



Open Access This file is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons license, and indicate if changes were made. In the cases where the authors are anonymous, such as is the case for the reports of anonymous peer reviewers, author attribution should be to 'Anonymous Referee' followed by a clear attribution to the source work. The images or other third party material in this file are included in the article's Creative Commons license, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons license and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this license, visit <http://creativecommons.org/licenses/by/4.0/>.

Reviewers' Comments:

Reviewer #1:

Remarks to the Author:

In this manuscript, Long et al. present a deep graph network-based method DeepST for 1) inference of spatial domains, 2) data integration between ST samples or between ST and scRNA-seq data, 3) deconvolution. The method utilized graph neural network together with self-contrastive learning techniques. The authors showed better performance than several benchmarked methods, and demonstrated its abilities on identifying meaningful spatial domains and integrating different samples or data modalities.

Overall, the topic is important and timely. While this method has strong potential in application, the manuscript lacks clarity on the details of the method, and it is unclear what unique advantages the method has compared with other existing methods in terms of its performance. Here are several specific major points in addition to multiple minor point for improving the work. Below are detailed comments:

Major concerns:

1. Lines 123-124, "To our knowledge, few existing methods use self-supervised contrastive learning on spatial transcriptomics". While the work may be new in using self-supervised contrastive learning on spatial transcriptomics, there are some papers that use such kind of methods to address similar questions. For example,
 - 1) Ren, H., Walker, B.L., Cang, Z. et al. Identifying multicellular spatiotemporal organization of cells with SpaceFlow. Nat Commun 13, 4076 (2022). <https://doi.org/10.1038/s41467-022-31739-w>
 - 2) conST: an interpretable multi-modal contrastive learning framework for spatial transcriptomics. Yongshuo Zong, Tingyang Yu, Xuesong Wang, Yixuan Wang, Zhihang Hu, Yu Li bioRxiv 2022.01.14.476408; doi: <https://doi.org/10.1101/2022.01.14.476408>. What is the novelty of the proposed method compared to those two methods?
2. I tried to test and reproduce results shown in this manuscript using the provided links in the manuscript (<https://deepst-tutorials.readthedocs.io/>), and observed the following problem:
 - a. DeepST/utlils.py:121, in refine_label(adata, radius, key)

```
120 for j in range(1, n_neigh+1):  
--> 121 neigh_type.append(old_type[index[j]])
```

IndexError: index 3602 is out of bounds for axis 0 with size 3583

Basically, in the clustering step (In refine_label function), the shape of old_type and index are inconsistent, which might be caused by the inconsistent size between the dimensions from adata and distance matrix. As a result, I couldn't reproduce the results shown in the paper. See the attached `DeepST_test.ipynb` for details.

3. Regarding to the method, especially the contrastive learning component, both the formula and the ideas are very similar to Deep Graph Infomax (DGI) (Veličković et al. 2018). What are the novel elements and major differences between current method and DGI? This needs to be addressed.
4. What are the meaning and motivation of formula (5), line 734? It's important to show the performance difference with and without adding this term by experiment, because DGI only contains (4) instead of (5). Does this term actually improve the performance?
5. The reason of using reconstructed expression data to cluster instead of using the latent embedding need to be justified. Moreover, why choosing mclust over graph-based methods such as Leiden, Louvain? It's important to justify such choices in terms of data analysis and results.
6. In line 684, do the authors augment data through creating corrupted graph by randomly adding or dropping edges? What is the effect of such procedure on the overall performance of the method.
7. Regarding the data integration performance shown in Figure 4, why did the authors not compare many other methods designed for nonspatial scRNA-seq data, such as scVI (Lopez et al. 2018) and Harmony (Korsunsky et al. 2019), because those classical methods have been well demonstrated for good performance for single-cell data.
8. It's important to show the STAGATE results that similar to Fig 4E to better demonstrate the data

integration performance.

Minor points

9. To better support the manual annotation result in Fig 6A, the spatial expression distribution of several marker genes for each panels in Fig 6A need to be added.
10. In line 836, the author mentioned the first loss term indicates contrastive loss, why is there only one instead of two terms? What is the meaning of the first term?
11. In line 712, which norm is used? L1 or L2 or others?
12. In Fig 3C, the titles of panels Mesenchyme and Dermomyotome seem misplaced.
13. All color bars need to be explained for their meanings.
14. Many typos and grammar errors in the manuscript, e.g., in line 28, "has" should be "have"; in line 59, "K-means" should be "k-means"
15. Lines 124-126, "Using self-supervised contrastive learning improves performance in learning relevant latent features and has the additional benefit of removing batch effects". This sentence occurs without any supporting evidence. It needs to be fixed.
16. In lines 158-162, the authors introduced "self-reconstruction loss" and "contrastive loss" and their effects. It's important to show what the two losses are in the context of biology.
17. It's unclear how the neighbor graph is constructed. In the caption of Fig. 1 (lines 1034-1035), the authors wrote "...neighbor graph constructed using spot coordinates (x,y) of that fall within a distance threshold". However, in the method section in lines 665-675, the authors wrote "Finally, we select the top k-nearest spots as its neighbors". It's unclear whether the authors used a distance threshold or a threshold for k.
18. Regarding the method (lines 655-659), the descriptions seem to be for the spatial transcriptomics data. However, this is not clear from the description, as two kinds of datasets (spatial transcriptomics data and scRNA-seq data) are mentioned in this paper.
19. In lines 684-688, "...while keeping the original graph structure unchanged": was the corrupted neighbor graph G' the same as the original G ?
20. In lines 708-709, " W_d and b_d represent the trainable weight matrix and bias vector, respectively, which are shared by all nodes in the graph". Please justify why W_d and b_d need to be shared by all nodes in the graph. Besides, is this the same case for W_e and b_e ?

Reference

- Korsunsky, Ilya, Nghia Millard, Jean Fan, Kamil Slowikowski, Fan Zhang, Kevin Wei, Yuriy Baglaenko, Michael Brenner, Po-Ru Loh, and Soumya Raychaudhuri. 2019. "Fast, Sensitive and Accurate Integration of Single-Cell Data with Harmony." *Nature Methods* 16 (12): 1289–96.
- Lopez, Romain, Jeffrey Regier, Michael B. Cole, Michael I. Jordan, and Nir Yosef. 2018. "Deep Generative Modeling for Single-Cell Transcriptomics." *Nature Methods* 15 (12): 1053–58.
- Veličković, Petar, William Fedus, William L. Hamilton, Pietro Liò, Yoshua Bengio, and R. Devon Hjelm. 2018. "Deep Graph Infomax." *ArXiv [Stat.ML]*. arXiv. <http://arxiv.org/abs/1809.10341>.

Reviewer #2:

Remarks to the Author:

In the manuscript "DeepST: A novel graph" by Long and coworkers the authors develop a new method for better describing spatial transcriptomics data and being able to integrate multiple studies. The method is based on graph neural networks and contrastive learning, which makes it possible to combine scRNA-seq of better resolution and spatial transcriptomics. The method is logically sound and makes a lot of sense. Moreover, authors show that it empirically identifies more relevant clusters and allows data integration for higher power. Although, I am not an expert in spatial transcriptomics these problems seem of great importance and authors spend good effort to show that it works a planned.

Having said that my expertise is in neural networks and translational bioinformatics I believe that the paper would be a good contribution to the spatial transcriptomics field. From my side, I have no concerns of the paper and like to see it published.

1 **Response to comments for paper NCOMMS-22-35863B “Spatially informed clustering,**
2 **integration, and deconvolution of spatial transcriptomics with GraphST”**

3
4 **Response to Reviewer #1**

5
6 **Summary:** *In this manuscript, Long et al. present a deep graph network-based method GraphST for 1)*
7 *inference of spatial domains, 2) data integration between ST samples or between ST and scRNA-seq*
8 *data, 3) deconvolution. The method utilized graph neural network together with self-contrastive learning*
9 *techniques. The authors showed better performance than several benchmarked methods, and*
10 *demonstrated its abilities on identifying meaningful spatial domains and integrating different samples or*
11 *data modalities.*

12
13 *Overall, the topic is important and timely. While this method has strong potential in application, the*
14 *manuscript lacks clarity on the details of the method, and it is unclear what unique advantages the*
15 *method has compared with other existing methods in terms of its performance. Here are several specific*
16 *major points in addition to multiple minor point for improving the work. Below are detailed comments:*

17 **Response:** We thank the reviewer for the positive comments. We summarize Graph ST’s technological
18 novelty and advantages over existing methods in the following paragraphs. We have carefully
19 addressed all the comments and suggestions when preparing this revision of the manuscript. Please
20 let us know if you have additional comments.

21 With rapid technological advances in spatial transcriptomics, it is now widely applied towards studying
22 tissue complexity and cell-cell communications. However, the current bottleneck still lies in data
23 analysis. Although multiple methods have been developed for spatial transcriptomics, there is still a
24 great need for developing novel tools that offer greater accuracy, robustness, and generalizability
25 towards a wide range of application on different tissue types and technology platforms. Furthermore,
26 the analysis pipeline for spatial transcriptomic data comprises three key tasks, namely spatial clustering,
27 multi-sample integration, and cell type deconvolution. However, there is no comprehensive tool that can
28 perform all these three tasks. To overcome this limitation, we developed GraphST, the first of its kind
29 that integrates these tasks into a streamlined process. Most importantly, GraphST outperforms existing
30 methods in each task. We achieved this by adopting and tailoring graph self-supervised contrastive
31 learning for spatial transcriptomics analysis.

32 In the spatial clustering task, we achieved higher accuracy and robustness with an average of 10%
33 improvement over the best of existing methods in a variety of datasets. Our GraphST clusters revealed
34 finer tissue structures and niches in complex tissues such as the brain, olfactory bulb, and embryo.
35 Although existing methods conST and SpaceFlow also adopted graph contrastive learning for spatial
36 clustering, there are major technical differences and performance advantages when comparing
37 GraphST to conST and SpaceFlow. Briefly, GraphST is different from DGI, conST and SpaceFlow in
38 three aspects: A) definition of positive/negative pairs, B) objective function and contrastive loss, and C)
39 training procedure. These differences enable GraphST to outperform the other methods in the spatial
40 clustering task. Furthermore, we have conducted several ablation studies to confirm that each of these
41 differences improves the effective integration of gene expression and spatial context to obtain
42 informative and discriminative latent representations. Please kindly refer to Response 1.1 for details of
43 comparison between GraphST and conST, SpaceFlow.

44 In the multi-sample integration task, GraphST can better correct batch effects when integrating serial
45 tissue slices than existing methods that have been developed for spatial (e.g., STAGATE) or non-spatial
46 batch integration (e.g., scVI and Harmony). Moreover, for the horizontal integration of mouse anterior
47 and posterior brain slices, GraphST outperformed SpaGCN and STAGATE in that GraphST could
48 assign the common cortical layers that aligned well across the shared edge and also reveal the dorsal
49 and ventral horns of the hippocampus regions.

50 In the final task, GraphST produced more accurate cell type deconvolution with simulation data than
51 existing methods, including cell2location that was recognized as the top performing method in a recent
52 benchmark. Moreover, when applied to 10x Visium acquired human lymph node data, GraphST’s

53 deconvolution was able to better capture the germinal centers and mapped the B cell subpopulations
54 with higher spatial coherence. Lastly, application on human breast cancer 10x Visium data revealed
55 immune cell distributions across healthy, tumor edge, invasive ductal carcinoma (IDC), and ductal
56 carcinoma in situ (DCIS) regions. In particular, the T cells enriched in the IDC regions showed
57 upregulation of known exhaustion markers including *LAG3*, *TIGIT*, *PD1*, *TIM3*, and *CTLA4*, suggesting
58 a tumor induced immune suppressive environment.

59

60 **Major concerns:**

61 **Comment 1.1.** Lines 123-124, “To our knowledge, few existing methods use self-supervised contrastive
62 learning on spatial transcriptomics”. While the work may be new in using self-supervised contrastive
63 learning on spatial transcriptomics, there are some papers that use such kind of methods to address
64 similar questions. For example,

65 1) Ren, H., Walker, B.L., Cang, Z. et al. Identifying multicellular spatiotemporal organization of cells with
66 SpaceFlow. *Nat Commun* 13, 4076 (2022). <https://doi.org/10.1038/s41467-022-31739-w>

67 2) conST: an interpretable multi-modal contrastive learning framework for spatial
68 transcriptomics. Yongshuo Zong, Tingyang Yu, Xuesong Wang, Yixuan Wang, Zhihang Hu, Yu Li
69 *bioRxiv* 2022.01.14.476408; doi: <https://doi.org/10.1101/2022.01.14.476408>.

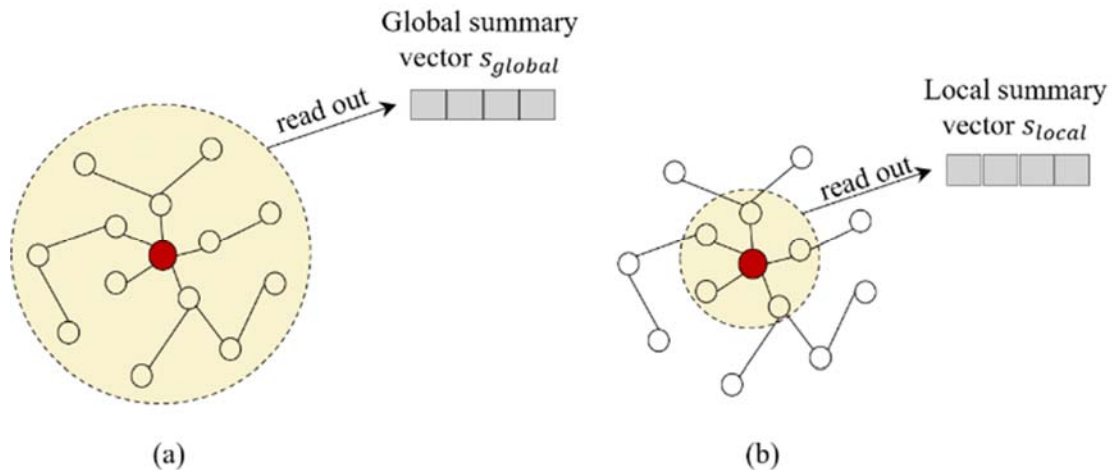
70 *What is the novelty of the proposed method compared to those two methods?*

71 **Response 1.1:** Thank you very much for raising this critical issue. Indeed, GraphST bears much
72 similarity to conST and SpaceFlow with the use of graph contrastive learning for spatial clustering.
73 However, there are several major differences between GraphST and the other two methods.

74 First, and most importantly, both conST and SpaceFlow were mainly developed for spatial clustering
75 only. In addition to spatial clustering, GraphST can be also applied to two other important ST data
76 analysis tasks, multi-sample integration and cell type deconvolution of ST. GraphST comprises three
77 modules with different network architectures tailored for each of the three tasks respectively.

78 Secondly, even for the spatial clustering task, there are also major differences when comparing
79 GraphST to conST and SpaceFlow, despite all three methods adopting graph contrastive learning
80 similar to DGI. Here we elaborate on their differences in three aspects: A) definition of positive/negative
81 pairs, B) objective function and contrastive loss, and C) training procedure.

82 A) GraphST’s contrastive learning is different from DGI, conST, and SpaceFlow in terms of their
83 definition of positive/negative pairs. DGI, conST, and SpaceFlow construct positive/negative pairs
84 by pairing spot embedding h_i/h'_i from the original/corrupted graph with a global summary vector
85 s_{global} (as shown in Figure R1 (a)). Therefore, the spot embedding learned by DGI, conST, and
86 SpaceFlow captures more of the global structure information but less spot-specific local
87 neighbourhood information. Such contrastive learning may result in feature overfitting and reduced
88 spot-to-spot variability. To deal with this issue, GraphST improves over DGI’s contrastive learning
89 by re-defining the positive/negative pairs. Specifically, motivated by the assumption that different
90 spots in a tissue sample have different local spatial contexts, we define positive/negative pairs by
91 pairing spot embedding h_i/h'_i with its local summary vector s_{local} (as shown in Figure R1(b)) instead
92 of the global summary vector. With the local summary vector, the model can better preserve local
93 context information and spot-to-spot variability. We demonstrate the effectiveness of local context
94 with an ablation study describe in Figure R4.



95
96
97

Figure R1. Illustrations of local and global summary vectors.

98 B) GraphST is also different from SpaceFlow, conST, and DGI in terms of the objective function and
99 contrastive loss formulations. The objective function of GraphST includes contrastive loss and
100 reconstruction loss, while DGI's objective function includes only contrastive loss. The objective
101 function of SpaceFlow includes both contrastive loss and a spatial consistency penalty term.
102 Addition of the penalty term helps SpaceFlow bring spatially adjacent spots closer in the latent
103 embedding. However, the lack of reconstruction loss in DGI and SpaceFlow may lead to insufficient
104 preservation of the original gene expression information. In contrast, GraphST adds reconstruction
105 loss to its objection function to ensure that the latent embedding preserves the original gene
106 expression information effectively. Furthermore, the contrastive loss functions are also different
107 between GraphST, DGI, and SpaceFlow. GraphST adopts symmetric contrastive loss (formulas (4)
108 and (5)) for model training while conST and SpaceFlow use single contrastive loss like DGI.
109 Symmetric contrastive loss can help stabilize the model and learn a better representation as
110 illustrated in Figure R5.

111
112
113
114
115
116
117
118
119
120

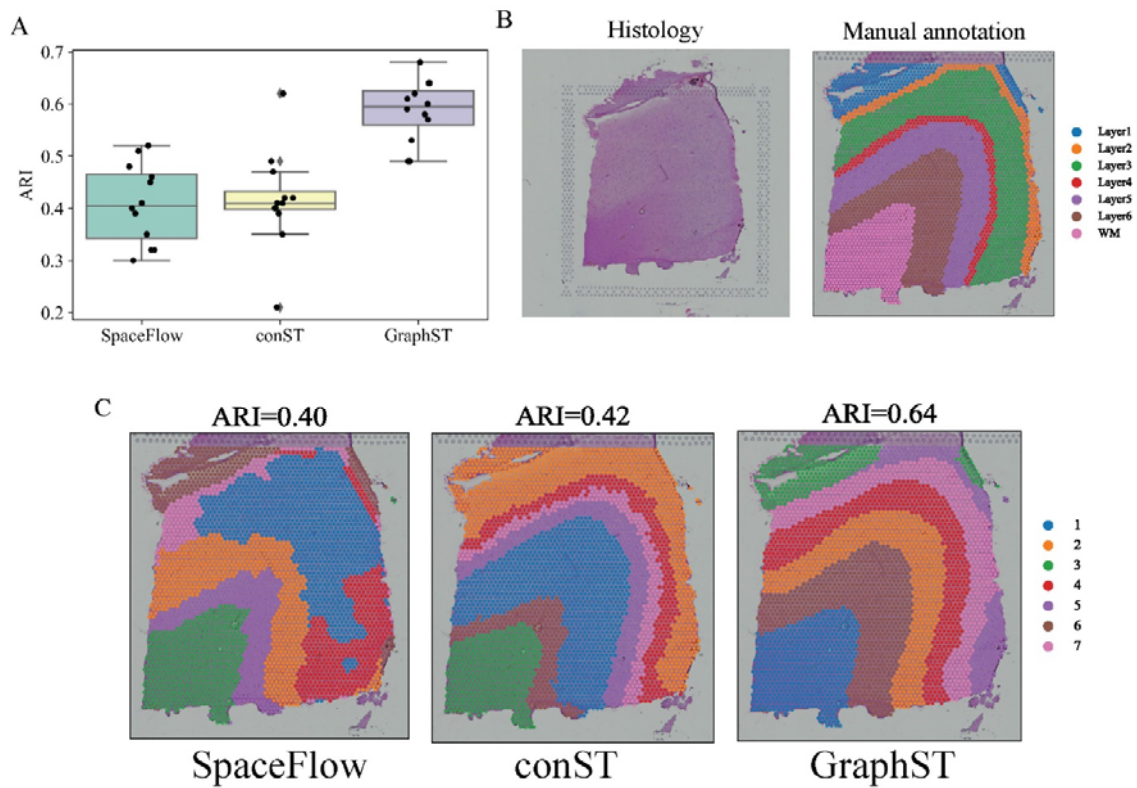
C) Lastly, although conST's objective function also contains contrastive loss and reconstruction loss
like GraphST, GraphST's training procedure is different from that of conST. conST splits the training
into pre-training and major training stages, where the model is trained with reconstruction loss in
the pre-training stage and contrastive loss in the major training stage. This two-stage training
procedure lacks mutual constraints on the contrastive and reconstructive loss, thus may fail to
identify the optimal combination of the two loss functions. In contrast, GraphST trains the model in
a single step by jointly optimizing for the reconstruction and contrastive losses. During this training,
GraphST can adaptively adjust the contributions of the different loss functions to achieve better
representation learning.

121
122
123
124
125
126
127
128
129
130
131
132
133
134

Based your comments, we compared the performance of conST, SpaceFlow, and GraphST in the
spatial clustering task with the DLPFC dataset. Figure R2 (A) shows the median ARI scores of the
different methods. We can see that GraphST achieves a much higher median ARI score of 0.60 over
conST (0.41) and SpaceFlow (0.41). Figure R2 (C) shows the results of SpaceFlow, conST, and
GraphST on slice #151673. Visually, SpaceFlow has the poorest performance among the three
methods. The domains identified by SpaceFlow are irregular though it can accurately recover the WM
domain and layer 1. conST performs slightly better than SpaceFlow with each identified domain being
continuous. However, most of the domains do not match the manual annotation well. In contrast,
GraphST's clusters are more continuous than conST and SpaceFlow, and are more consistent with the
manual annotation. Quantitatively, GraphST achieves the highest ARI score of 0.64 among the three
methods. Overall, GraphST outperforms conST and SpaceFlow in spatial clustering. The results of all
12 DPLFC slices are shown in Figure R3, which again illustrates GraphST's advantages over conST
and SpaceFlow. We have added these results to Figure 2 in the revised manuscript and Supplementary
Figure S1.

135 To evaluate the effectiveness of local context over global context, we conducted an ablation study by
 136 comparing GraphST with a variant that uses a global summary vector instead of local summary vectors.
 137 We ran GraphST and the variant on the 12 DLPFC slices and evaluated their performance using their
 138 median ARI scores. Figure R4 shows that GraphST outperforms the variant (median ARI score of 0.51)
 139 with a significantly higher median ARI score of 0.60. This ablation study demonstrated that local context
 140 does help GraphST perform better than with the global context. We have added these results to
 141 Supplementary Figure S14A in the revised manuscript.

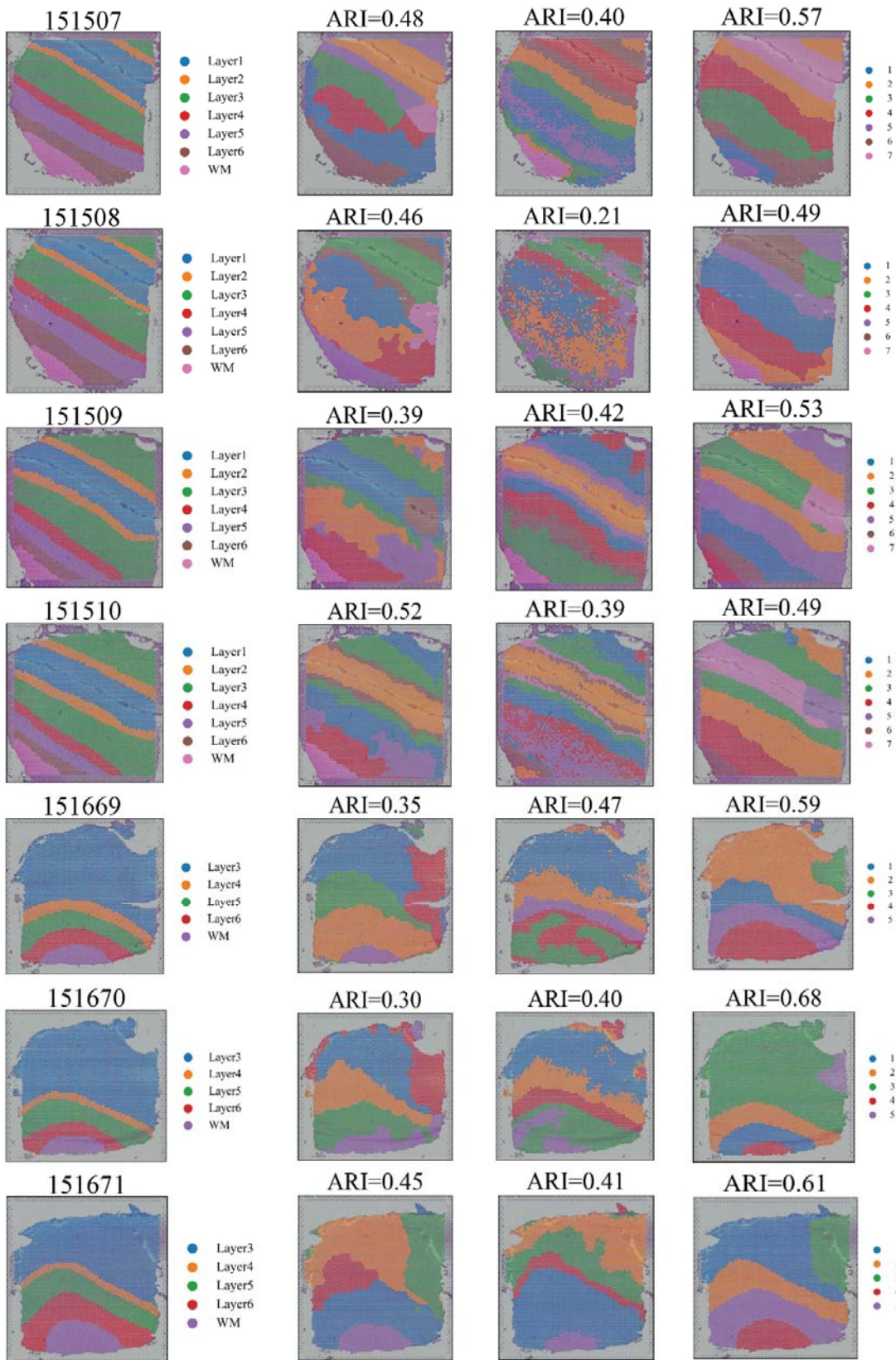
142 To demonstrate the advantage of symmetric contrastive loss over single contrastive loss, we conducted
 143 another ablation study to compare GraphST with a variant that does not use the contrastive corrupted
 144 loss (formula (5) in the manuscript). We tested GraphST and this variant on the 12 DLPFC samples
 145 and evaluated their performance with the ARI metric. Figure R5 shows that GraphST achieved much
 146 better performance than the variant, showing that the contrastive corrupted loss contributes to better
 147 embedding learning. We have added these results to Supplementary Figure S14B of the revised
 148 manuscript.

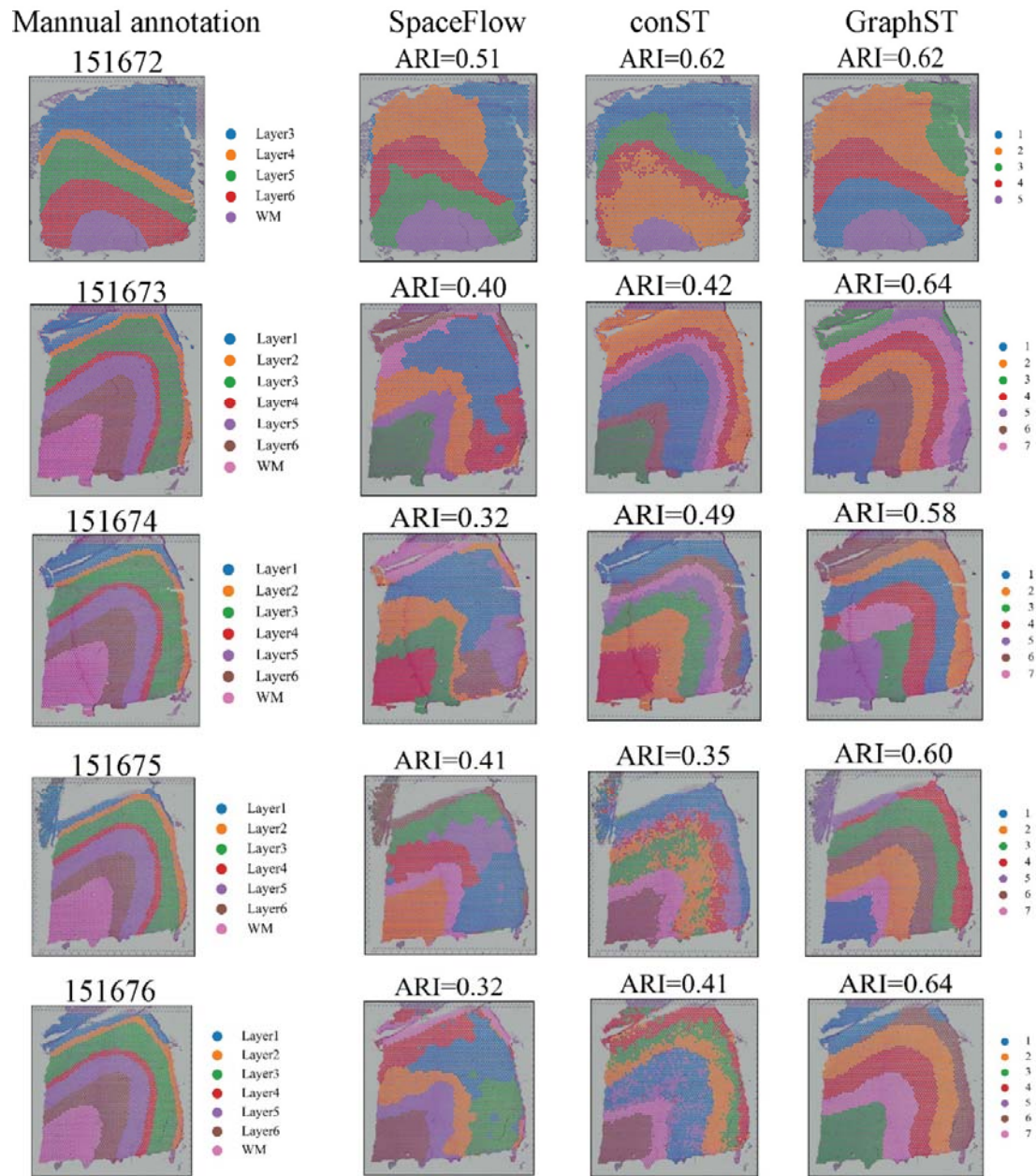


149

150 Figure R2. Comparison between SpaceFlow, conST and GraphST on DLPFC.

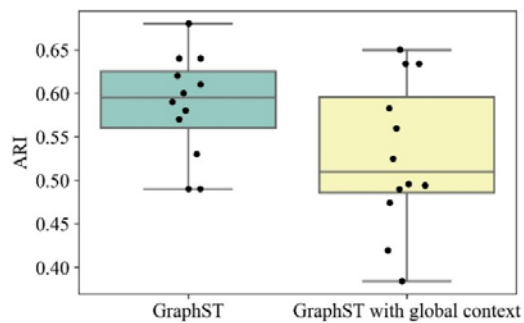
Manual annotation





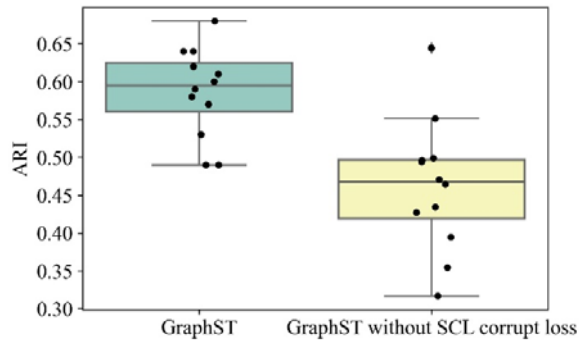
152

153 Figure R3. Comparison between SpaceFlow, conST, and GraphST on DLPFC.



154

155 Figure R4. Comparison analysis between GraphST and the variant that uses a global summary vector
 156 on DLPFC.



157

158 Figure R5. Comparison analysis between GraphST and its variant GraphST without contrastive
 159 corrupted loss on DLPFC.

160

161 **Comment 1.2.** I tried to test and reproduce results shown in this manuscript using the provided links in
 162 the manuscript (<https://GraphST-tutorials.readthedocs.io/>), and observed the following problem:
 163 a. `GraphST/utills.py:121, in refine_label(adata, radius, key)`

164 `120 for j in range(1, n_neigh+1):`

165 `--> 121 neigh_type.append(old_type[index[j]])`

166 `IndexError: index 3602 is out of bounds for axis 0 with size 3583`

167 Basically, in the clustering step (In `refine_label` function), the shape of `old_type` and `index` are
 168 inconsistent, which might be caused by the inconsistent size between the dimensions from `adata` and
 169 distance matrix. As a result, I couldn't reproduce the results shown in the paper. See the attached
 170 `GraphST_test.ipynb` for details.`

171 **Response 1.2:** We are sorry for this error. We have revised and updated our codes that are openly
 172 accessible (<https://github.com/JinmiaoChenLab/GraphST>). Furthermore, we provide a detailed tutorial
 173 to guide users on using our tool. The tutorial is available at [https://deepst-
 174 tutorials.readthedocs.io/en/latest/](https://deepst-tutorials.readthedocs.io/en/latest/). We welcome you to test our codes again according to the tutorial.

175

176 **Comment 1.3.** Regarding to the method, especially the contrastive learning component, both the
 177 formula and the ideas are very similar to Deep Graph Infomax (DGI) (Veličković et al. 2018). What are
 178 the novel elements and major differences between current method and DGI? This needs to be
 179 addressed.

180 **Response 1.3:** Thank you for the comments. As discussed above, the main differences between
 181 GraphST and DGI include the definition of positive/negative pairs and the contrastive loss functions,
 182 which we elaborate in the following paragraphs. In addition to the contrastive loss, GraphST employs
 183 reconstruction loss to preserve the original gene expressions in the latent embedding. DGI in contrast
 184 only uses contrastive loss.

185 DGI constructs positive/negative pairs by pairing each spot embedding h_i/h'_i from the original/corrupted
 186 graph with a global summary vector S_{global} (as shown in Figure R1 (a)). Therefore, the spot embedding
 187 learned by DGI captures more of the global structure information but less spot-specific local
 188 neighbourhood information. Such contrastive learning may result in feature overfitting and reduced
 189 spot-to-spot variability. To deal with this issue, GraphST improves over DGI's contrastive learning by
 190 re-defining the positive/negative pairs. Specifically, motivated by the assumption that different spots in
 191 a tissue sample have different local spatial contexts, we define positive/negative pairs by pairing each
 192 spot embedding h_i/h'_i with its local summary vector S_{local} (as shown in Figure R1(b)) instead of the
 193 global summary vector. With local summary vectors, the model can better preserve local context
 194 information and spot-to-spot variability. To evaluate the effectiveness of using local context over the
 195 global context, we conducted an ablation study by comparing GraphST with a variant that uses the

196 global summary vector instead of local summary vectors. We ran GraphST and the variant on the 12
 197 DLPFC slices and evaluated their performance with their median ARI scores. Figure R4 shows that
 198 GraphST outperformed the variant (median ARI score of 0.51) with a significantly higher median ARI
 199 score of 0.60. This demonstrated that using local context does help GraphST perform better than with
 200 the global context. We have added these results to Supplementary Figure S14A in the revised
 201 manuscript.

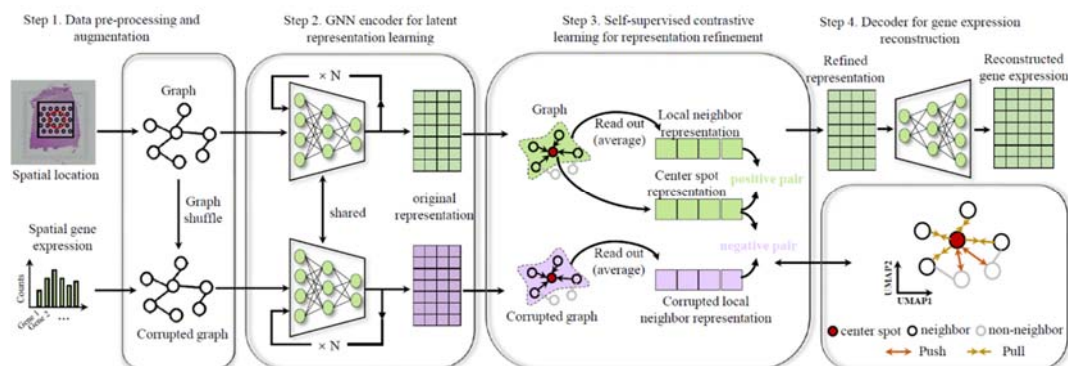
202 Moreover, GraphST is also different from DGI in their contrastive loss functions. DGI uses single
 203 contrastive loss, while GraphST employs symmetric contrastive loss by adding a contrastive corrupted
 204 loss term (formula (5)). Symmetric contrastive loss can help stabilize the model and learn a better
 205 representation. To demonstrate the advantage of symmetric contrastive loss over single contrastive
 206 loss, we conducted an ablation study to compare GraphST with a variant that does not include
 207 contrastive corrupted loss. We tested GraphST and the variant on the 12 DLPFC samples and
 208 evaluated their performance with the ARI metric. Figure R5 shows that GraphST achieved a median
 209 ARI score of 0.60, an improvement of 28% compared to the variant (median ARI score of 0.47). Thus,
 210 we conclude that symmetric contrastive loss does improve the model's performance. We have added
 211 the results in Supplementary Figure S14B of the revised manuscript.

212 **Comment 1.4.** What are the meaning and motivation of formula (5), line 734? It's important to show the
 213 performance difference with and without adding this term by experiment, because DGI only contains (4)
 214 instead of (5). Does this term actually improve the performance?

215 **Response 1.4:** Thank you very much for your question and suggestion. Formula (5) is a contrastive
 216 corrupted loss function that is symmetric to formula (4). The combination of loss functions (4) and (5)
 217 forms a symmetric contrastive loss that can make GraphST's model training more stable and robust,
 218 thus improving the spatial clustering performance.

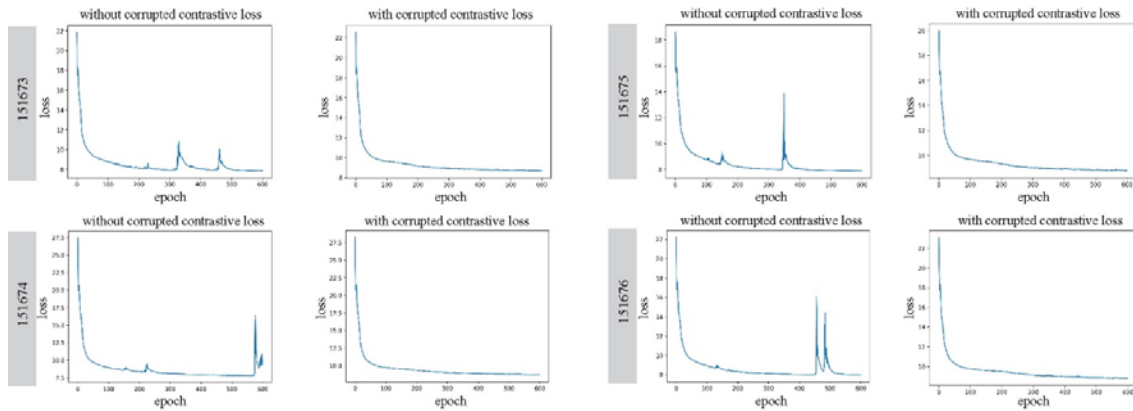
219 As shown in our GraphST workflow (Figure R6), we use the original graph as input to create a corrupted
 220 graph by randomly shuffling features across spots while keeping the adjacency matrix of the graph
 221 unchanged. The original and corrupted graphs are thus structurally identical. At first, we followed DGI
 222 in constructing the contrastive learning with only a single contrastive loss (i.e., formula (4) in the revised
 223 manuscript). However, during model training, we found that the loss curve was unstable, as shown in
 224 Figure R7. Motivated by that and the fact that the original and corrupted graphs are structurally
 225 symmetric, we added a symmetric (corrupted) contrastive loss function to make the model training more
 226 stable and robust.

227 Based on your comments, we conducted an ablation study to validate the effectiveness of symmetric
 228 contrastive loss. We compared GraphST with a variant without contrastive corrupted loss on the 12
 229 DLPFC samples and used the ARI metric for evaluation. The results in Figure R5 show that GraphST
 230 without contrastive corrupted loss achieved a lower median ARI score (0.47) than the original GraphST
 231 (0.60), supporting the idea that symmetric contrastive loss helps our model achieve better performance.
 232 Furthermore, Figure R7 shows that the model training curve is stabilized with symmetric contrastive
 233 loss. These results have been included in Figure S14B in the Supplementary.



234

235 Figure R6. Workflow of GraphST for spatial clustering.



236

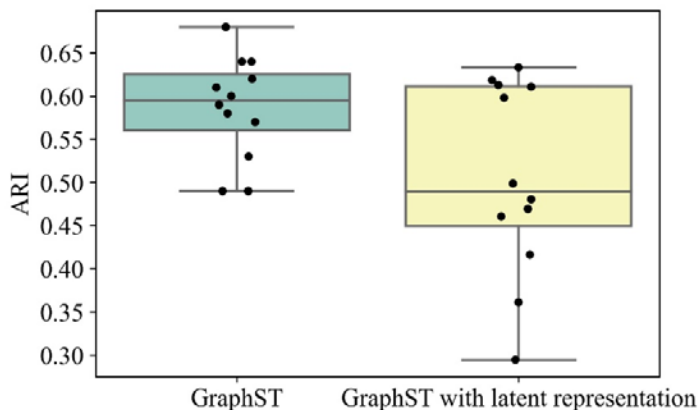
237 Figure R7. Training loss curves with and without contrastive corrupted loss on four DLPFC samples.

238

239 **Comment 1.5.** The reason of using reconstructed expression data to cluster instead of using the latent
 240 embedding need to be justified. Moreover, why choosing mclust over graph-based methods such as
 241 Leiden, Louvain? It's important to justify such choices in terms of data analysis and results.

242 **Response 1.5:** Thank you very much for your very constructive comments. In our framework, the
 243 reconstructed expression is more informative than the latent representation for two reasons. Firstly, as
 244 shown in Figure R6, our GraphST framework consists of a GCN (graph convolutional network)-based
 245 encoder and a GCN-based decoder. The encoder and decoder have symmetrical structures and equal
 246 numbers of GCN layers. In our model, the number of layers for encoder and decoder are set to 1. The
 247 basic principle of GCN is to update the node representation by iteratively aggregating information from
 248 the neighbours. Therefore, as the output of the decoder, the reconstructed expression contains more
 249 information about the local context than the latent representation, as the reconstructed expression
 250 aggregates feature information of two-hop neighbours while the latent representation only aggregates
 251 one-hop neighbours' feature information. Secondly, compared to the latent representation, the
 252 reconstructed expression captures more topological structure and semantic information. This is
 253 because the reconstructed expression is obtained through two GCN layers, meaning that the adjacency
 254 matrix is used twice.

255 In response to your comments, we compared the clustering performance of using the latent
 256 representation and the reconstructed expression on the 12 DLPFC samples. The results in Figure R8
 257 show that GraphST achieved much a higher median ARI score when using the reconstructed
 258 expression for clustering than the latent representation, suggesting that the former contains more useful
 259 information than the latter. These results have been included in Supplementary Figure S14D.



260

261 Figure R8. Comparison analysis between latent representation and reconstructed expression on
 262 DLPFC.

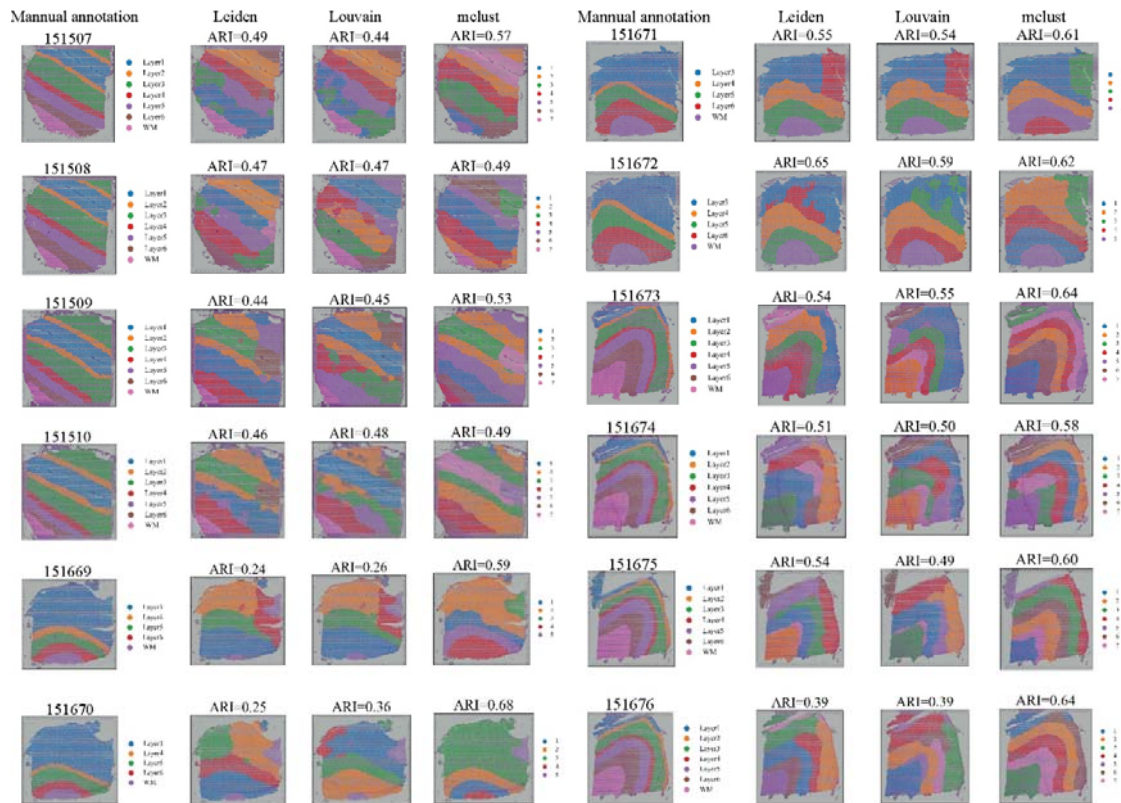
263 We chose mclust as the default clustering method because our assessment showed that mclust
 264 performs better than Leiden and Louvain in most cases. Figure R9 shows the clustering results on the
 265 12 DLPFC samples using Leiden, Louvain, and mclust. mclust consistently outperformed Leiden and
 266 Louvain on all 12 samples in terms of the ARI metric, with a much higher median ARI score (Figure
 267 R10). Visually, the clusters identified by mclust are more continuous. These results have been included
 268 in Supplementary Figure S15. Here, we would like to mention that several previously published spatial
 269 clustering methods such as BayesSpace (Zhao et al., 2021) and STAGATE (Dong and Zhang, 2022)
 270 also use mclust. Nevertheless, we have now added Leiden and Louvain to GraphST as alternative
 271 options.

272 **Reference**

273 Zhao et al. Spatial transcriptomics at subspot resolution with BayesSpace. Nature Biotechnology,
 274 39(11), 1375-1384.

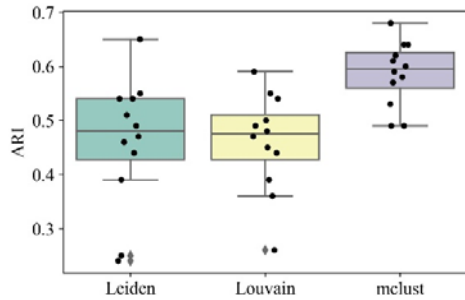
275 Kangning Dong and Shihua Zhang. Deciphering spatial domains from spatially resolved transcriptomics
 276 with an adaptive graph attention auto-encoder. Nature Communications. 2022.

277



278

279 Figure R9. Comparison analysis between Leiden, Louvain, and mclust with the output of GraphST as
 280 input on DLPFC.



281

282 Figure R10. Boxplots of clustering accuracy of Leiden, Louvain, and mclust with the output of GraphST
 283 as input on DLPFC in terms of ARI.

284 **Comment 1.6.** *In line 684, do the authors augment data through creating corrupted graph by randomly*
 285 *adding or dropping edges? What is the effect of such procedure on the overall performance of the*
 286 *method.*

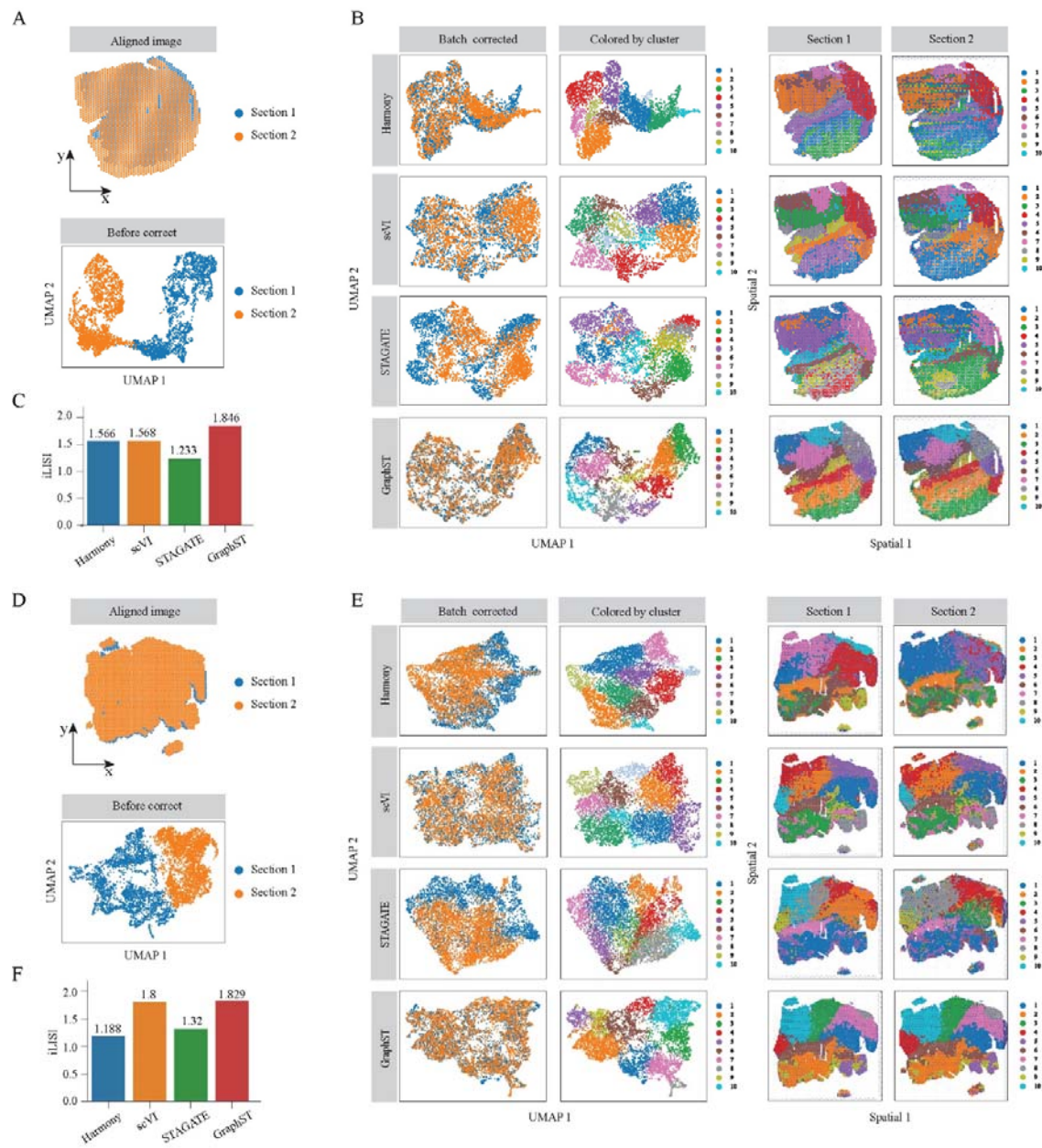
287 **Response 1.6:** We are sorry for the confusion. In our framework, instead of randomly adding or
 288 dropping edges, we create a corrupted graph by randomly shuffling gene expression vectors between
 289 spots while leaving the adjacency matrix of the graph unchanged. In contrastive learning, data
 290 augmentation aims to increase the diversity of training data and thus enhance the model's learning
 291 capability. Therefore, the data augmentation procedure plays an important role in contrastive learning.
 292 During data augmentation, the distribution of the augmented data should be distinguishable from the
 293 original data. Otherwise, it can easily cause the model to overfit. Therefore, we adopted random feature
 294 swapping to perturb the original data as much as possible, such that the model can learn more useful
 295 information from the spatial data.

296

297 **Comment 1.7.** *Regarding the data integration performance shown in Figure 4, why did the authors not*
 298 *compare many other methods designed for nonspatial scRNA-seq data, such as scVI (Lopez et al. 2018)*
 299 *and Harmony (Korsunsky et al. 2019), because those classical methods have been well demonstrated*
 300 *for good performance for single-cell data.*

301 **Response 1.7:** Thank you for your insightful comments. Following your comments, we added scVI and
 302 Harmony to our tests on the two mouse breast cancer datasets for vertical integration (Figure R11).
 303 Both scVI and Harmony were able to mix the two slices, but some batch differences were still visible
 304 post integration (Figure R11B). In comparison, GraphST evenly mixed the two slices, achieving better
 305 batch mixing than scVI and Harmony. We also quantitatively evaluated batch mixing with the iLISI metric
 306 where the higher iLISI score, the better the batch mixing. GraphST achieved a much higher iLISI score
 307 than Harmony and scVI (Figure R11C), confirming our visual observations. In the post-integration
 308 clustering, Harmony failed to align clusters across the two slices. scVI performed better than Harmony
 309 but some clusters were still not accurately mapped, such as clusters 1, 5, and 10. In contrast, GraphST's
 310 clusters highly overlapped between the two slices.

311 We also tested scVI and Harmony on one more mouse breast cancer sample (Figure R11D-F). Batch
 312 differences remained visible on the Harmony UMAP plot (Figure R11E). Comparatively, scVI removed
 313 the batch effects much better than Harmony. GraphST performed the best by evenly mixing the two
 314 slices. In terms of iLISI, Harmony significantly underperformed GraphST while scVI was comparable to
 315 GraphST (Figure R11F). Most of Harmony's clusters did not match across the two slices. While scVI
 316 generated clusters that were more consistent than Harmony, some clusters were fragmented,
 317 especially in section 2. GraphST again identified clusters that were spatially coherent and aligned well
 318 across the two slices. These results have been added to Figure 4 of the revised manuscript.



319
320

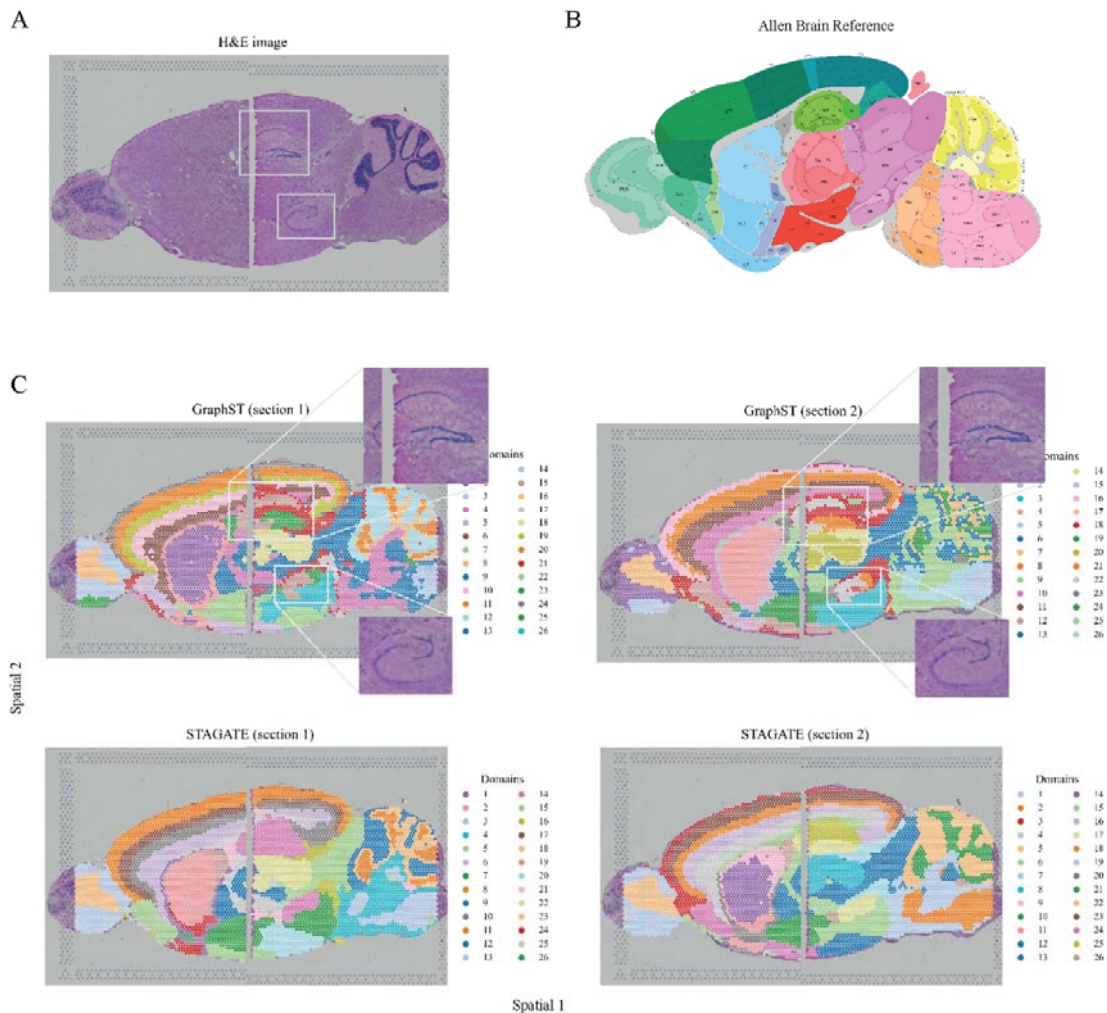
Figure R11. Vertical integration results of different methods on the two mouse breast cancer datasets.

321
322
323

Comment 1.8. It's important to show the STAGATE results that similar to Fig 4E to better demonstrate the data integration performance.

324
325
326
327
328
329
330
331
332
333

Response 1.8: Thank you for your great suggestions. Following your suggestions, we tested the performance of STAGATE on the two mouse brain samples. For fair comparison, we set the same number of clusters for all methods, i.e., 26 clusters. As shown in Figure R12, both GraphST and STAGATE produced continuous clusters that match the Allen brain reference well. Most importantly, like GraphST, the STAGATE's clusters were aligned along the edges of the anterior and posterior sections. However, some key brain regions were not represented in STAGATE's clusters. For example, STAGATE failed to identify the dorsal (top) and ventral (bottom) horn of the hippocampus regions highlighted with white boxes on the H&E images. In contrast, GraphST was able to reveal these regions. Overall, compared with STAGATE, GraphST performed slightly better in the horizontal integration task. We have added STAGATE's results to Figure 4 of the revised manuscript.



334

335 Figure R12. Horizontal integration results of different methods on two mouse brain samples.

336

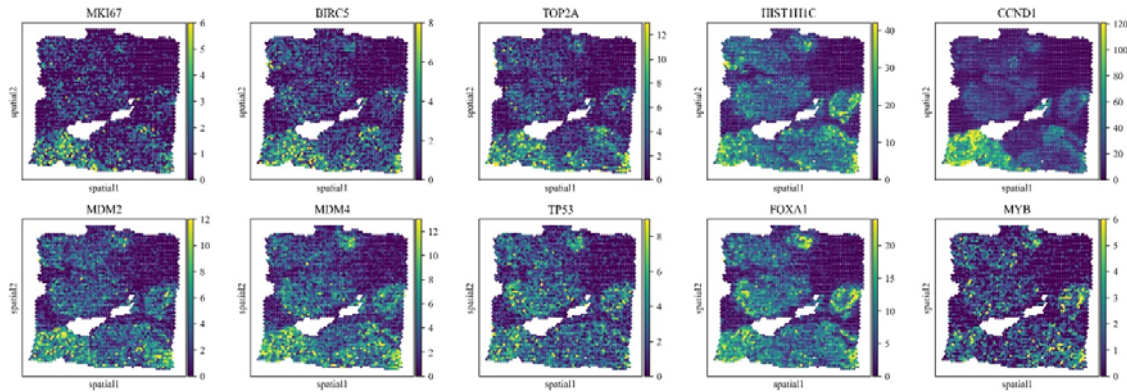
337 **Minor points:**

338 **Comment 1.9.** To better support the manual annotation result in Fig 6A, the spatial expression
 339 distribution of several marker genes for each panel in Fig 6A need to be added.

340 **Response 1.9:** The breast cancer tissue was ER positive, PR negative, Her2 positive, and diagnosed
 341 with ductal carcinoma in situ, lobular carcinoma in situ, and invasive carcinoma. Our pathologist
 342 collaborator produced the manual annotation based on the H&E image. Morphologically, it is easier
 343 to discern the IDC, DCIS, healthy, and tumour edge regions. As tumours usually harbour high cellular
 344 heterogeneity, it is challenging to find known gene expression markers that can distinguish IDC from
 345 DCIS. Here, we selected several reported marker genes of breast cancer from literature and plotted
 346 their spatial expression to support our manual annotation (Figure R13). We have added the marker
 347 genes to Supplementary Figure S13.

348 **Reference**

349 The Cancer Genome Atlas Network. Comprehensive molecular portraits of human breast tumours,
 350 Nature 2012.



351
352 Figure R13. Spatial expression distribution of reported breast cancer markers

353
354 **Comment 1.10.** *In line 836, the author mentioned the first loss term indicates contrastive loss, why is*
355 *there only one instead of two terms? What is the meaning of the first term?*

356 **Response 1.10:** Thank you for your comments. Formula (10) is the overall objective function of
357 GraphST' third module for cell type deconvolution of ST data. In our framework, we use two different
358 contrastive learning methods for the first and third modules, respectively. Motivated by Zhang et al.
359 (2022), we use an augmentation-free contrastive learning method for the third module. Therefore, there
360 is only one term for contrastive loss. The first term (i.e., contrastive loss) in formula (10) is an InfoNCE
361 objective function that aims to maximize the similarities of positive pairs and minimize those of negative
362 pairs. We have added more details in the revised manuscript to describe the contrastive learning
363 method of the third module.

364 **Reference**

365 Zhang et al. Dual Temperature Helps Contrastive Learning Without Many Negative Samples:
366 Towards Understanding and Simplifying MoCo. CVPR2022.

367
368 **Comment 1.11.** *In line 712, which norm is used? L1 or L2 or others?*

369 **Response 1.11:** Thank you for your comments. The norm term in formula (3) is L2-norm. We have
370 revised formula (3) in the revised manuscript.

371
372 **Comment 1.12.** *In Fig 3C, the titles of panels Mesenchyme and Dermomyotome seem misplaced.*

373 **Response 1.12:** Thank you very much for pointing out this. We have changed the titles of these two
374 panels. Please refer to Figure 3C in the revised manuscript.

375
376 **Comment 1.13.** *All color bars need to be explained for their meanings.*

377 **Response 1.13:** Thank you for your comments. We have added legends to all colour bars in the revised
378 manuscript.

379
380 **Comment 1.14.** *Many typos and grammar errors in the manuscript, e.g., in line 28, "has" should be*
381 *"have"; in line 59, "K-means" should be "k-means"*

382 **Response 1.14:** Thank you very much for your careful reading. We have carefully gone through our
383 manuscript and corrected the typographical and grammatical errors.

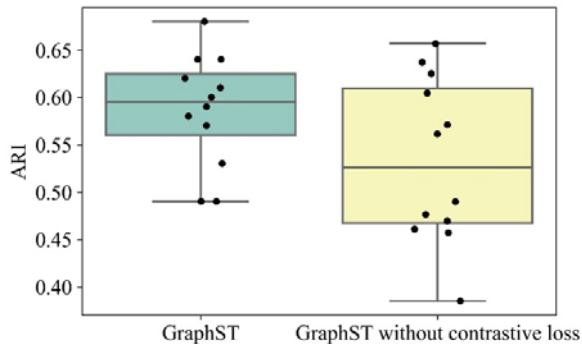
384
385 **Comment 1.15.** *Lines 124-126, "Using self-supervised contrastive learning improves performance in*

386 *learning relevant latent features and has the additional benefit of removing batch effects”. This sentence*
387 *occurs without any supporting evidence. It needs to be fixed.*

388 **Response 1.15:** Thank you very much for pointing out this. We have revised this sentence as follows.

389 *“Using self-supervised contrastive learning improves performance in learning relevant latent features.”*

390 To demonstrate the contribution of self-supervised contrastive learning, we conducted an ablation study
391 by comparing GraphST with a variant of GraphST without contrastive loss on the DLPFC dataset.
392 Without contrastive loss, the performance of GraphST is significantly reduced (Figure R14), indicating
393 that contrastive loss contributes to the performance improvement of our GraphST model. These results
394 have been added to Supplementary Figure S14C.



395

396 Figure R14. Comparison between GraphST and its variant GraphST without contrastive loss on the
397 DLPFC dataset.

398

399 **Comment 1.16.** *In lines 158-162, the authors introduced “self-reconstruction loss” and “contrastive loss”*
400 *and their effects. It’s important to show what the two losses are in the context of biology.*

401 **Response 1.16:** Thank you very much for highlighting this point. In our GraphST framework, we take
402 the gene expression matrix and spatial graph as inputs. The gene expression matrix contains the
403 feature information of spots while the spatial graph stores the spatial adjacency of spots. Our model is
404 a GNN-based model that aims to integrate the gene expression of spots with their corresponding spatial
405 information for spatial clustering. The key feature information is in the gene expression matrix which
406 should be retained. Therefore, we design the self-reconstruction loss to enforce the preservation of the
407 original gene expression information in the reconstructed expression.

408 The contrastive loss design is based on the assumption that a spot in the spatial data usually has a cell
409 type label similar to its local context, e.g., one-hop or two-hop neighbours. As discussed above in the
410 contrastive learning part, we define positive/negative pairs by pairing spot embedding h_i/h'_i from the
411 original/corrupted graph with its local summary vector s_{local} . The local summary vector s_{local} represents
412 the local context of a spot and is obtained by a sigmoid of the mean of all its neighbours’ embeddings.
413 The main goal of contrastive learning is to make spot embedding h_i close to its local context s_{local}
414 from the original graph. Therefore, trained with contrastive loss, spatially adjacent spots will have similar
415 embeddings while non-adjacent spots will have dissimilar embeddings.

416 Based on your comment, we have added more biological contexts when describing self-reconstruction
417 loss and contrastive loss in the revised manuscript.

418

419 **Comment 1.17.** *It’s unclear how the neighbor graph is constructed. In the caption of Fig. 1 (lines 1034-*
420 *1035), the authors wrote “...neighbor graph constructed using spot coordinates (x,y) of that fall within a*
421 *distance threshold”. However, in the method section in lines 665-675, the authors wrote “Finally, we*
422 *select the top k-nearest spots as its neighbors”. It’s unclear whether the authors used a distance*
423 *threshold or a threshold for k.*

424 **Response 1.17:** We are sorry for the confusion. In our framework, we use a threshold for k when
425 constructing the neighbourhood graph. We have carefully gone through the manuscript and ensured
426 consistency in the revised manuscript.

427

428 **Comment 1.18.** *Regarding the method (lines 655-659), the descriptions seem to be for the spatial*
429 *transcriptomics data. However, this is not clear from the description, as two kinds of datasets (spatial*
430 *transcriptomics data and scRNA-seq data) are mentioned in this paper.*

431 **Response 1.18:** Thank you for your comments. GraphST was developed for three analysis tasks,
432 spatial clustering, multiple ST data integration, and cell type deconvolution of spatial data. Cell type
433 deconvolution is achieved by projecting single-cell RNA-seq data onto the spatial data. Therefore, in
434 addition to spatial data, single cell RNA-seq data is also used in the third module of our framework.

435

436 **Comment 1.19.** *In lines 684-688, "...while keeping the original graph structure unchanged": was the*
437 *corrupted neighbor graph G' the same as the original G ?*

438 **Response 1.19:** Thank you for your comments. Yes, the adjacency matrix of the corrupted
439 neighbourhood graph G' is the same as the original G . When creating the corrupted neighbourhood
440 graph G' , we only randomly shuffle feature vectors between spots while keeping the graph's topological
441 structure unchanged. For example, the entire feature vector of node A is assigned to node B and vice
442 versa.

443

444 **Comment 1.20.** *In lines 708-709, " W_d and b_d represent the trainable weight matrix and bias vector,*
445 *respectively, which are shared by all nodes in the graph". Please justify why W_d and b_d need to be*
446 *shared by all nodes in the graph. Besides, is this the same case for W_e and b_e ?*

447 **Response 1.20:** Thank you very much for your insightful comments. In our model, without loss of
448 generality, the trainable weight matrices W_e, W_d and bias vectors b_e, b_d are shared by all nodes in the
449 graph. In the GNN model, the dimensions of the weight matrices and bias vectors are usually very large
450 depending on the number of nodes of the input graph. For example, in the datasets used in our
451 manuscript, the smallest (i.e., DLPFC slice #151676) has 3460 spots and the largest (i.e., Mouse
452 embryo E14.5) has 92,928 bins. If we use a different weight matrix and bias vector for each node, it will
453 be very challenging to train the model. Sharing the weight and bias significantly reduce the number of
454 weight and bias terms used, making it easier to train the model. It also helps reduce the running time.

455

456 **Reference**

457 Korsunsky, Ilya, Nghia Millard, Jean Fan, Kamil Slowikowski, Fan Zhang, Kevin Wei, Yuriy Baglaenko,
458 Michael Brenner, Po-Ru Loh, and Soumya Raychaudhuri. 2019. "Fast, Sensitive and Accurate
459 Integration of Single-Cell Data with Harmony." *Nature Methods* 16 (12): 1289–96.
460 Lopez, Romain, Jeffrey Regier, Michael B. Cole, Michael I. Jordan, and Nir Yosef. 2018. "Deep
461 Generative Modeling for Single-Cell Transcriptomics." *Nature Methods* 15 (12): 1053–58.
462 Veličković, Petar, William Fedus, William L. Hamilton, Pietro Liò, Yoshua Bengio, and R. Devon Hjelm.
463 2018. "Deep Graph Infomax." *ArXiv [Stat.ML]*. arXiv. <http://arxiv.org/abs/1809.10341>.

464

465

466

467

468

469

470

471 **Reviewer #2** (Remarks to the Author):

472

473 *In the manuscript "GraphST: A novel graph" by Long and coworkers the authors develop a new*
474 *method for better describing spatial transcriptomics data and being able to integrate multiple studies.*

475 *The method is based on graph neural networks and contrastive learning, which makes it possible to*
476 *combine scRNA-seq of better resolution and spatial transcriptomics. The method is logically sound and*
477 *makes a lot of sense. Moreover, authors show that it empirically identifies more relevant clusters and*

478 *allows data integration for higher power. Although, I am not an expert in spatial transcriptomics these*
479 *problems seem of great importance and authors spend good effort to show that it works a planned.*

480 *Having said that my expertise is in neural networks and translational bioinformatics I believe that the*
481 *paper would be a good contribution to the spatial transcriptomics field. From my side, I have no concerns*
482 *of the paper and like to see it published.*

483 **Response:** We thank the reviewer for the above positive comments. We welcome you to provide
484 additional valuable comments.

485

486

Reviewers' Comments:

Reviewer #1:

Remarks to the Author:

The revision has fully addressed my comments and concerns. Well done.