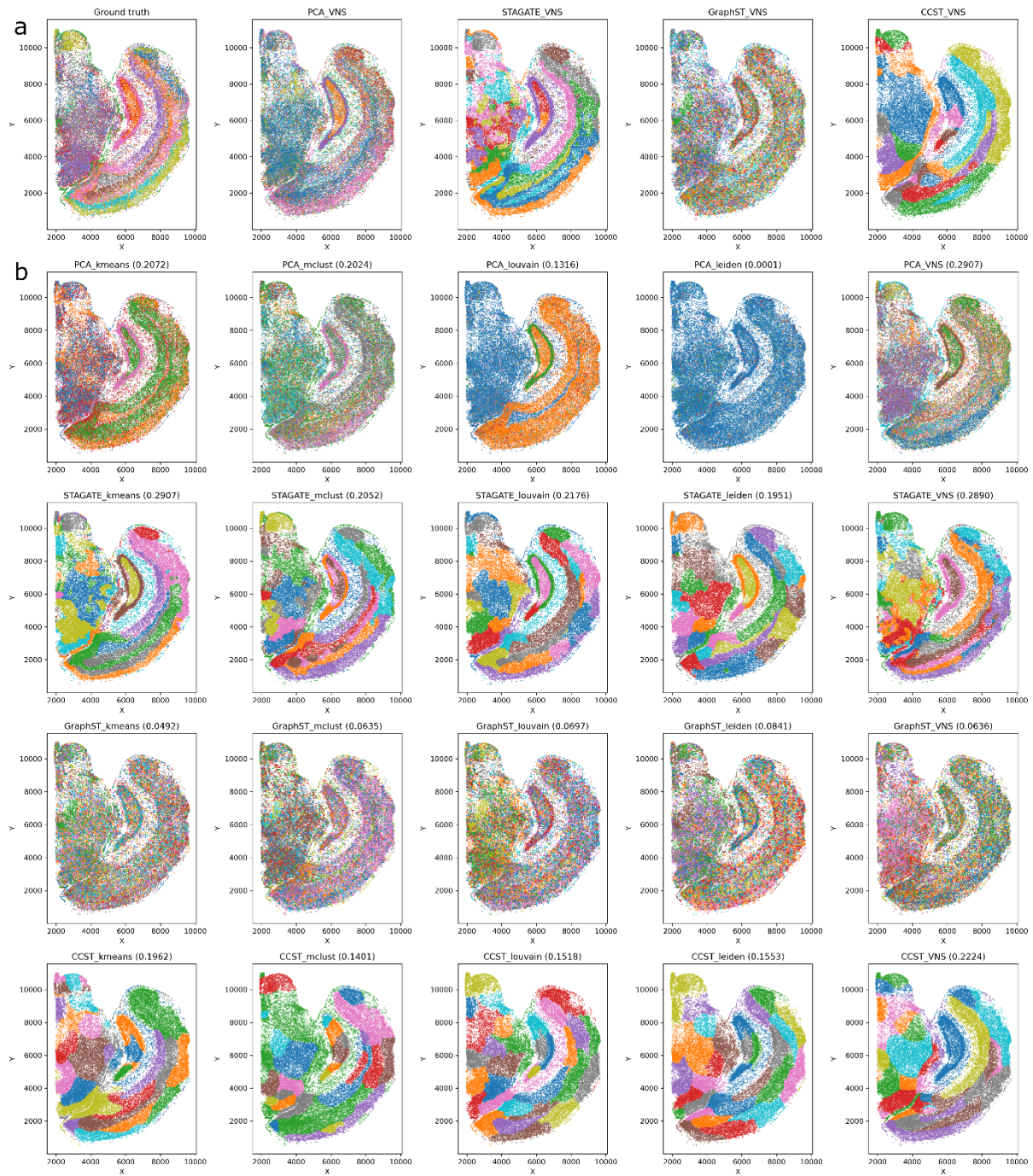
## Acknowledgements

513

**Figure 1.** (a) Results of the VNS clustering on the E2 dataset. The first figure on the left presents the ground truth data. These results were obtained using the VNS method with PCA, STAGATE, GraphST, and CCST embeddings. (b) Clustering results for the E2 dataset. Each row presents the clustering results obtained by $k$-Means, MClust, Louvain, Leiden, and VNS over a certain embedding method. Therefore, the first row presents the results obtained by all clustering methods when using PCA embedding. The next three rows used STAGATE, GraphST, and CCST embeddings.
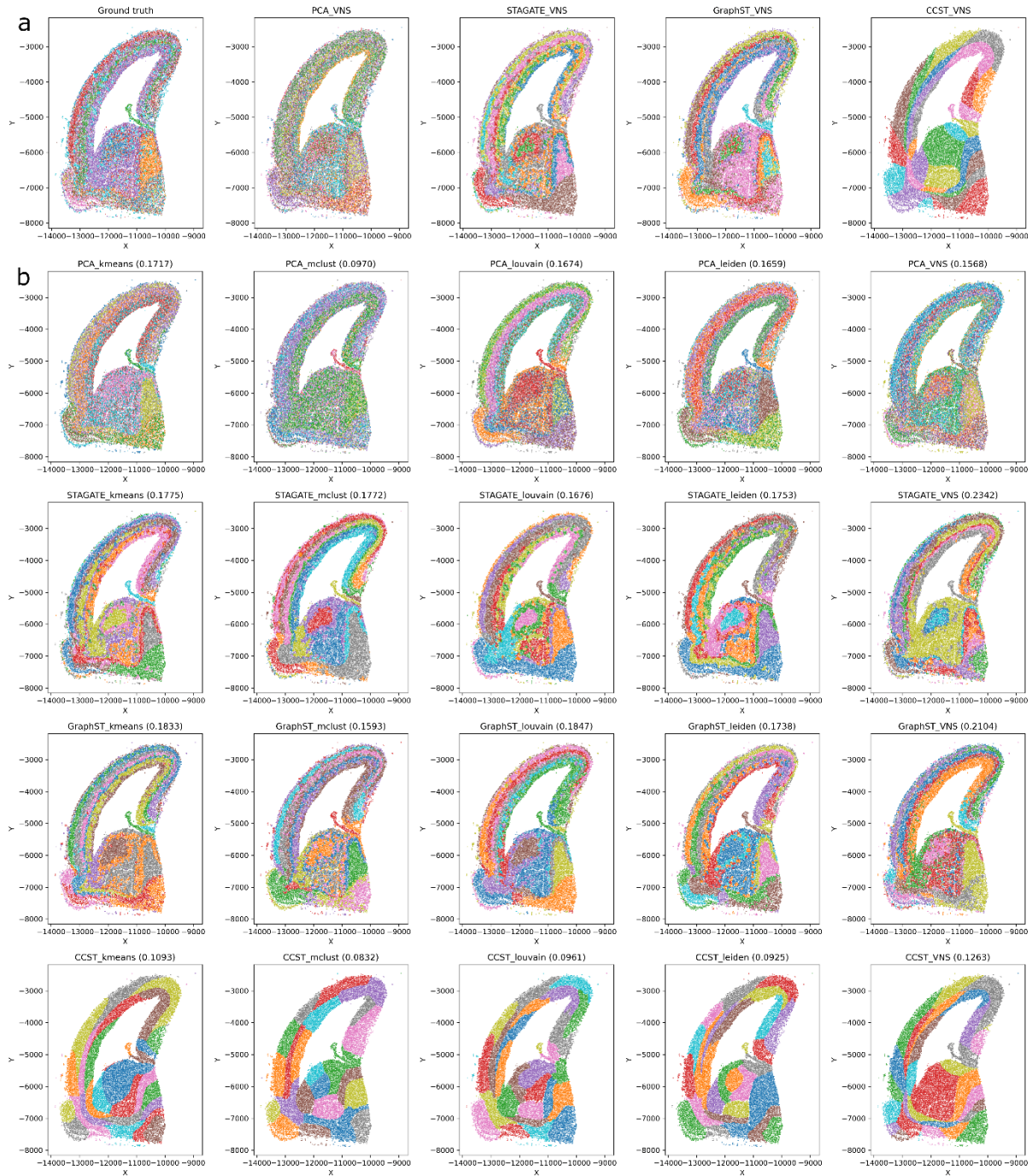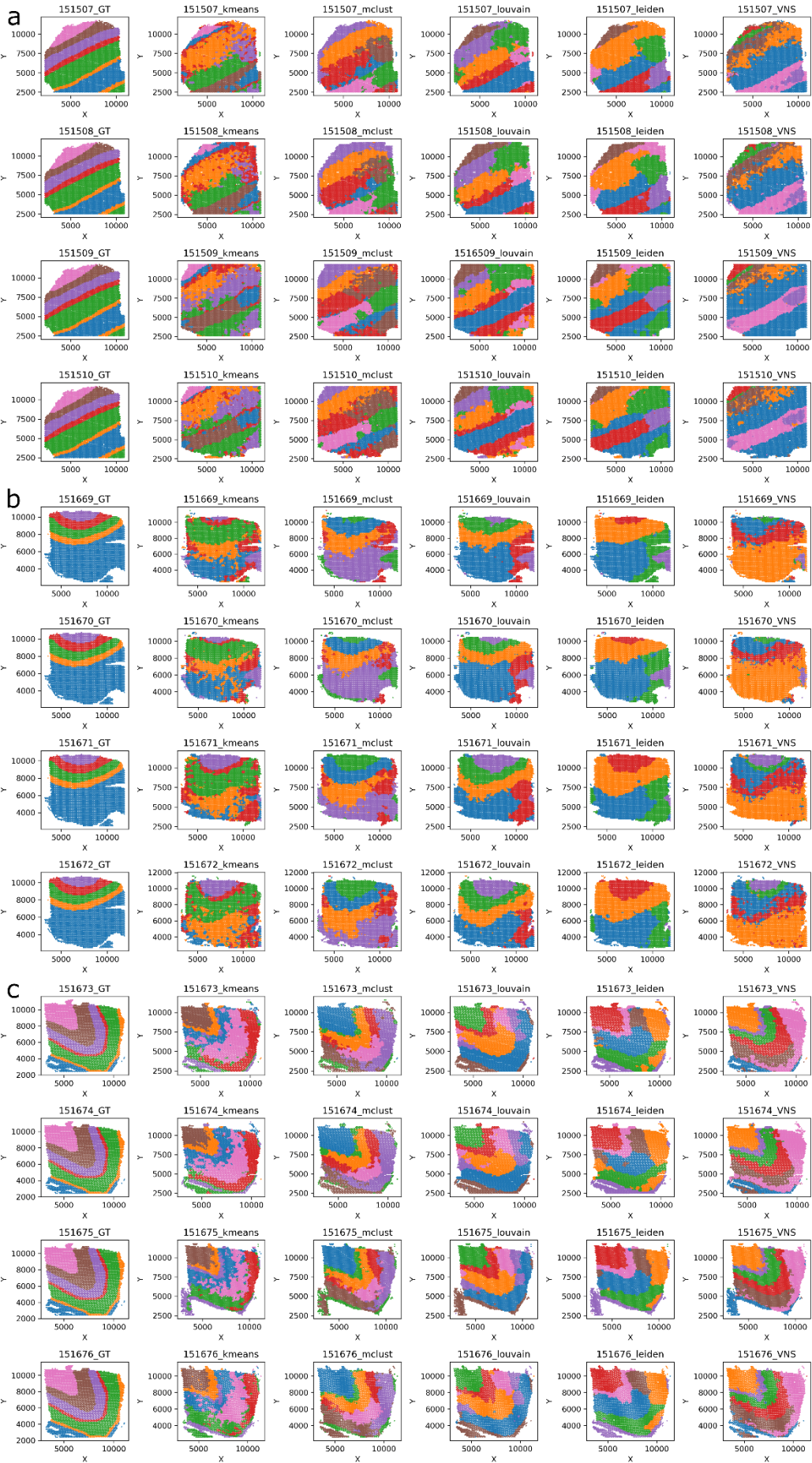
**Figure 2.** (a) Results of the VNS clustering on the Forebrain dataset. The first figure on the left presents the ground truth data. These results were obtained using the VNS method with PCA, STAGATE, GraphST, and CCST embeddings. (b) Clustering results for the Forebrain dataset. Each row presents the clustering results obtained by $k$-Means, MClust, Louvain, Leiden, and VNS, over a certain embedding method. Therefore, the first row presents the results obtained by all clustering methods when using PCA embedding. The next three rows used STAGATE, GraphST, and CCST embeddings.

a

| 151507_GT | 151507_kmeans | 151507_mclust | 151507_louvain | 151507_leiden | 151507_VNS |
| 151508_GT | 151508_kmeans | 151508_mclust | 151508_louvain | 151508_leiden | 151508_VNS |
| 151509_GT | 151509_kmeans | 151509_mclust | 1516509_louvain | 151509_leiden | 151509_VNS |
| 151510_GT | 151510_kmeans | 151510_mclust | 151510_louvain | 151510_leiden | 151510_VNS |

b

| 151669_GT | 151669_kmeans | 151669_mclust | 151669_louvain | 151669_leiden | 151669_VNS |
| 151670_GT | 151670_kmeans | 151670_mclust | 151670_louvain | 151670_leiden | 151670_VNS |
| 151671_GT | 151671_kmeans | 151671_mclust | 151671_louvain | 151671_leiden | 151671_VNS |
| 151672_GT | 151672_kmeans | 151672_mclust | 151672_louvain | 151672_leiden | 151672_VNS |

c

| 151673_GT | 151673_kmeans | 151673_mclust | 151673_louvain | 151673_leiden | 151673_VNS |
| 151674_GT | 151674_kmeans | 151674_mclust | 151674_louvain | 151674_leiden | 151674_VNS |
| 151675_GT | 151675_kmeans | 151675_mclust | 151675_louvain | 151675_leiden | 151675_VNS |
| 151676_GT | 151676_kmeans | 151676_mclust | 151676_louvain | 151676_leiden | 151676_VNS |

531

**Figure 3.** The clustering results on the DLPFC datasets 151507-151510, 151669-151672, and 151673-151676 are presented in panels (a), (b), and (c), respectively. The first column shows the ground truth data, while the subsequent columns display the results obtained using $k$-Means, MClust, Louvain, Leiden, and the VNS method with STAligner embeddings.

# References

1. Giladi, A., Amit, I. Single-cell genomics: a stepping stone for future immunology discoveries. *Cell*, 2018, 172.1: 14-21.
2. Trapnell, C. et al. The dynamics and regulators of cell fate decisions are revealed by pseudo temporal ordering of single cells. *Nature biotechnology*, 2014, 32.4: 381-386.
3. Stuart, T. et al. Comprehensive integration of single-cell data. *Cell*, 2019, 177.7: 1888-1902. e21.
4. Satija, R. et al. Spatial reconstruction of single-cell gene expression data. *Nature biotechnology*, 2015, 33.5: 495-502.
5. Vickovic, S. et al. High-definition spatial transcriptomics for in situ tissue profiling. *Nature methods*, 2019, 16.10: 987-990.
6. Williams C.G., Lee H.J., Asatsuma T., Vento-Tormo R., Haque A. An introduction to spatial transcriptomics for biomedical research. Genome Medicine 2022;14(1):1–18.
7. Hansen, P., Mladenović, N., Brimberg, J., Pérez, J.A.M. (2010). Variable Neighborhood Search. In: Gendreau, M., Potvin, JY. (eds) Handbook of Metaheuristics. International Series in Operations Research & Management Science, vol 146. Springer, Boston, MA. https://doi.org/10.1007/978-1-4419-1665-5_3
8. MacQueen J, et al. Some methods for classification and analysis of multivariate observations. In: Proceedings of the fifth Berkeley symposium on mathematical statistics and probability, vol. 1 Oakland, CA, USA; 1967. p. 281–297.
9. Blondel V.D., Guillaume J.L., Lambiotte R., Lefebvre E. Fast unfolding of communities in large networks. Journal of statistical mechanics: theory and experiment 2008;2008(10): P10008.
10. Traag V.A., Waltman L., Van Eck N.J. From Louvain to Leiden: guaranteeing well-connected communities. Scientific reports 2019;9(1):5233.
11. Fraley C., Raftery A. MCLUST: Software for model-based cluster and discriminant analysis. Department of Statistics, University of Washington: Technical Report 1998; 342:1312.
12. Reynolds D.A., et al. Gaussian mixture models. Encyclopaedia of biometrics 2009; 741(659-663).
13. Moon TK. The expectation-maximization algorithm. IEEE Signal processing magazine 1996;13(6):47–60.
14. Dong K., Zhang S. Deciphering spatial domains from spatially resolved transcriptomics with an adaptive graph attention autoencoder. Nature communications 2022;13(1):1739.
15. Karamizadeh S., Abdullah S.M., Manaf A.A., Zamani M., Hooman A. An overview of principal component analysis. Journal of Signal and Information Processing 2013;4(3B):173.
16. Long Y, Ang KS, Li M, Chong KLK, Sethi R, Zhong C, et al. Spatially informed clustering, integration, and deconvolution of spatial transcriptomics with GraphST. Nature Communications 2023;14(1):1155.
17. Li J, Chen S, Pan X, Yuan Y, Shen HB. Cell clustering for spatial transcriptomics data with graph neural networks. Nature Computational Science 2022;2(6):399–408.
18. Zhou, X., Kangning D., and Shihua Z. "Integrating spatial transcriptomics data across different conditions, technologies and developmental stages." Nature Computational Science (2023): 1-13.

580    19. Davidovic T, Glišovic N, Raškovic M. Bee colony optimization for clustering incomplete data.
581          In: The 7th International Conference on Optimization Problems and Their Applications,
582          OPTA-2018; 2018. https://ceur-ws.org/Vol-2098/paper8.pdf

583    20. Mladenovic N. A variable neighborhood algorithm-a new metaheuristic for combinatorial
584          optimization. In: papers presented at Optimization Days, vol. 12; 1995.

585    21. Mladenovic N, Hansen P, Variable neighbourhood search, computer and operations
586          research. Computers & Operations Research 1997; 24(11) p1097-1100
587          https://doi.org/10.1016/S0305-0548(97)00031-2

588    22. Hansen P, Mladenović N. (1999). An Introduction to Variable Neighborhood Search. In: Voß,
589          S., Martello, S., Osman, I.H., Roucairol, C. (eds) Meta-Heuristics. Springer, Boston, MA.
590          https://doi.org/10.1007/978-1-4615-5775-3_30

591    23. Chen A, Liao S, Cheng M, Ma K, Wu L, Lai Y, et al. Spatiotemporal transcriptomic atlas of
592          mouse organogenesis using DNA nanoball-patterned arrays. Cell 2022;185(10):1777–1792.

593    24. Maynard KR, Collado-Torres L, Weber LM, Uytingco C, Barry BK, Williams SR, et al.
594          Transcriptome-scale spatial gene expression in the human dorsolateral prefrontal cortex.
595          Nature neuroscience 2021;24(3):425–436.

596    25. Shen R, Liu L, Wu Z, Zhang Y et al. (2022). Data from: Application of Spatial-ID to large field
597          mouse brain hemisphere dataset measured by Stereo-seq [Data set]. Zenodo.
598          https://doi.org/10.5281/zenodo.7340795

599    26. STOMICS database MOSTA download https://db.cngb.org/stomics/mosta/download/

600    27. Spatial LIBD GitHub https://github.com/LieberInstitute/spatialLIBD

601    28. Hubert L, Arabie P. Comparing partitions. Journal of Classification 1985; 2: 193–218.

602    29. Djordjevic A; Li J; Fang S; Cao L; Ivanovic M (2024): Supporting data for "A Novel Variable
603          Neighborhood Search Approach for Cell Clustering for Spatial Transcriptomics" GigaScience
604          Database. http://dx.doi.org/10.5524/102498

605