

SpaMGCN:基于多视图图神经网络 识别空间域的方法*

刘和鑫¹,尚文婧¹,赵翔宇¹,郑一帆¹,张佳²,冯振兴¹

(1. 内蒙古工业大学理学院,呼和浩特 010051; 2. 内蒙古大学生命科学学院,呼和浩特 010021)

摘要:作为空间转录组学关键任务之一,空间域识别旨在精确划分基因表达与空间关联的组织区域。但多数方法难以兼顾局部邻域与全局结构特征。为此,本研究提出了多视图图卷积网络方法 SpaMGCN,通过引入双模态图卷积网络,同步提取基因表达相关性与空间邻域的局部与全局特征。在多个公开数据集的验证中,SpaMGCN实现了对复杂生物学组织区域精细的边界划分,在空间域识别方面展现优越的性能,为理解生物组织的空间异质性提供了高效的计算工具。

关键词:空间转录组学;空间域识别;多视图图神经网络

中图分类号:Q811.4; TP391.4; O29 **文献标志码:**A

空间转录组学(Spatial transcriptomics, ST)技术通过整合基因表达谱、空间坐标信息与组织形态图像,为精确解析生物组织内基因表达的空间分布模式提供了机会^[1]。然而,ST数据具有高维、高稀疏、高噪声等特点,需要高效准确的算法和计算资源进行处理。作为解析基因表达空间异质性和细胞互作关系的核心要素,空间位置信息不仅明确了检测位点的组织解剖定位,更为揭示组织微环境中的功能分区提供了几何坐标基础。

空间域识别旨在划分具有一致性基因表达和组织形态特征的空间连续区域。精准识别空间结构域是解析组织异质性和细胞功能机制的重要前提。近年来,基于深度学习的空间域识别方法在处理高维大规模数据时展现出显著优势^[2]。其中,图卷积网络(Graph convolutional network, GCN)通过对空间位点的特异性拓扑关系进行建模,成为了当前处理复杂ST数据的主流工具^[3]。SpaGCN算法通过融合组织学特征将二维坐标扩展到三维来构建无向加权图,并结合迭代聚类算法实现对空间异质性的解析^[4]。conST通过引入了深度图信息框架,利用对比学习技术提高了图卷积网络编码层所捕获的低维潜在嵌入的质量^[5]。STAGATE通过自适应学习邻域斑点的注意力权重的编码器实现空间信息的无监督聚合^[6]。stLearn能够通过非负矩阵分解和自监督学习方法利用图卷积网络提取空间潜在分布模式,后结合位点的对应基因表达信息实现了对组织微环境的精细解析^[7]。DeepST构建多模态融合框架,利用预训练卷积网络提取组织图像特征生成增强矩阵,结合变分图自编码器与去噪自编码器提升潜在表征的判别能力^[8]。GraphST通过设计多层次的神经网络结构实现了以端到端的学习方式自动提取复杂的空间模式^[9]。SEDR可以在无基础真相的情况下,通过特

* 收稿日期:2025-05-08; 修回日期:2025-06-24

基金项目:内蒙古自然科学基金面上项目(2024MS06027);自治区直属高校基本科研业务费(JY20230067);自治区级大学生创新创业训练计划项目(S202410128012)

作者简介:刘和鑫(1999—),男,内蒙古包头人,2023级硕士研究生。E-mail:2927493054@qq.com

通信作者:冯振兴(1988—),男,河南遂平人,副教授,博士。主要从事生物数学、生物统计学算法开发等方面研究。E-mail:zxfeng@imut.edu.cn

定的对比学习目标函数约束网络学习到鲁棒的空间特征表示^[10]。EfNST通过EfficientNet架构来优化空转数据的图像学信息,通过结合图像信息来辅助增强模型对空间域识别的精度^[11]。

尽管上述方法通过图神经网络实现了多模态数据的整合,但这些算法采用的单视图模型提取特征的视角较为单一,不能充分利用全局结构信息。针对上述挑战,本研究提出了一种基于多视图图卷积网络的空间域识别方法SpaMGCN。该方法通过融合空间邻近依赖与基因表达相关性的双尺度结构建模,实现对复杂ST数据的多维度特征解析,从全局视角揭示组织结构的层级化信息。SpaMGCN以共享多视图图卷积网络为基础框架,同步学习ST数据的多模态结构特征,并嵌入深度聚类模块通过增强潜在表征的紧凑性优化特征空间,从而提升模型对数据复杂关联性的捕捉能力。在多个公开数据集上的基准实验表明,相较于现有先进方法,SpaMGCN在空间域识别精度、聚类表征的几何紧凑性及类别可分性等关键指标上均展现出显著优势。

1 模型介绍和方法

1.1 SpaMGCN模型介绍

SpaMGCN是一个用于解读ST数据的多视图图卷积框架(图1),通过整合空间位置信息与基因表达谱,实现对ST数据的联合特征表示的深度学习。模型工作流程包含四个核心模块:图构建、特征提取、特征整合与特征优化(图1-A)。首先,在图构建模块,SpaMGCN构建的空间邻接图和特征邻接图从不同角度刻画了斑点间的相关关系。对于空间邻接图,采用 r 半径定义点间的空间依赖性,为局部特征提取提供空间结构约束;对于特征邻接图,通过余弦相似性计算点间的表达谱相关性,构建不依赖空间位置的全局关联网,用于捕获空间距离较远但基因表达模式相似的斑点特征。其次,在特征提取模块,模型基于多视图图卷积网络架构,将空间邻接图与特征邻接图分别映射至潜在表示空间,生成包含原位表达模式的空间特征与反映全局表达相关性的基因特征,实现特征信息的并行提取。接着,在特征整合模块,模型采用线性融合策略对空间特征与基因特征进行深度整合生成可为后续分析提供兼具局部特异性与全局关联性的特征表示。最后,在特征优化模块,模型通过解码器网络重构原始基因表达谱,以重构损失约束潜在特征空间,迫使模型学习到能够准确还原原始表达信息的紧凑特征表示^[12];并结合深度嵌入聚类模块^[13],将无监督聚类目标融入训练过程,通过优化聚类损失函数,进一步提升嵌入特征的类间区分度与类内紧凑性。在下游分析方面(图1-B),SpaMGCN能够对最终的嵌入特征进行多任务ST数据分析:包括空间域识别、可视化表征、标记基因功能验证、聚类结果细化及空间功能域深度解析等。

1.2 空间邻接图与特征图的构建

1.2.1 空间邻接图构建

通过采用组织中斑点的空间位置,来衡量相邻斑点间的空间相似性。对于第 i 个斑点 S_i ,其在组织切片中的空间位置用二维坐标 (x_i, y_i) 表示,为充分利用空间位置信息,构造空间邻接图 $G(A_s, X)$,该图由空间邻接矩阵 A_s 和基因表达矩阵 X 组成。其中空间邻接矩阵 A_s 由每个斑点的空间坐标确定,具体而言,利用 r 半径方法,将空间位置坐标转化为矩阵形式,参数 r 用于确定邻接图的紧凑性。为准确描述空间关系,令每个斑点与最近的六个邻斑点连接,距离度量方式采用欧氏距离。 $A_s \in R^{N \times N}$ 为包含 N 个斑点的空间邻接矩阵, $X \in R^{N \times M}$ 表示标准化的基因表达矩阵,其中 M 为过滤后的基因个数,如果斑点 S_i 和斑点 S_j 之间的欧氏距离小于预定义的半径 r ,则设置 $A_s^{ij} = A_s^{ji} = 1$;否则设置 $A_s^{ij} = A_s^{ji} = 0$ 。欧式距离计算公式如下:

$$d(S_i, S_j) = 1 / \left(1 + \sqrt{(x_i - x_j)^2 + (y_i - y_j)^2} \right)$$

1.2.2 特征邻接图构建

为充分利用基因表达的潜在结构信息,构建特征邻接图 $G(A_f, X)$ 。特征邻接图由基因表达矩阵

A_f 和基因表达谱 X 组成。其中,基因表达矩阵 A_f 采用 K 近邻方法构建,邻居斑点间距离度量方式采用余弦相似性来确定,默认设置 $k=6$ 。若斑点 S_i 是斑点 S_j 的邻居,则设置 $A_f^{ij}=1$; 否则设置 $A_f^{ij}=0$ 。为进一步简化计算过程,并有效提取主要特征以增强分析的准确性,采用主成分分析法(Principal component analysis, PCA)降低基因表达谱的维度。余弦相似度计算公式如下:

$$\text{sim}(S_i, S_j) = (S_i \cdot S_j) / (|S_i| \cdot |S_j|)$$

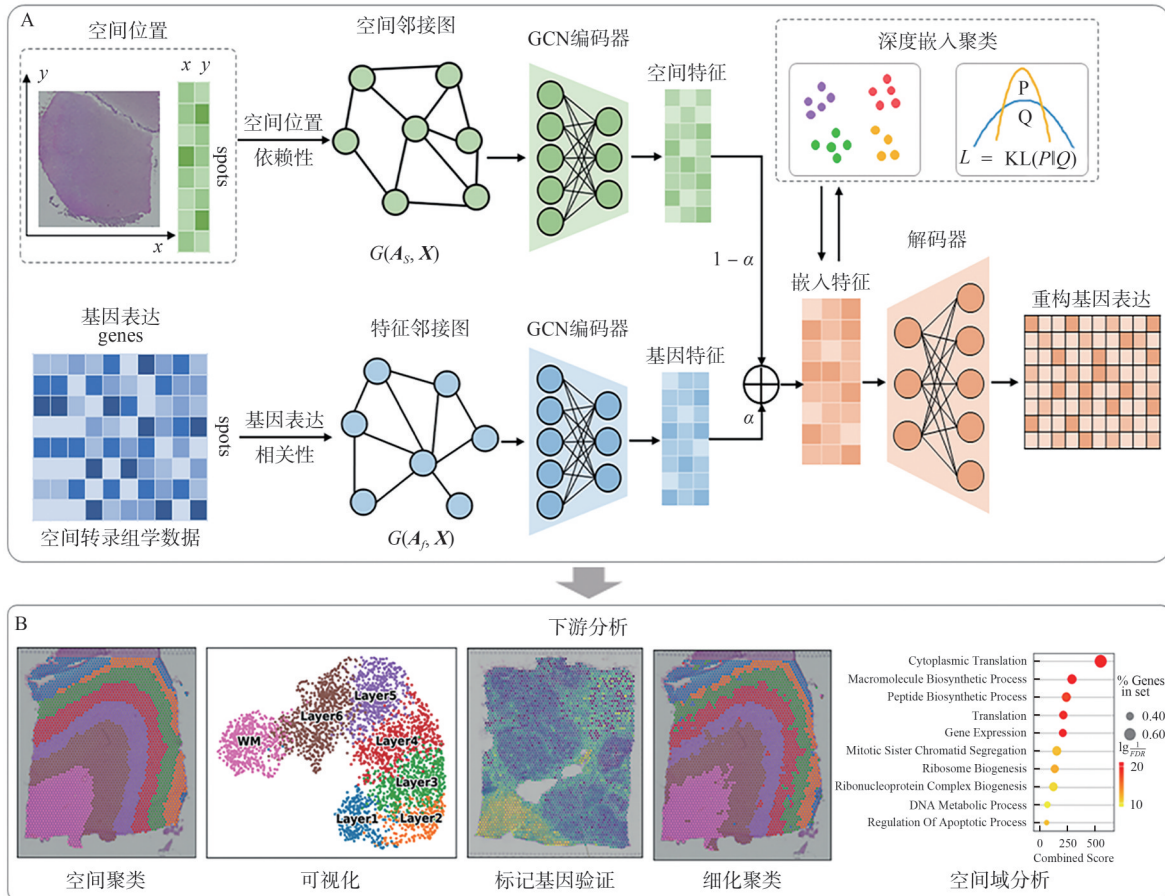


图 1 SpaMGCN 的流程

Fig. 1 Overview of the SpaMGCN

1.3 多视图图卷积自编码器

图卷积网络作为一种功能强大的图神经网络,能够直接处理图数据,并有效利用图结构信息。它可以聚合邻居的信息,捕获节点之间的依赖关系,进而生成富有信息的嵌入表示。为了从基因表达和空间结构中精准提取最为相关的信息,本研究采用多视图 GCN 编码器,对空间邻接图和特征邻接图进行卷积操作。多视图 GCN 编码器由 4 个部分构成:

1)空间卷积 为了实现基因表达信息与空间位置信息的有机结合,同时捕捉空间邻接信息,对空间邻接矩阵 A_s 和基因表达 X 执行卷积操作,以此聚合邻居节点的空间信息。接下来,多层空间卷积网络遵循以下层次传播规则:

$$H_s^{(l+1)} = \text{ReLU}(\tilde{D}_s^{-\frac{1}{2}} \tilde{A}_s \tilde{D}_s^{-\frac{1}{2}} H_s^{(l)} W_s^{(l)})$$

式中: $H_s^{(l)}$ 表示斑点在第 l 层的特征,初始 $H_s^{(0)} = X$; $\text{ReLU}(\cdot)$ 为激活函数,用于进行非线性变换; $\tilde{A}_s = A_s + I$ 为添加自循环的空间邻接矩阵,其中 A_s 为斑点 i 对应的空间邻接矩阵, I 为单位矩阵; \tilde{D}_s 为 A_s 对应的度矩阵; $W_s^{(l)}$ 为第 l 层的权重。

2)特征卷积 与空间卷积类似,为了获得更全面的基因表达信息,并推算出特征图中斑点的基因表达,对特征邻接矩阵 A_f 和基因表达 X 进行特征卷积,其运算公式如下:

$$H_f^{(l+1)} = \text{ReLU}(\tilde{D}_f^{-\frac{1}{2}} \tilde{A}_f \tilde{D}_f^{-\frac{1}{2}} H_f^{(l)} W_f^{(l)}).$$

式中: $H_f^{(l)}$ 表示斑点在第 l 层的特征,初始值 $H_f^{(0)} = X$,用于实现非线性变换; $\tilde{A}_f = A_f + I$ 为添加自循环的空间邻接矩阵, A_f 为斑点 i 对应的空间邻接矩阵, I 为单位矩阵; \tilde{D}_f 为 \tilde{A}_f 对应的度矩阵;其中 $W_f^{(l)}$ 为第 l 的权重。

此外,为了实现空间特征和基因表达特征的有效整合,采用线性相加得到最终的嵌入特征:

$$H = (1 - \alpha) H_f^{(l+1)} + \alpha H_s^{(l+1)}.$$

式中: H 为嵌入特征; $H_f^{(l+1)}$ 为空间特征; $H_s^{(l+1)}$ 为基因特征; α 表示基因特征权重的超参数。

3)解码器 解码器的运算公式如下:

$$X' = H W_d + b_d.$$

式中: X' 是重构的基因表达矩阵; H 是最终的嵌入特征; W_d 和 b_d 分别是解码器的权重和偏置。

4)重构损失 重构损失用于衡量重构基因表达与原始基因表达之间的差异。通过最小化如下重构损失来确保所得的嵌入捕获了足够的生物学信息:

$$L_{rec} = \frac{1}{N} \sum_{i=1}^N (X_i - X'_i)^2.$$

式中: N 是斑点的数量; X_i 与 X'_i 分别是原始基因表达矩阵与重构基因表达矩阵中的第 i 个样本。

1.4 深度嵌入聚类

SpaMGCN运用深度嵌入聚类方法,以迭代的方式将斑点划分到不同的组别中,旨在优化聚类结果。具体而言,采用scikit-learn库中的K-means^[14]算法,依据学习得到的潜在表征来初始化聚类中心。设 C 为初始簇的个数。使 $i = 1, 2, \dots, N$ 来索引这些点, $u, k = 1, 2, \dots, C$ 则用于标记初始簇。

第一步计算软赋值,通过Student's t-distribution计算软赋值 q_{iu} ,用来评估点嵌入 H_i 和聚类中心嵌入 μ_u 之间的相似度:

$$q_{iu} = \frac{(1 + \|H_i - \mu_u\|^2)^{-1}}{\sum_{k=1}^C (1 + \|H_i - \mu_k\|^2)^{-1}}.$$

第二步迭代细化聚类,其所采用的关键技术是辅助目标分布(Auxiliary target distribution)方法^[15]。该方法考虑到在聚类过程中,直接使用原始的软分配可能会导致聚类中心不稳定或聚类质量不高的问题。算法通过引入一个更合理的目标分布,引导模型优化聚类分配,即对于那些已经被模型分配到某个聚类且置信度较高的样本,给予更大的权重,以强化这些分配。通过强调高置信度样本,辅助目标分布有助于模型学习更清晰的聚类边界。避免了某些聚类中心因样本过多而主导学习过程,导致其他聚类中心学习不足的问题。相比直接使用软分配,辅助目标分布可以减少噪声样本对聚类中心更新的干扰,提高算法的稳定性。

具体的,我们通过基于 q_{iu} 的辅助目标分布 p_{iu} ,从具有高置信度的分配中进行学习,从而对聚类进行迭代细化。辅助目标分布 p_{iu} 的计算方式为:

$$p_{iu} = \frac{q_{iu}^2 / \sum_{i=1}^N q_{iu}}{\sum_{k=1}^C (q_{ik}^2 / \sum_{i=1}^N q_{ik})}.$$

第三步定义目标函数,利用软分配 q_{iu} 和辅助目标分布 p_{iu} ,通过KL散度定义目标函数 L_{DEC} ,其表达式为:

$$L_{DEC} = \text{KL}(P \| Q) = \sum_{i=1}^N \sum_{u=1}^C p_{iu} \ln \frac{p_{iu}}{q_{iu}}.$$

这一过程通过不断迭代优化该目标函数,实现对斑点的有效聚类,使同一簇内的斑点在潜在表征空间中具有更高的相似性,不同簇之间的差异更加显著,使得嵌入特征更适用于聚类任务。

2 结果与分析

2.1 SpaMGCN 有效检测了人类背外侧前额叶皮层的层级结构

ST 数据的结构解析精度高度依赖模型的表征学习能力。由于人类背外侧前额叶皮层(DLPFC)^[16]数据集拥有清晰的层次结构,所以本研究基于 DLPFC 数据的 12 个切片(图 2-A 为切片#151673),对 SpaMGCN 的表示学习性能开展系统性评估。实验选择 7 个算法 STAGATE、SpaGCN、Scanpy^[17]、conST、stLearn、Seurat^[18]、EfnST 进行对比,采用调整兰德指数 ARI^[19]量化 12 个切片的空间域识别精度(图 2-B)。结果显示,SpaMGCN 以平均 $ARI=0.50$ 的性能显著优于对比方法。

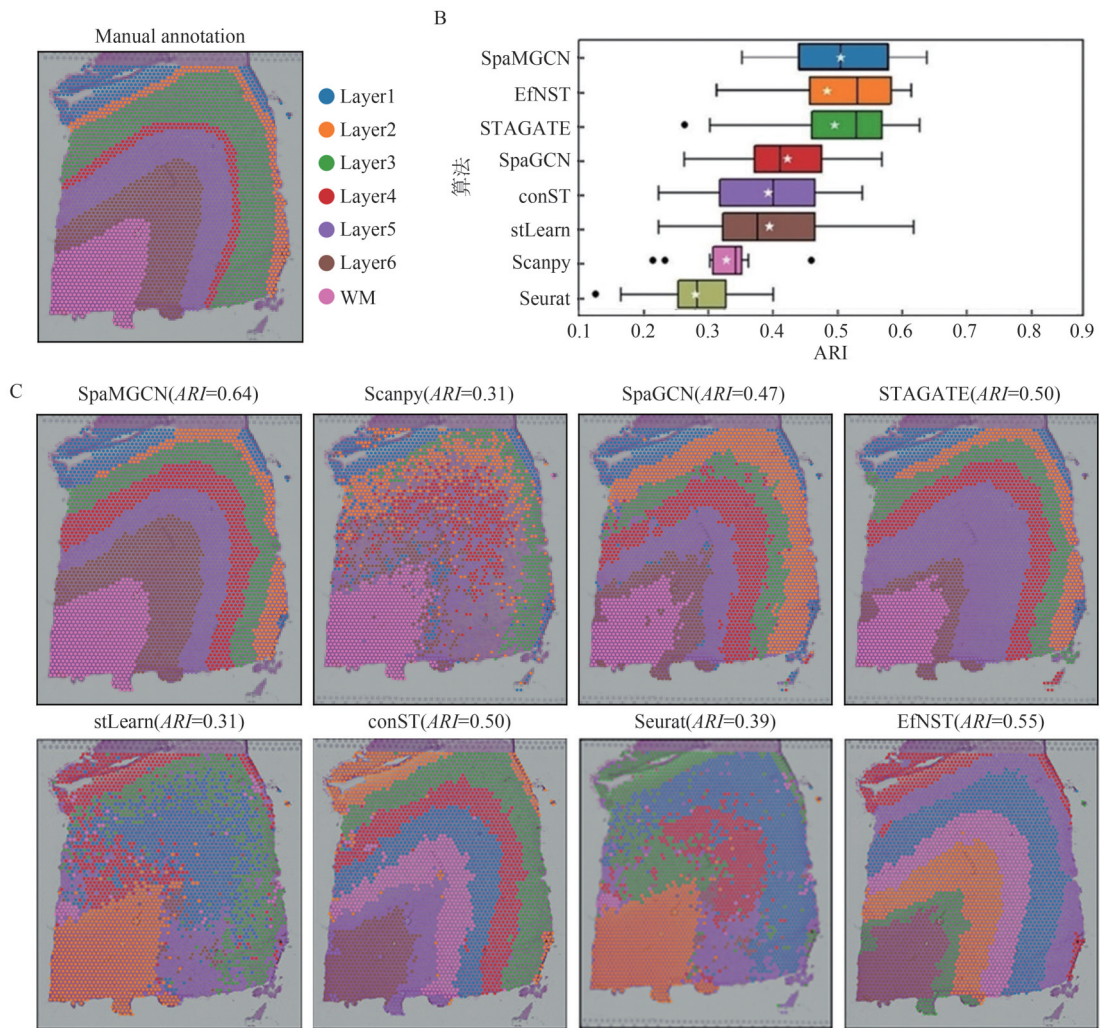


图 2 SpaMGCN 在 DLPFC 数据上的聚类结果与比较

Fig. 2 Clustering results and comparison of SpaMGCN on DLPFC dataset

以包含完整七层结构(6个皮质层+WM层)的切片#151673为例,SpaMGCN 的域识别性能呈现显著优势(图 2-C),且唯一实现了与基准真相高度一致的 Layer_5、Layer_6 与 WM 区域的划分。反观 Scanpy,因缺乏空间约束,聚类结果呈现离散分布,未能检测到 2~6 层的结构边界,严重阻碍了空间域的有效解析;stLearn 的 SME 聚类高度依赖 HE 染色图像质量,图像信息的噪声会导致形态学特征提取偏差,所以 stLearn 的聚类结果 Layer_1~Layer_6 以及 WM 的层级结构边界模糊;SpaGCN 一定程

度反映了切片的层次结构,但是层间划分不明确且存在严重的斑点混杂现象;STAGATE与EfnST的聚类结果有明显的层次结构,但与SpaMGCN相比STAGATE与EfnST的聚类结果层与层的边界模糊,这可能是由于EfnST过度依赖于图像质量,STAGATE没有充分利用空间信息所引起的;Seurat的降维聚类过分依赖于基因表达信息,对空间坐标的利用停留在“邻域平均”层面,难以捕捉细胞间复杂的空间依赖关系,这使得除了WM区域外,Seurat的聚类结果无法看出明显的层级结构。

我们对切片#151673进行了UMAP^[20]可视化分析(图3)进一步印证了这一结论。其中UMAP是常用的数据降维和可视化工具,通过对高维数据进行非线性映射,保持数据点之间的邻近关系。在空间转录组学中,每个观测位点的基因表达谱被视作一个高维向量,UMAP方法在保留数据局部邻域关系以及全局拓扑结构的前提下,将这些向量投影到共享的低维空间中,形成易于观察的聚类模式,使得我们可以在2维空间来观察数据的聚类结果。在UMAP图中,Scanpy、stLearn以及Seurat的皮质层斑点呈混杂聚集;SpaGCN仅能区分WM与Layer_6;STAGATE的Layer_1与Layer_2边界模糊;conST的Layer_2到Layer_4边界混杂;STAGATE与EfnST都能很好地区分Layer_4到Layer_6以及WM,但是Layer_1到Layer_3边界混杂;而SpaMGCN清晰呈现各层的有序排列,精准反映了皮层从浅至深的发育层级^[21-22]。

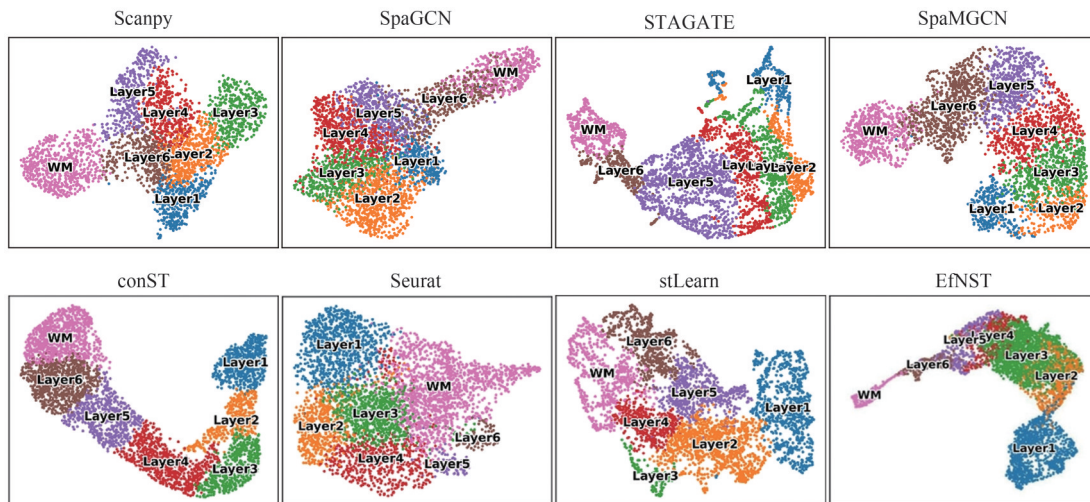


图3 #151673切片UMAP可视化效果图

Fig. 3 UMAP visualization of slice #151673

2.2 SpaMGCN深入解析了人类乳腺癌的组织空间异质性

乳腺癌组织具有复杂的结构和异质性,利用空间转录组学技术解析肿瘤微环境,可以更好地理解不同细胞类型之间的空间关系,尤其在揭示肿瘤内/间异质性及组织紊乱机制中具有不可替代的价值^[23]。本研究基于人类乳腺癌组织的10x Visium数据集,验证SpaMGCN对复杂癌症组织的解析能力。该数据集包含20个精细标注的空间域,可归为四种形态学类别:浸润性导管癌(IDC)、健康组织(Healthy)、导管原位癌/小叶原位癌(DCIS/LCIS)和肿瘤周围低恶性特征区域(Tumor_edge)^[24](图4-A)。

在多指标量化性能对比实验中,SpaMGCN在多项聚类指标如调整互信息(Adjusted mutual information, AMI)^[25]、标准化互信息(Normalized mutual information, NMI)^[26]和福克斯·马洛斯指数(Fowlkes-Mallows index, FMI)^[27]中显著高于其他方法(图4-B),在空间域识别可视化中(图4-C),SpaMGCN展现出与人工标注极高的空间一致性,其识别的区域1(IDC_4)、区域8(Tumor_edge_2)等关键病灶区域,不仅边界划分清晰连续,且域内斑点呈现生物学意义上的空间聚集性;值得注意的是EfnST在分析人类乳腺癌异质性任务上取得了不错的效果,可以识别出部分较为复杂的组织结

构,这可能得益于组织学图像对于复杂组织结构识别的辅助作用,相较之下,SpaGCN、stLearn与Seurat的聚类结果存在严重的域间混淆,健康组织与肿瘤边缘区域的点混杂分布。Scanpy、STAGATE和conST虽能识别主要肿瘤区域,但存在少量离群点干扰,且不同区域间边界划分为较粗糙。

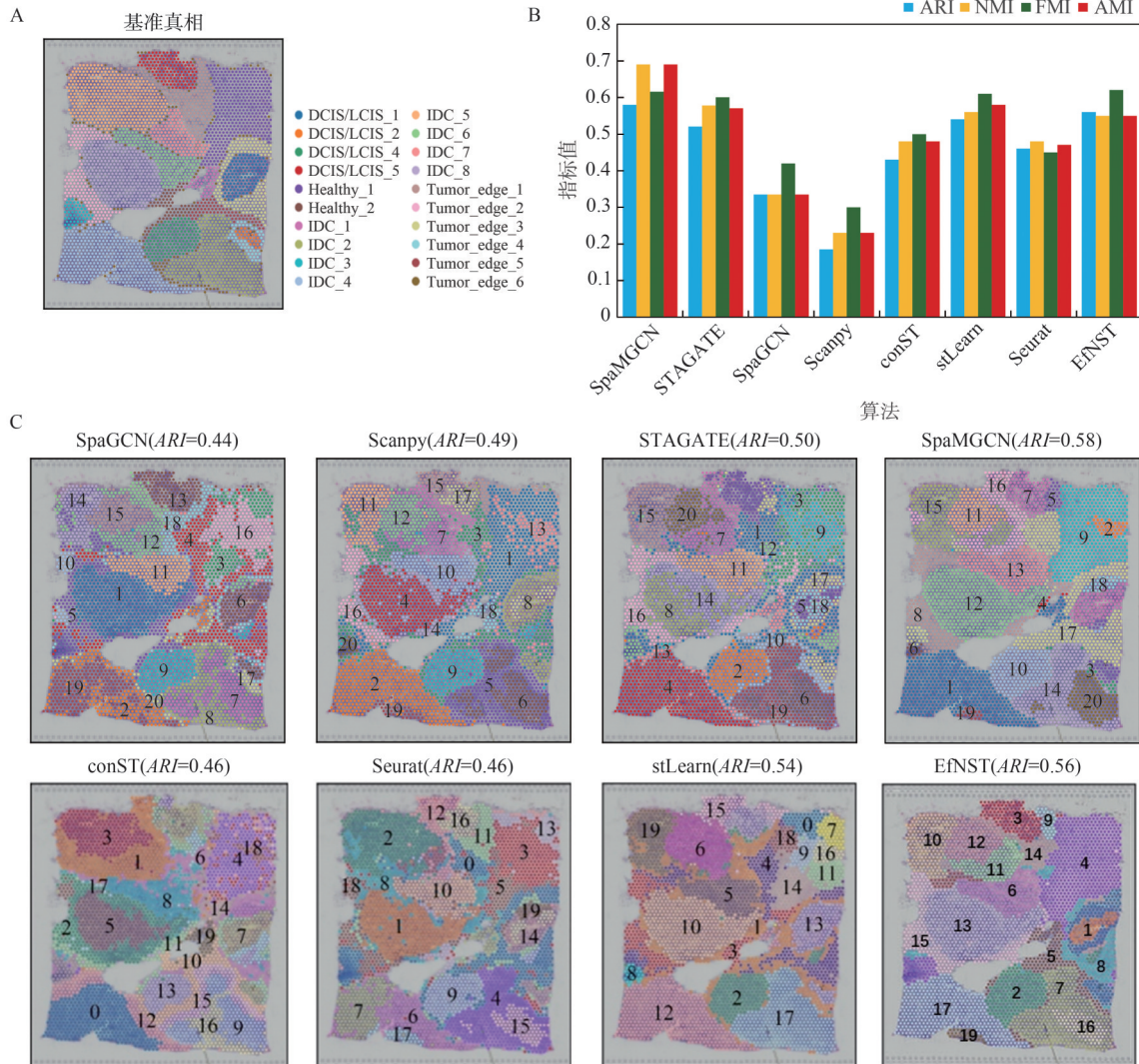


图 4 SpaMGCN在乳腺癌数据上的结果与比较

Fig. 4 Clustering results and comparison of SpaMGCN on breast cancer dataset

为深入剖析人类乳腺癌组织的空间结构异质性与基因表达模式关联,本研究从空间距离变异、空间可变基因和区域标记基因三个维度,对人类乳腺癌组织的空间聚类结果展开深度解析。首先基于Ripley's L函数^[28]量化各空间域的聚类或分散点分布的空间格局。发现区域1(IDC_4)、区域9(Healthy_1)及区域12(IDC_8)呈现显著的空间聚集特征(图5),而其他域多表现为随机或离散分布。这种聚集性差异暗示了肿瘤微环境中不同功能区域的空间组织特性。

针对上述高聚集性区域,进一步筛选前10个差异表达基因并开展功能富集分析^[29]。分析结果表明,空间域1与细胞增殖和炎症反应密切相关,在免疫细胞的迁移与定位过程中发挥作用,这可能涉及肿瘤微环境中免疫细胞的活性以及其对肿瘤生长的抑制作用。空间域9在免疫反应、肿瘤免疫监视以及抑制肿瘤生长方面具有重要意义;空间域12则与免疫反应、细胞外基质重塑和细胞黏附相关,参与细胞分化与转移过程免疫系统的功能紧密相连,特别是在抗原呈递、T细胞激活和免疫反应调节方面(图6)。

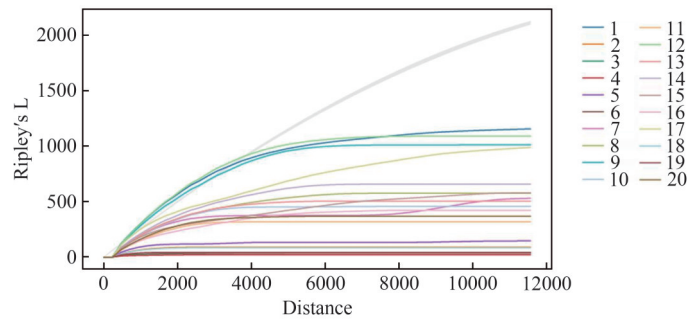


图5 Ripley's L函数计算乳腺癌的分布特征

Fig. 5 Calculation of distribution characteristics of breast cancer using Ripley's L function

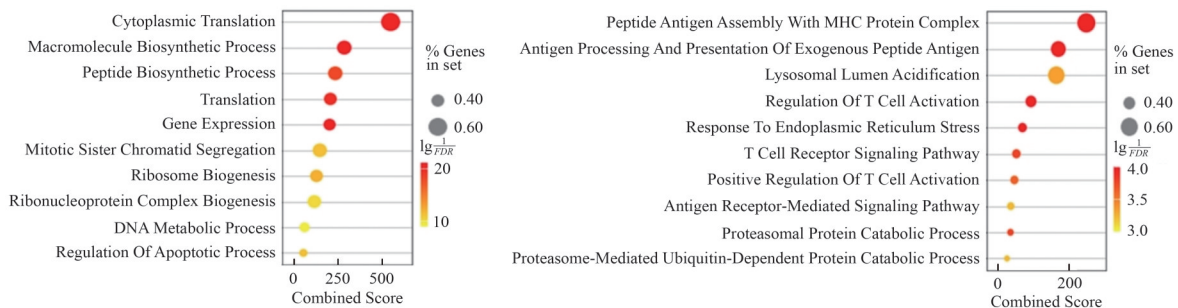


图6 乳腺癌区域1和12的功能富集分析

Fig. 6 Enrichment analysis of region 1 and region 12 in breast cancer

在肿瘤微环境的空间异质性研究中,基因表达的空间分布模式是解析健康-肿瘤组织边界、肿瘤侵袭路径及免疫细胞募集机制的关键线索。为深度阐释乳腺癌组织的区域特异性特征并验证空间域识别效能,本研究对4种典型形态区域:浸润性导管癌、健康组织、导管/小叶原位癌及肿瘤边缘区的标记基因表达模式进行空间可视化(图7),结合生物学功能分析揭示其临床意义。标记基因CXCL14^[30-32]在区域1(浸润性导管癌)显著富集,作为趋化因子通常与炎症反应和免疫细胞的招募密切相关。CXCL14在乳腺癌组织中极有可能参与肿瘤微环境中免疫细胞的调节过程,进而影响肿瘤的生长与转移,与肿瘤的侵袭性以及免疫逃逸机制存在关联。标记基因AC087379.2^[33-35]在区域10(导管/小叶原位癌)特异性高表达,AC087379.2是一个非编码RNA,其可能涉及基因表达调控、细胞增殖或分化,通过影响关键的信号通路来促进肿瘤的发展,与原位癌向浸润癌转化的关键分子事件相关。标记基因APOE^[36]在区域15(肿瘤边缘区)中富集,作为脂质代谢枢纽分子,APOE参与脂质代谢和运输过程。在乳腺癌中,其高表达可能支持肿瘤细胞的膜脂合成与能量代谢重构,为侵袭性细胞提供物质基础。标记基因MALAT1^[37-38]显著富集在区域9(健康组织),这是一种长链非编码RNA,在多种细胞过程,如细胞增殖、迁移和分化中发挥作用。该区域的高表达模式提示MALAT1的生理功能可能与健康组织的结构维持及肿瘤发生的早期防御机制相关。

2.3 SpaMGCN精准揭示了小鼠冠状大脑的细微组织结构

为验证SpaMGCN在无人工注释的复杂数据集上的有效性,本研究将其应用于包含细微解剖结构的小鼠冠状脑组织10x Visium数据集。小鼠冠状大脑数据集是研究大脑结构和功能的重要资源,覆盖多个解剖区域,例如皮层和海马体等,它为全面理解整个大脑结构提供了支持。该数据集虽缺乏基准真相,但成年小鼠冠状大脑结构注释文献^[39]为空间域识别的准确性验证提供了重要参考。定量分析显示,SpaMGCN识别出的空间域与标注结果高度吻合,在空间平滑度、区域准确性以及完整性等关键指标如轮廓系数(Silhouette coefficient, SC)和DB指数(Davies-Bouldin index, DB)中均显著优于对比方法(图8)。这一结果表明,SpaMGCN通过聚类分配策略,在特征空间中实现了不同区域的有效分离。

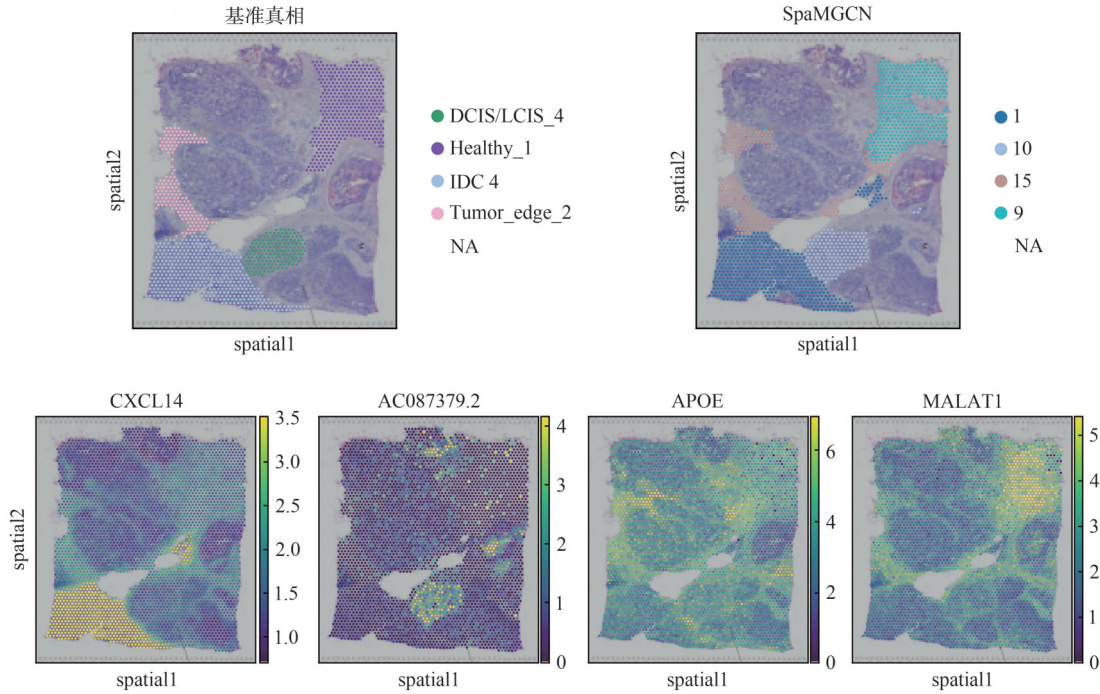


图 7 乳腺癌标记基因的空间分布

Fig. 7 Spatial distribution of breast cancer signature genes

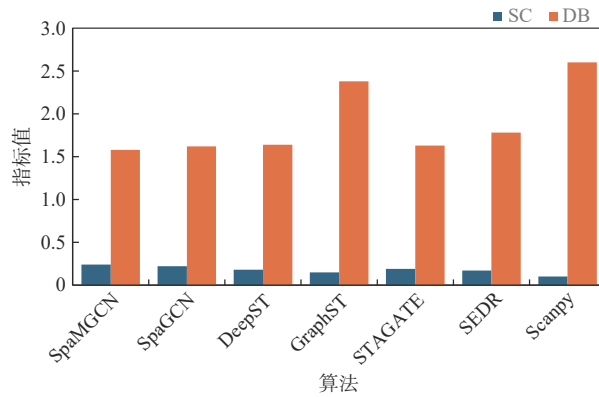


图 8 SpaMGCN 在小鼠冠状大脑数据上聚类指标的比较

Fig. 8 Comparison of SpaMGCN's clustering metrics on mouse coronal brain dataset

基于 Allen 参考图谱进一步表明, SpaMGCN 能够精准识别注释中的精细空间结构(图 9)。以小鼠海马体这一核心脑区为例, 其特征性的绳索状结构对应阿蒙角锥体层, 可细分为 CA1/CA2、CA3 区域, 箭头状结构则代表齿状回(DG)层。对比实验显示: DeepST、GraphST 以及 SpaGCN 皆未能准确识别连续的海马体组织, 导致 DG 区域与 CA3 区域分离且漏检 CA1 区域; STAGATE、SEDR 和 Scanpy 虽能识别海马体轮廓, 但 Scanpy 识别的 CA1 区域存在范围狭小、平滑度不足及点级噪声问题, STAGATE 与 SEDR 界定的 CA3 区域则存在过度扩展现象。相较之下, SpaMGCN 清晰表征了海马体的精细组织结构, 准确划分出 DG(区域 15)、CA3(区域 6)和 CA1(区域 12)等结构。值得注意的是, 通过与小鼠脑图谱的比对分析, 仅 SpaMGCN 成功识别出功能复杂的下丘脑后核区域(区域 5), 该区域在体温调节和能量平衡等生理过程中发挥关键作用。

综上所述, SpaMGCN 通过高效整合基因共表达依赖性和空间相邻依赖性信息, 在复杂 ST 数据的识别精度和细微生物结构解析能力上均达到较高水平, 展现出强大的方法学优势和临床转化应用潜力。

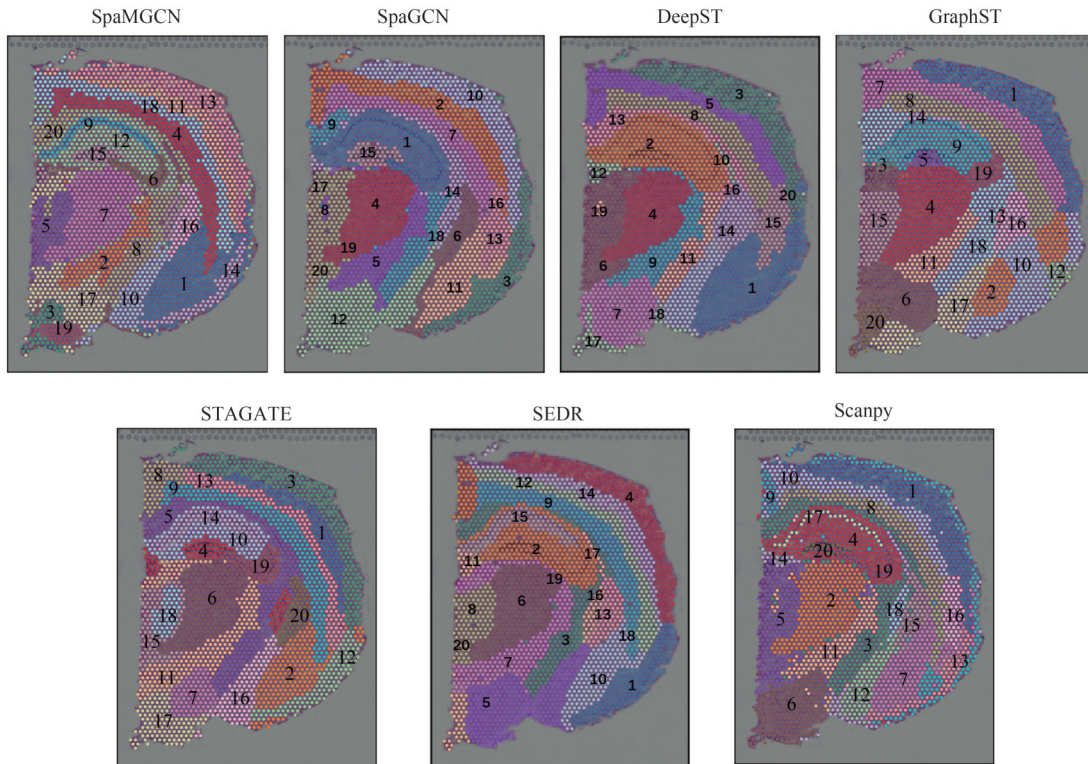


图9 SpaMGCN在小鼠冠状大脑数据上的结果与比较

Fig. 9 Clustering results and comparison of SpaMGCN on mouse coronal brain dataset

2.4 参数优化及消融实验

SpaMGCN模型在特征提取阶段所采用的策略是基于多视图的特征融合,通过对空间转录组学数据的基因表达信息以及空间位置信息来构建多视图模型,利用共享的编码器对斑点信息进行学习。并且为了实现空间特征和基因表达特征的有效整合,采用线性相加的方式,得到最终的嵌入特征。其中涉及表示基因特征权重的超参数 α ,我们在DLPFC数据集的多个切片上进行验证,可知当 α 取值范围在0.2左右时,所得的最终嵌入在进行下游任务时能取得相对稳定且良好的效果(图10),这可能是由于相较于基因特征,空间特征在执行空间域识别下游任务时拥有更加重要的作用。

为验证 SpaMGCN 算法模型各模块的作用机制,我们在 DLPFC 数据集的 #151673 切片上进行了消融实验。该实验系统地去除了深度嵌入聚类模块、空间视图、特征视图和细化聚类模块,以评估它们各自对模型性能的贡献(图11)。w/o-DEC表示无深度嵌入聚类模块;w/o-S表示无空间视图;w/o-F表示无特征视图;w/o-X表示无细化聚类模块。SpaMGCN

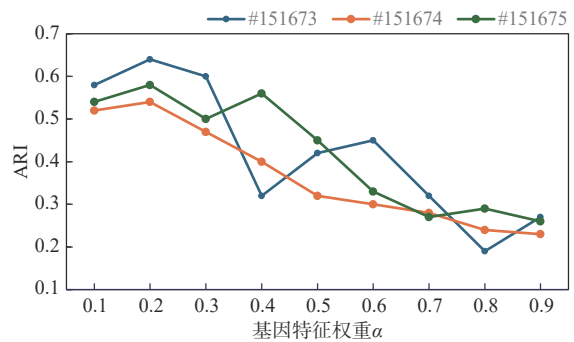


图10 基因特征权重 α 不同取值下模型效果比较

Fig. 10 Comparison of model performance under different values of gene feature weight α

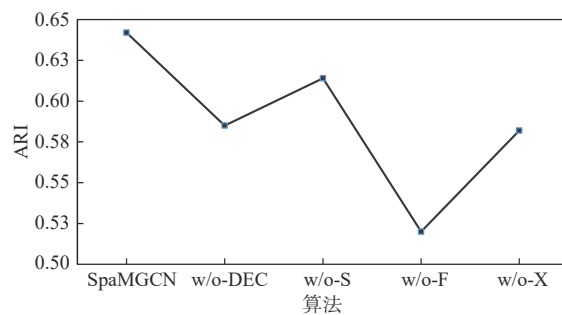


图11 #151673切片的消融实验

Fig. 11 Ablation experiment of the slice #151673

SpaMGCN

的ARI值明显优于其他变体(图11),表明这些模块的集成有利于空间域的识别。

此外,在ST数据分析中,初始聚类结果可能会受到噪声或局部不一致性的影响,导致某些区域的空间域识别不够准确。因此,SpaMGCN在算法中加入了细化聚类模块,通过对初始聚类结果进行后处理,确保每个点的类别标签与其周围大多数邻居的标签一致,从而提高聚类的准确性和稳定性。在12个切片的消融实验中(图12),有细化聚类模块的SpaMGCN算法的ARI值,显著高于无细化聚类的值。

3 结论

本研究提出一种基于多视图图卷积网络的空间域识别方法SpaMGCN,该方法利用多视图图卷积模块,从拓扑邻接图中提取空间位置相关的全局结构特征,并且从特征邻接图中捕获基因表达驱动的局部功能特征,通过线性融合策略实现两类特征的有机整合,使模型能够全面捕捉数据中隐含的空间依赖性与生物学功能关联性。此外,方法引入深度嵌入聚类模块对融合后的特征表示进行微调,有效捕捉基因表达模式与空间组织形态之间的复杂关系,进一步优化空间域聚类的性能。为了验证这些模块的有效性,本研究进行了消融实验,证实了各功能模块对SpaMGCN聚类性能的增强效果。

在多个公开数据集上,SpaMGCN与其他方法进行了系统性对比实验。结果表明,SpaMGCN在空间域划分的多项评估指标上均显著优于对比方法。尤其在处理边界清晰的DLPFC组织样本和细胞类型高度混杂的复杂数据集(如人类乳腺癌组织)时,SpaMGCN识别的空间域与已知组织学注释的一致性均达到更高水平,展现出对不同数据特征的强适应性。对比分析进一步证实,无论组织样本的细胞异质性高低或目标区域的尺度大小,SpaMGCN均能稳定地识别具有生物学意义的空间功能单元,展现出鲁棒的复杂结构解析能力,为解析ST数据复杂组织的空间域,揭示分子调控网络的空间异质性提供了高效的计算工具。当前的研究聚焦于对ST数据空间信息以及位点分子信息进行特征提取,而浪费了ST数据所自带的组织影像学信息,未来我们将针对ST数据的组织图像信息进行进一步分析,用以辅助增强模型对于数据的分析能力。

参考文献:

- [1] BLONDEL D V, GUILLAUME J, LAMBIOTTE R, et al. Fast unfolding of communities in large networks[J]. Journal of Statistical Mechanics-Theory and Experiment, 2008, 2008(10):476408.
- [2] 冯振兴, 尚文婧, 司佳宝, 等. 基于空间转录组学数据的空间域识别算法综述[J]. 内蒙古大学学报(自然科学版), 2024, 55(6):652-662.
- [3] KIPF T N, WELING M. Semi-supervised classification with graph convolutional networks[J]. Transactions on Machine Learning Research, 2016, 41(10):2577-2591.
- [4] HU J, LI X J, COLEMAN K, et al. SpaGCN: Integrating gene expression, spatial location and histology to identify spatial domains and spatially variable genes by graph convolutional network[J]. Nature Methods, 2021, 18(11):1342-1351.
- [5] ZONG Y, YU T, WANG X, et al. conST: An interpretable multi-modal contrastive learning framework for spatial

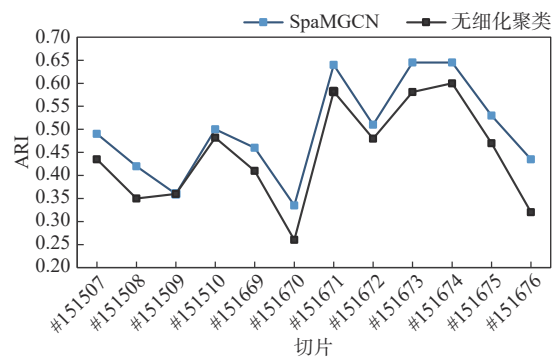


图12 有无细化聚类模块算法的比较

Fig. 12 Comparison of the algorithms with and without the refined clustering module

- transcriptomics[J]. *bioRxiv*, 2022:476408. <https://doi.org/10.1101/2022.01.14.476408>.
- [6] DONG K N, ZHANG S H. Deciphering spatial domains from spatially resolved transcriptomics with an adaptive graph attention auto-encoder[J]. *Nature Communications*, 2022, 13(1):1739.
- [7] PHAM D, TAN X, BALDERSON B, et al. Robust mapping of spatiotemporal trajectories and cell-cell interactions in healthy and diseased tissues[J]. *Nature Communications*, 2023, 14(1):7739.
- [8] XU C, JIN X Y, WEI S R, et al. DeepST: Identifying spatial domains in spatial transcriptomics by deep learning[J]. *Nucleic Acids Research*, 2022, 50(22):e131.
- [9] LONG Y H, ANG K S, LI M W, et al. Spatially informed clustering, integration, and deconvolution of spatial transcriptomics with GraphST[J]. *Nature Communications*, 2023, 14(1):1155.
- [10] XU H, FU H Z, LONG Y H, et al. Unsupervised spatially embedded deep representation of spatial transcriptomics[J]. *Genome Medicine*, 2024, 16(1):12.
- [11] ZHAO Y N, LONG C S, SHANG W J, et al. A composite scaling network of EfficientNet for improving spatial domain identification performance[J]. *Communications Biology*, 2024, 7(1):1567.
- [12] LEZAMA J, QIANG Q, PABLO M, et al. OLE: Orthogonal low-rank embedding, a plug and play geometric loss for deep learning[C]//2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition. Salt Lake City, UT, USA: IEEE, 2018:8109-8118.
- [13] XIE J, ROSS G, ALI F. Unsupervised deep embedding for clustering analysis[C]//Proceedings of the 33rd International Conference on Machine Learning. New York: JMLR, 2016:478-487.
- [14] MACQUEEN J. Some methods for classification and analysis of multivariate observations[J]. *Computing Research Repository*, 1967, 1:281-297.
- [15] SELVI A, KREACIC E, GHASSEMI M, et al. Distributionally and adversarially robust logistic regression via intersecting wasserstein balls[C]//Proceedings of the 29th International Conference on Neural Information Processing Systems. Cambridge: UAI, 2024:1576-1584.
- [16] MAYNARD K R, COLLADO-TORRES L, WEBER L M, et al. Transcriptome-scale spatial gene expression in the human dorsolateral prefrontal cortex[J]. *Nature Neuroscience*, 2021, 24(3):425-436.
- [17] WOLF F A, ANGERER P, THEIS F J. SCANPY: Large-scale single-cell gene expression data analysis[J]. *Genome Biology*, 2018, 19(1):15.
- [18] HAO Y H, HAO S, ANDERSEN-NISSEN E, et al. Integrated analysis of multimodal single-cell data[J]. *Cell*, 2021, 184(13):3573-3587.
- [19] STEINLEY D. Properties of the hubert-arabic adjusted rand index[J]. *Psychological Methods*, 2004, 9(3):386-396.
- [20] BECHT E, MCINNES L, HEALY J, et al. Dimensionality reduction for visualizing single-cell data using UMAP[J]. *Nature Biotechnology*, 2019, 37(1):38-44.
- [21] REN H L, WALKER B L, CANG Z X, et al. Identifying multicellular spatiotemporal organization of cells with SpaceFlow[J]. *Nature Communications*, 2022, 13(1):4076.
- [22] WAN L, FU Z, SUN L, et al. Self-supervised teaching and learning of representations on graphs[C]//Proceedings of the ACM Web Conference 2023. New York, NY, USA: Association for Computing Machinery, 2023:489-498.
- [23] POLYAK K I. Heterogeneity in breast cancer[J]. *J Clin Invest*, 2011, 121(10):3786-3788.
- [24] YEAP B H, Muniandy S, Lee S K, et al. Specimen shrinkage and its influence on margin assessment in breast cancer[J]. *Asian Journal of Surgery*, 2007, 30(3):183-187.
- [25] SIMONE R, BAILEY J, VINH N X, et al. Standardized mutual information for clustering comparisons: One step further in adjustment for chance[C]//Proceedings of the 31st International Conference on Machine Learning. Beijing, China: JMLR, 2014:1143-1151.
- [26] ESTÉVEZ P A, TESMER M, PEREZ C A, et al. Normalized mutual information feature selection[J]. *IEEE Transactions on Neural Networks*, 2009, 20(2):189-201.
- [27] FOWLKES E B, MALLOWS C L. A method for comparing two hierarchical clusterings[J]. *Journal of the American Statistical Association*, 1983, 78(383):553-569.
- [28] PALLA G, SPITZER H, KLEIN M, et al. Squidpy: A scalable framework for spatial omics analysis[J]. *Nature*

- Methods, 2022, 19(2): 171-178.
- [29] BABA S A, LABHSETWAR S, KLEMKE R, et al. A dataset of chromosomal instability gene signature scores in normal and cancer cells from the human breast[J]. Data in Brief, 2023, 51: 109647.
- [30] WALDEMER-STREYER R J, REYES-ORDOÑEZ A, KIM D, et al. Cxcl14 depletion accelerates skeletal myogenesis by promoting cell cycle withdrawal[J]. NPJ Regenerative Medicine, 2017, 2: 16017.
- [31] TANG L, CHEN X, HOU J A, et al. CXCL14 in prostate cancer: Complex interactions in the tumor microenvironment and future prospects[J]. Journal of Translational Medicine, 2025, 23(1): 9.
- [32] PARIKH A, SHIN J H, FAQUIN W, et al. Malignant cell-specific CXCL14 promotes tumor lymphocyte infiltration in oral cavity squamous cell carcinoma[J]. Journal for Immunotherapy of Cancer, 2020, 8(2): e001048.
- [33] KIM J, PIAO H L, KIM B J, et al. Long noncoding RNA MALAT1 suppresses breast cancer metastasis[J]. Nature Genetics, 2018, 50(12): 1705-1715.
- [34] YUE X, WU W Y, DONG M, et al. LncRNA MALAT1 promotes breast cancer progression and doxorubicin resistance via regulating miR-570-3p[J]. Biomedical Journal, 2021, 44(6): S296-S304.
- [35] HUANG X J, XIA Y, HE G F, et al. MALAT1 promotes angiogenesis of breast cancer[J]. Oncology Reports, 2018, 40(5): 2683-2689.
- [36] TAVAZOIE M F, POLLACK I, TANQUECO R, et al. LXR/ApoE activation restricts innate immune suppression in cancer[J]. Cell, 2018, 172(4): 825-840.
- [37] GOYAL B, YADAV SRM, AWASTHEE N, et al. Diagnostic, prognostic, and therapeutic significance of long non-coding RNA MALAT1 in cancer[J]. Biochim Biophys Acta Rev Cancer, 2021, 1875(2): 188502.
- [38] TUFAIL M. The MALAT1-breast cancer interplay: Insights and implications[J]. Expert Review of Molecular Diagnostics, 2023, 23(8): 665-678.
- [39] SUNKIN S M, NG L, LAU C, et al. Allen brain Atlas: An integrated spatio-temporal portal for exploring the central nervous system[J]. Nucleic Acids Research, 2013, 41(Database issue): D996-D1008.

(责任编辑 那顺布和)

SpaMGCN: A Method for Identifying Spatial Domains Based on Multi-View Graph Neural Networks

LIU Hexin¹, SHANG Wenjing¹, ZHAO Xiangyu¹, ZHENG Yifan¹,
ZHANG Jia², FENG Zhenxing¹

(1. College of Science, Inner Mongolia University of Technology, Hohhot 010051, China;

2. College of Life Sciences, Inner Mongolia University, Hohhot 010021, China)

Abstract: Spatial domain identification, a critical task in spatial transcriptomics, aims to accurately delineate tissue regions by integrating gene expression profiles with spatial information. However, existing methods often fall short in capturing both local neighborhood and global structural features. To overcome this limitation, we present SpaMGCN, a novel multi-view graph convolutional network approach that leverages a dual-modality architecture to concurrently extract local and global features from gene expression correlations and spatial proximity. Evaluated across multiple public datasets, SpaMGCN achieves fine-grained boundary delineation of complex tissue architectures, demonstrating superior performance in spatial domain identification. This method offers an effective computational framework for elucidating the spatial heterogeneity of biological tissues.

Key words: spatial transcriptomics; identifying spatial domain; multi-view graph neural network