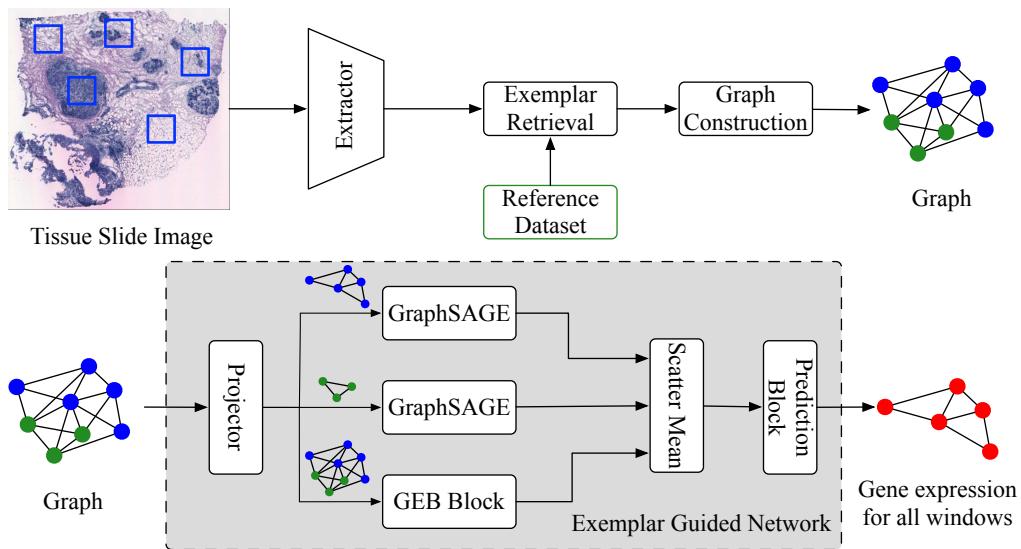

Graphical Abstract

Spatial Transcriptomics Analysis of Gene Expression Prediction using Exemplar Guided Graph Neural Network

Yan Yang, Md Zakir Hossain, Eric Stone, Shafin Rahman



In this paper, we aim to predict gene expression of each **window** in a tissue slide image. Given a tissue slide image, we encode the **windows** to feature space, retrieve its **exemplars** from the reference dataset, construct a graph, and then dynamically predict **gene expression** of each window with our exemplar guided graph network.

Email addresses: u6169130@anu.edu.au (Yan Yang), zakir.hossain@anu.edu.au (Md Zakir Hossain), eric.stone@anu.edu.au (Eric Stone), shafin.rahman@northsouth.edu (Shafin Rahman)

Highlights

Spatial Transcriptomics Analysis of Gene Expression Prediction using Exemplar Guided Graph Neural Network

Yan Yang, Md Zakir Hossain, Eric Stone, Shafin Rahman

- We propose an exemplar guided graph network to accurately predict gene expression from a slide image window.
- We design a graph construction strategy to connect windows and exemplars for performing exemplar learning of gene expression prediction.
- We propose a graph exemplar bridging block to revise the window feature by using its nearest exemplars.
- Experiments on two standard benchmark datasets demonstrate our superiority when compared with state-of-the-art approaches.

Spatial Transcriptomics Analysis of Gene Expression Prediction using Exemplar Guided Graph Neural Network

Yan Yang^a, Md Zakir Hossain^{a,b}, Eric Stone^{a,*}, Shafin Rahman^c

^a*Biological Data Science Institute, The Australian National University, Canberra, Australia*

^b*Optus Centre for AI, Curtin University, Perth, Australia*

^c*Department of Electrical and Computer Engineering, North South University, Bangladesh*

Abstract

Spatial transcriptomics (ST) is essential for understanding diseases and developing novel treatments. It measures the gene expression of each fine-grained area (i.e., different windows) in the tissue slide with low throughput. This paper proposes an exemplar guided graph network dubbed EGNN to accurately and efficiently predict gene expression from each window of a tissue slide image. We apply exemplar learning to dynamically boost gene expression prediction from nearest/similar exemplars of a given tissue slide image window. Our framework has three main components connected in a sequence: i) an extractor to structure a feature space for exemplar retrievals; ii) a graph construction strategy to connect windows and exemplars as a graph; iii) a graph convolutional network backbone to process window and exemplar fea-

*Corresponding author.

Email addresses: u6169130@anu.edu.au (Yan Yang), zakir.hossain@anu.edu.au (Md Zakir Hossain), eric.stone@anu.edu.au (Eric Stone), shafin.rahman@northsouth.edu (Shafin Rahman)

tures, and a graph exemplar bridging block to adaptively revise the window features using its exemplars. Finally, we complete the gene expression prediction task with a simple attention-based prediction block. Experiments on standard benchmark datasets indicate the superiority of our approach when compared with past state-of-the-art methods.

Keywords: Spatial transcriptomics, Gene expression prediction, Deep learning, Graph convolution, Tissue slide image

1. Introduction

Based on an editorial report of the *Natural Methods* [1], spatial transcriptomics (ST) is the future of studying disease because of its capabilities in measuring gene expression of fine-grained areas (i.e., different windows) of tissue slides. However, ST is in low throughput due to limitations in current analysis for the candidate windows [2]. To accurately predict gene expression from each window of a tissue slide image (Fig. 1), this paper proposes a solution named exemplar guided graph network (EGGN), allowing efficient and concurrent analysis.

Previous works adopt end-to-end neural networks, namely STNet [3] and NSL [4], to independently establish a mapping between gene expression and the slide image window. STNet is a transfer learning-based approach that finetunes a pretrained DenseNet [5] for the gene expression prediction task. On the contrary, NSL maps the color intensity of the slide image window to gene expression by a single convolution operation. Though amenable to high throughput because of using neural networks, their prediction performance is inferior.

We investigate two important limitations of the existing approaches [3, 4].

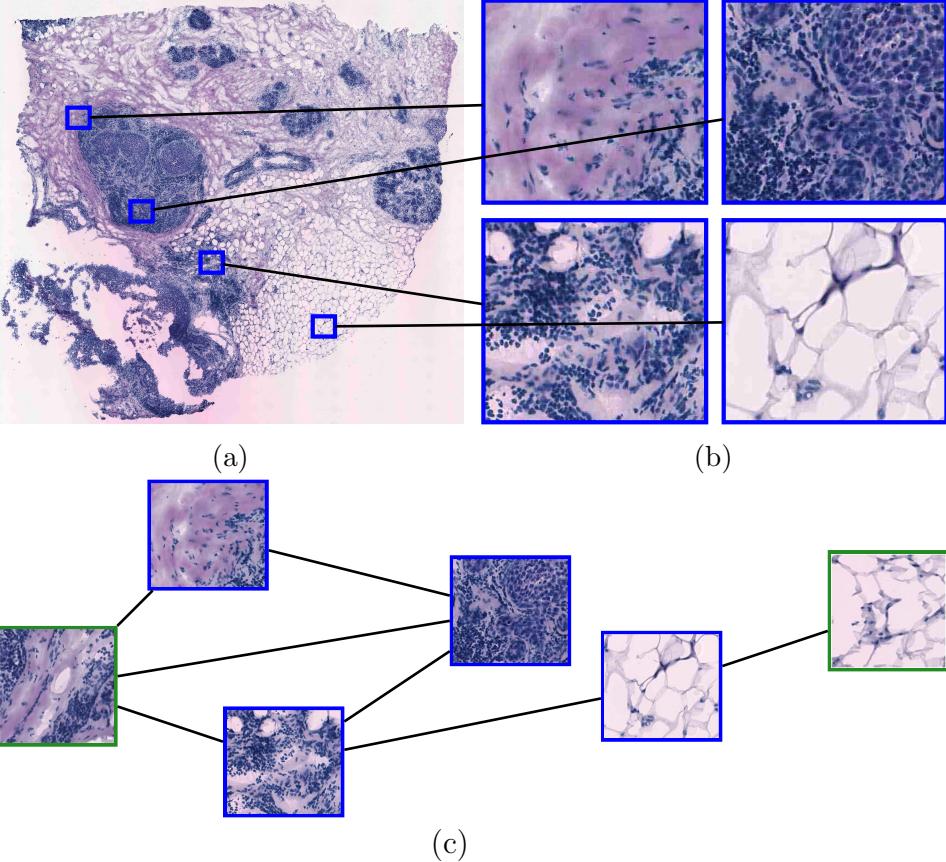


Figure 1: Overview of fields. Each fine-grained area (i.e., window) of (a) a tissue slide image has a distinct expression of the same gene types. For example, we have a tissue slide image with (b) four windows and each of the windows corresponds with the different expression of the same gene types. Our goal is to predict the gene expression of each window. Our key idea is to build (c) a graph connecting **windows** spatially close with each other and the **exemplars** for each window, which is used by our proposed framework to dynamically benefit gene expression prediction.

19 i) Local feature aggregation: gene expression prediction can be considered as
 20 individually aggregating and identifying the feature of each gene type for the
 21 slide image window. The long-range dependency, i.e., global context, among
 22 identified features is needed to reason about complex scenarios [6, 7], as those
 23 features are generally non-uniformly distributed across the slide image (see

24 Sec. 3 for details). STNet (i.e., a pure convolution approach) emphasizes
25 local context during feature aggregation within a limited slide image win-
26 dow. It fails to bring interaction among features that are apart from each
27 other. By experimenting with extensive state-of-the-art (SOTA) network ar-
28 chitectures, we show that models with local feature aggregation achieve low
29 performance when compared to models with long-range dependencies on the
30 slide image. ii) Vulnerable assumption: identifying gene expression by di-
31 rectly detecting the color intensity of the slide image window is vulnerable.
32 By experimenting on standard benchmark datasets, we show that NSL only
33 works in extreme cases. For example, in the STNet dataset [3], tumor ar-
34 eas of slide image windows are usually purple, which benefits tumor-related
35 gene expression prediction. This method finds a negative Pearson correla-
36 tion coefficient (PCC) when evaluating the model with the least reliable gene
37 expression prediction, i.e., PCC@F in Tab. 1.

38 Furthermore, as shown by [3], windows of a slide image spatially close
39 with each other usually exhibit similar features, suggesting correlated gene
40 expression among them. However, both of the above two approaches in-
41 dependently predict the gene expression of each window and ignore inter-
42 actions among these windows. Windows distributed over the slide image
43 can be connected to construct a graph for gene expression prediction with a
44 graph convolutional network (GCN) [8]. Compared to independently mod-
45 eling the window, graph-based modeling of the task allows reasoning depen-
46 dency among windows for predicting gene expression, considering the local
47 structure/context/spatial relations among windows.

48 In this paper, we propose an EGNN framework to address the above lim-

49 itations. EGNN uses GraphSAGE [9] as a GCN backbone and incorporates
50 exemplar learning concepts for the gene expression prediction task. To enable
51 exemplar retrieval of a given slide image window, we use a feature extractor
52 for defining a feature space to calculate the similarity between two slide im-
53 age windows. Then, we construct a graph, considering the local context (i.e.,
54 nearby windows) and global context (i.e., shared exemplars act as anchors
55 for information propagation) for the gene expression prediction task. We use
56 GraphSAGE layers as our backbone to allow interactions among windows
57 and exemplars. Meanwhile, we have a graph exemplar bridging (GEB) block
58 to update window features by the exemplars and the gene expression of ex-
59 emplars. Allowing dynamic information propagation, the exemplar feature
60 also receives and is updated with the status of the window features. Se-
61 mantically, the former update corresponds with ‘the known gene expression’,
62 and the latter corresponds with ‘the gene expression the model wants to be
63 known’. Finally, we have an attention-based prediction block to aggregate
64 exemplars of each window and the exemplar-revised window features, for
65 predicting gene expression.

66 Our contributions are summarised as follows:

- 67 • We propose an EGNN framework, a GCN-based exemplar learning
68 approach, to accurately predict gene expression from the slide image
69 window;
- 70 • We design a graph construction strategy to connect windows and ex-
71 emplars for performing exemplar learning of gene expression prediction
72 under GCNs;
- 73 • We propose a GEB block to revise the window feature by using its

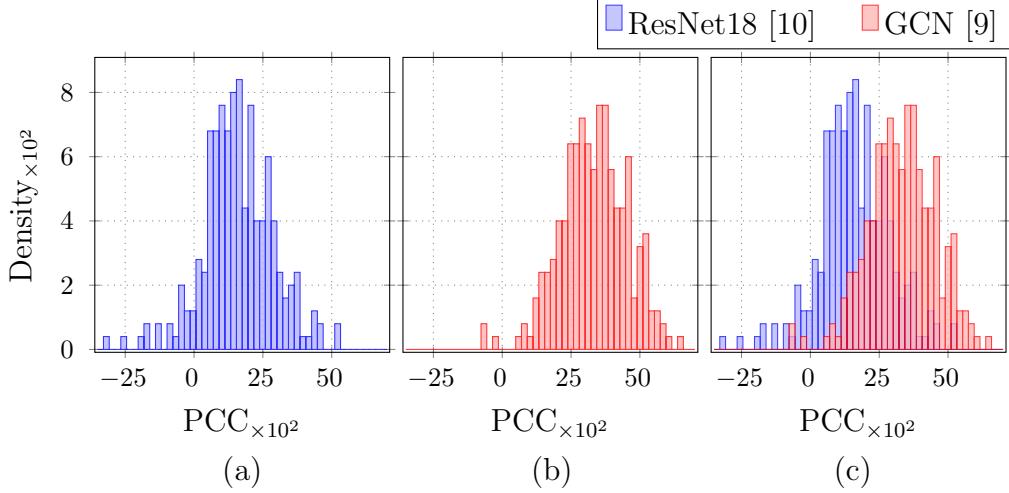


Figure 2: We compare the PCC distribution obtained by using (a) ResNet18 model and (b) GCN model for gene expression prediction, where window features used in the GCN [9] are extracted by the ImageNet-1K pre-trained ResNet18 [10]. The x-axis and y-axis respectively denote the PCC and the density at different PCC. We also combine the plot (a) and plot (b) into the plot (c), for a better comparison. As shown, modeling the spatial relations among windows out-weights feature learning of each window.

74 nearest exemplars;

75 • Experiments on two standard benchmark datasets demonstrate our su-
76 periority when compared with SOTA approaches.

77 A preliminary version of this paper has been published previously [11].

78 However, it considers interactions between exemplars and each window in-
79 dependently, and ignores the exemplar shared by multiple different windows
80 potentially serve as anchors for facilitating information propagation among
81 spatial apart windows of a slide image. In this paper, we extend our previous
82 version as follows: i) we demonstrate that reasoning spatial relations among
83 windows out-weights feature learning in our problem, which motivates us to
84 design a GCN-based exemplar learning framework (Fig. 2); ii) we address
85 the above independent interaction and information propagation bottleneck

86 issues in our new exemplar learning framework; iii) we analyze the proposed
87 framework with careful ablation designs, *e.g.*, we show that the proposed
88 exemplar learning strategy benefits general GCN frameworks; iv) we explore
89 SOTA GCN-based methods in extensive experimental comparisons.

90 **2. Related Work**

91 This section first reviews the study of gene expression prediction. Then,
92 we summarise recent exemplar learning achievements in natural language
93 processing and computer vision domains. Finally, we briefly introduce the
94 background of GCNs.

95 *Gene Expression Prediction.* Measuring gene expression is a fundamental
96 process in developing novel treatments and monitoring human diseases [12].
97 Recently, deep learning methods have been introduced to this task. Existing
98 methods predict the gene expression either from DNA sequences [12] or slide
99 images [3, 4, 13]. This paper explores the latter approach which is divided
100 into two streams. First, Schmauch *et al.*[13] employs a multi-stage method in-
101 cluding pretrained ResNet feature extractions [10] and a K-Means algorithm
102 to model the Bulk RNA-Seq technique [14]. It measures the gene expression
103 across cells within a large predefined area that is up to $10^5 \times 10^5$ pixels in
104 a corresponding slide image [13]. However, this approach is ineffective for
105 studies that require fine-grained gene expression information, such as tumor
106 heterogeneity [15]. Second, on the contrary, He *et al.* [13, 5] and Dawood *et*
107 *al.* [4] design a STNet and an NSL framework to predict the gene expression
108 for each window (*i.e.*, fine-grained area) of a slide image. This corresponds
109 with the ST technique [16]. We model ST to predict gene expression, as this

110 potentially solves the bulk RNA-Seq prediction task simultaneously [3]. For
111 example, aggregation of gene expression predictions for each window across
112 a slide image results in a bulk RNA-Seq prediction.

113 *Exemplar Learning.* The K-nearest neighborhood classifier is the most straight-
114 forward case of exemplar learning. It classifies the input by considering the
115 class labels of the nearest neighbors. Exemplar learning is a composition
116 of retrieval tasks and learning tasks [17]. It has been widely employed to
117 increase the model capability by bringing in extra knowledge of similar ex-
118 emplars. The applications of exemplar learning include visual question an-
119 swering [18, 19, 20, 21], language generation [22, 23, 24, 25], real-time visual
120 object tracking [24], fact-checking [26], fact completion [27, 28], and dialogue
121 [29]. However, most of the exemplar learning approaches do not apply to our
122 task because of the domain shift. This paper investigates an application of
123 exemplar learning in gene expression prediction from slide image windows.
124 As a result, we devise a GEB block to adapt exemplar learning to our task.

125 *Graph Convolutional Network.* GCNs are an extension of Convolutional Neu-
126 ral Networks (CNNs) designed for modeling non-Euclidean structured data
127 such as graphs and reasoning the underlying complex relationships [8]. GCNs
128 can be categorized into two main groups: spectral-based and spatial-based
129 approaches. Spectral-based GCNs rely on the graph Laplacian and perform
130 convolutions of spatial domain in the frequency domain [30, 31], while spatial-
131 based GCNs utilize message-passing frameworks to progressively propagate
132 information from source nodes to their neighbors (target nodes) along the
133 edges of the graph. After information propagation, target nodes aggregate
134 all received information, potentially bringing the perception of the graph

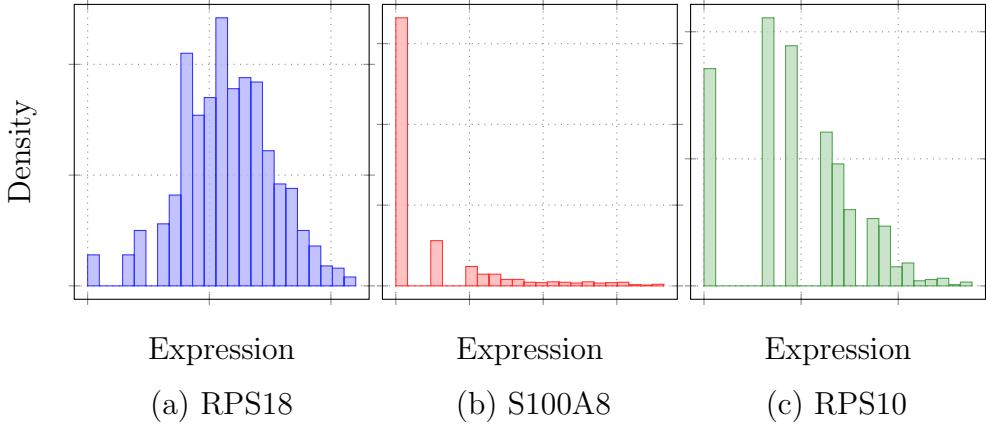


Figure 3: Gene expression distributions of STNet dataset [3]. Each gene expression is log-transformed. (a) is the well-distributed expression of gene RPS18. (b) and (c) are the long-tail distributed expression of gene S100A8 and gene RPS10.

structure into final node features [9, 32]. Some works also explore edge attributes during the propagation process [33, 34, 35]. GCNs with different information propagation and aggregation strategies have been proposed for diverse vision tasks [36, 37, 9, 32, 38]. In this paper, we propose a GEB to manage information propagation between windows and exemplars for the gene expression prediction task. The GEB block dynamically revises the window features using the nearest exemplars and updates the exemplars to improve the gene expression prediction.

3. EGNN Framework

Problem Formulation. We have a tissue slide image containing multiple windows; each window is annotated with gene expression. We denote the slide image as pairs of slide image windows \mathbf{w}_i and gene expression \mathbf{y}_i , i.e., $\{(\mathbf{w}_i, \mathbf{y}_i)\}_{i=1}^p$, where p is the set size. We aim to train a neural network model to predict \mathbf{y}_i from \mathbf{w}_i .

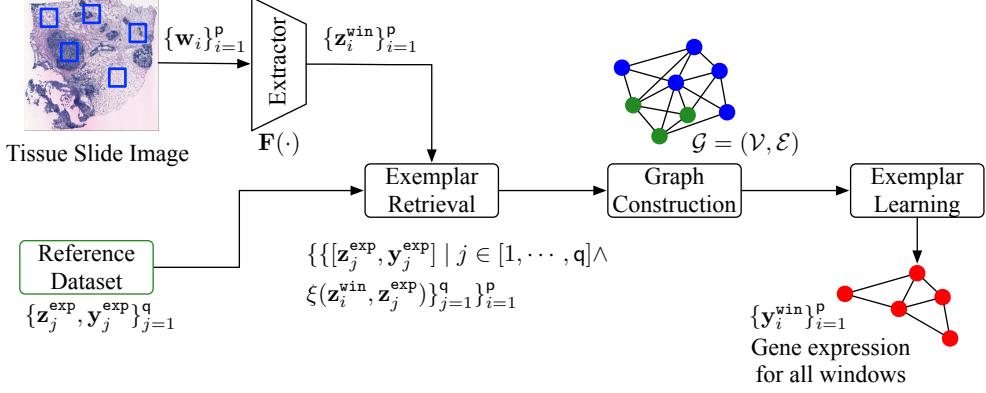


Figure 4: Framework overview. Given a slide images containing p windows \mathbf{w}_i , i.e., $\{\mathbf{w}_i\}_{i=1}^p$, we embed each window \mathbf{w}_i into features by using the feature extractor $\mathbf{F}(\cdot)$. We have $\mathbf{z}_i^{\text{win}} = \mathbf{F}(\mathbf{w}_i)$. Meanwhile, there is a reference database $\{\mathbf{z}_j^{\text{exp}}, \mathbf{y}_j^{\text{exp}}\}_{j=1}^q$ that collects q pairs of exemplar embedded by the same $\mathbf{F}(\cdot)$ and the gene expression of the exemplar. For each window \mathbf{w}_i , we then perform exemplar retrieval from the reference database, resulting in $\{[\mathbf{z}_j^{\text{exp}}, \mathbf{y}_j^{\text{exp}}] \mid j \in [1, \dots, q] \wedge \xi(\mathbf{z}_i^{\text{win}}, \mathbf{z}_j^{\text{exp}})\}$, where $\xi(\cdot, \cdot)$ determines if $\mathbf{z}_j^{\text{exp}}$ is the nearest exemplar of $\mathbf{z}_i^{\text{win}}$. We then construct a graph $\mathcal{G} = (\mathcal{V}, \mathcal{E})$. \mathcal{V} is a node set of **windows** and **exemplars**, and \mathcal{E} is edges that connect the nodes. We finally perform exemplar learning on the graph \mathcal{G} , obtaining **gene expression** for all window nodes, i.e., $\{\mathbf{y}_i^{\text{win}}\}_{i=1}^p$.

149 From the ST study [3], two main challenges exist. i) Long-range dependency:
 150 gene expression-related features are non-uniformly distributed over
 151 the slide image (refer to Figure 3 of [3] for evidence). The interactions
 152 among these features are needed to group expression of the same gene type;
 153 ii) Skewed gene expression distribution: the expression of some gene types
 154 has a skewed distribution, similar to the imbalance class distribution prob-
 155 lem. This skewed distribution (see Fig. 3 for an example) poses challenges
 156 in predicting the expression of these gene types. In this paper, we attempt
 157 to mitigate them by learning from similar exemplars.

158 *Model Overview.* With these motivations, we have designed the EGNN frame-
 159 work containing three main modules in sequence. The overview of our frame-

work is shown in Fig. 4. i) Exemplar retrieval (Sec. 3.1): we have a feature extractor $\mathbf{F}(\cdot)$ to embed the slide image window \mathbf{w}_i into feature $\mathbf{z}_i^{\text{win}}$, i.e., $\mathbf{z}_i^{\text{win}} = \mathbf{F}(\mathbf{w}_i)$. Meanwhile, there is a reference dataset $\{\mathbf{z}_j^{\text{exp}}, \mathbf{y}_j^{\text{exp}}\}_{j=1}^q$ containing q pairs of exemplars $\mathbf{z}_j^{\text{exp}}$ embedded by the feature extractor $\mathbf{F}(\cdot)$ and the exemplar gene expression $\mathbf{y}_j^{\text{exp}}$. In the feature space of $\mathbf{F}(\cdot)$, we then retrieve the nearest exemplars of \mathbf{w}_i from the reference dataset to form a set $\{[\mathbf{z}_j^{\text{exp}}, \mathbf{y}_j^{\text{exp}}] \mid j \in [1, \dots, q] \wedge \xi(\mathbf{z}_i^{\text{win}}, \mathbf{z}_j^{\text{exp}})\}$, where $\xi(\cdot, \cdot)$ is the kNN algorithm used to determine if $\mathbf{z}_j^{\text{exp}}$ is the nearest exemplar of $\mathbf{z}_i^{\text{win}}$. ii) Graph construction (Sec. 3.2): we construct an exemplar-based graph $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ for the slide image. Each window or each exemplar can be considered as a node, and they form two sets of nodes, a window node set \mathcal{V}^{win} and an exemplar node set \mathcal{V}^{exp} . The union of \mathcal{V}^{win} and \mathcal{V}^{exp} is \mathcal{V} , i.e., $\mathcal{V} = \mathcal{V}^{\text{win}} \cup \mathcal{V}^{\text{exp}}$. The edge set is denoted as $\mathcal{E} = \mathcal{E}_{\text{win} \rightarrow \text{win}} \cup \mathcal{E}_{\text{exp} \rightarrow \text{exp}} \cup \mathcal{E}_{\text{win} \rightarrow \text{exp}} \cup \mathcal{E}_{\text{exp} \rightarrow \text{win}}$, exploring relations of window to window $\mathcal{E}_{\text{win} \rightarrow \text{win}}$, exemplar to exemplar $\mathcal{E}_{\text{exp} \rightarrow \text{exp}}$, window to exemplar $\mathcal{E}_{\text{win} \rightarrow \text{exp}}$, and exemplar to window $\mathcal{E}_{\text{exp} \rightarrow \text{win}}$. In our formulation, \mathcal{G} is a heterogeneous graph. iii) Exemplar learning (Sec. 3.3): we train a model $\mathbf{C}(\cdot, \cdot, \cdot)$ that maps \mathcal{G} , $\{\mathbf{z}_i^{\text{win}}\}_{i=1}^p$, and $\{(\mathbf{z}_i^{\text{exp}}, \mathbf{y}^{\text{exp}_i})\}_{i=1}^q$ to $\{\mathbf{y}_i^{\text{win}}\}_{i=1}^p$ by a single forward pass. Our model uses a GraphSAGE-based backbone. We bring interactions between $\mathbf{e}_i^{\text{win}}$ and $\mathbf{e}_j^{\text{exp}}$ to leverage $\mathbf{y}_j^{\text{exp}}$ with a proposed GEB block, whenever $e_{ji} \in \mathcal{E}^{\text{exp} \rightarrow \text{win}}$. Meanwhile, for facilitating their interactions, the $\mathbf{e}_i^{\text{win}}$ is propagated back to revise $\mathbf{e}_j^{\text{exp}}$. With the introduction of exemplars, our framework dynamically benefits when predicting gene expression.

3.1. Exemplar Retrieval

To retrieve the exemplar of a given window \mathbf{w}_i , we have an extractor (i.e., an encoder) that is coupled with a distance metric to amount the similar-

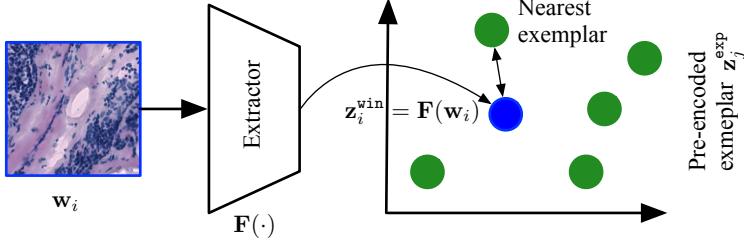


Figure 5: Overview of the exemplar retrieval. The green cycles are pre-embedded exemplar features, e.g., z_j^{exp} . Given a window w_i , we extract z_i^{win} with $F(\cdot)$ and retrieve the nearest exemplar.

185 ity between the given window and exemplar for a reference dataset for the
 186 exemplar retrieval.

187 *Feature Extractor.* Unlike our previous work [11] that trains a StyleGAN-
 188 based autoencoder [39, 40] for retrieving the exemplars, we found that an
 189 ImageNet-1K pre-trained ResNet brings a more accurate gene expression
 190 prediction in the graph-based setting. We denote $F(\cdot)$ as a ResNet feature
 191 extractor, where the classification layer is removed from a standard ResNet.
 192 Though $F(\cdot)$ is trained from a different domain, we show that $F(\cdot)$ effectively
 193 encodes textures for gene expression prediction, by comparing it with other
 194 possible encoders in a later section.

195 *Method.* We use the extractor $F(\cdot)$ for the exemplar retrieval. The window w_i
 196 is embedded into $z_i^{\text{win}} = F(w_i)$. In the feature space of $F(\cdot)$, we measure the
 197 similarity between z_i^{win} and each exemplar z_j^{exp} from the reference dataset by
 198 using the Manhattan distance \mathcal{L}_1 . We then execute the kNN algorithm $\xi(\cdot, \cdot)$
 199 for retrieving the nearest exemplar set. We empirically verify the optimal
 200 number of used exemplars in Sec. 4.2. To generalize the model performance,
 201 we restrict that the candidate image pairs are from different patients. The

202 overall process is presented in Fig. 5. In the experiment section, we compare
 203 the proposed exemplar retrieval with alternative retrieval approaches.

204 *3.2. Graph Construction*

205 We construct a graph $\mathcal{G} = (\mathcal{V}, \mathcal{E})$. Each window and exemplar is treated
 206 as a node, forming a window node set $\mathcal{V}^{\text{win}} = \{v_i^{\text{win}}\}_{i=1}^p$ and an exemplar node
 207 set $\mathcal{V}^{\text{exp}} = \{v_j^{\text{exp}}\}_{j=1}^q$. We take an union of the two set, $\mathcal{V} = \mathcal{V}^{\text{win}} \cup \mathcal{V}^{\text{exp}}$, to
 208 obtain a node set of our graph \mathcal{G} . As described, we consider four edge types.
 209 They are

$$\mathcal{E}_{\text{win} \rightarrow \text{win}} = \{e_{ij} \mid v_i^{\text{win}}, v_j^{\text{win}} \in \mathcal{V}^{\text{win}} \times \mathcal{V}^{\text{win}} \wedge \phi(v_i^{\text{win}}, v_j^{\text{win}})\} , \quad (1)$$

$$\mathcal{E}_{\text{exp} \rightarrow \text{exp}} = \{e_{ij} \mid v_i^{\text{exp}}, v_j^{\text{exp}} \in \mathcal{V}^{\text{exp}} \times \mathcal{V}^{\text{exp}} \wedge \xi(v_i^{\text{exp}}, v_j^{\text{exp}})\} , \quad (2)$$

$$\mathcal{E}_{\text{win} \rightarrow \text{exp}} = \{e_{ij} \mid v_i^{\text{win}}, v_j^{\text{exp}} \in \mathcal{V}^{\text{win}} \times \mathcal{V}^{\text{exp}} \wedge \xi(v_i^{\text{win}}, v_j^{\text{exp}})\} , \quad (3)$$

$$\mathcal{E}_{\text{exp} \rightarrow \text{win}} = \{e_{ij} \mid e_{ji} \in \mathcal{E}_{\text{win} \rightarrow \text{exp}}\} . \quad (4)$$

210 $\phi(\cdot, \cdot)$ explores the spatial relation of windows, determining if the spatial
 211 distance of two input windows in the slide is below a given threshold. Note
 212 that $\xi(\cdot, \cdot)$ is a kNN-based function in the feature space of $\mathcal{F}(\cdot)$ for connecting
 213 k-nearest neighbors. Thus, we have the edge set of the graph \mathcal{G} as $\mathcal{E} =$
 214 $\mathcal{E}_{\text{win} \rightarrow \text{win}} \cup \mathcal{E}_{\text{exp} \rightarrow \text{exp}} \cup \mathcal{E}_{\text{win} \rightarrow \text{exp}} \cup \mathcal{E}_{\text{exp} \rightarrow \text{win}}$.

215 *3.3. Exemplar Learning*

216 Our model $\mathbf{C}(\cdot, \cdot, \cdot)$ is composed of a projector, a L -layer GraphSAGE-
 217 based backbone, GEB blocks, and a prediction block, where we use the GEB
 218 block in each of the backbone layers. With graph \mathcal{G} , our model maps node
 219 features ($\{\mathbf{z}_i^{\text{win}}\}_{i=1}^p$ of a slide image and exemplars $\{\mathbf{z}_j^{\text{exp}}, \mathbf{y}_j^{\text{exp}}\}_{j=1}^q$) to gene
 220 expression $\{\mathbf{y}_i^{\text{win}}\}_{i=1}^p$ of each window. Our model allows interactions between

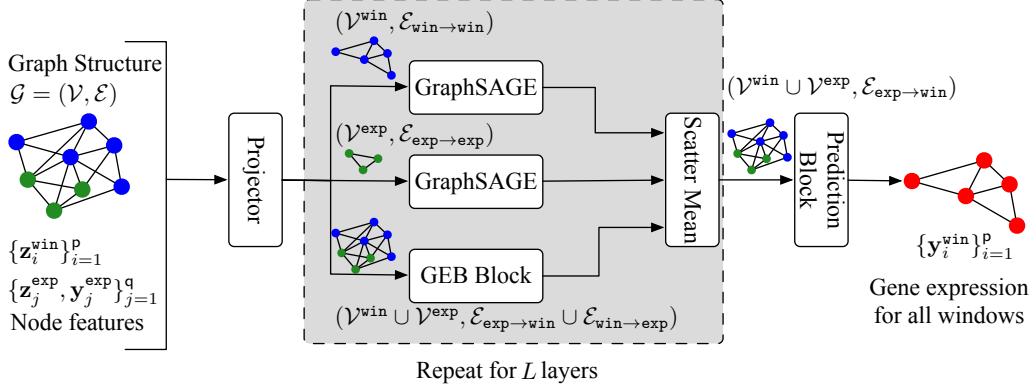


Figure 6: Architectures of our model. We have a graph $\mathcal{G} = (\mathcal{V}, \mathcal{E})$, where the node set $\mathcal{V} = \mathcal{V}^{\text{win}} \cup \mathcal{V}^{\text{exp}}$ is a union of **windows** and **exemplars**, and the edge set $\mathcal{E} = \mathcal{E}_{\text{win} \rightarrow \text{win}} \cup \mathcal{E}_{\text{exp} \rightarrow \text{exp}} \cup \mathcal{E}_{\text{win} \rightarrow \text{exp}} \cup \mathcal{E}_{\text{exp} \rightarrow \text{win}}$ contains four types of relations among nodes. We first refine features (i.e., \mathcal{G} , $\{\mathbf{z}_i^{\text{win}}\}_{i=1}^p$ and $\{\mathbf{z}_j^{\text{exp}}, \mathbf{y}_j^{\text{exp}}\}_{j=1}^q$) of each node by using a projector. We then respectively allow interactions among \mathcal{V}^{win} and \mathcal{V}^{exp} by two GraphSAGE block [9]. Meanwhile, we have a GEB block that revises the features of nodes contained in \mathcal{V}^{win} and \mathcal{V}^{exp} by considering edges between the two node sets. The scatter mean operation [41] average the features computed by the two GraphSAGE blocks and the GEB block for each node. Finally, after repeating the above operations for L layers, there is a prediction block that weights the contribution of each exemplar to a window for the gene expression prediction task. Finally, we have **gene expression** of all window nodes as $\{\mathbf{y}_i^{\text{win}}\}_{i=1}^p$. Note that the **blue window nodes** turn **red**, carrying gene expression prediction at each of the window nodes.

221 $\{\mathbf{z}_i^{\text{win}}\}_{i=1}^p$ and $\{\mathbf{z}_j^{\text{exp}}, \mathbf{y}_j^{\text{exp}}\}_{j=1}^q$ to progressively revise their intermediate fea-
 222 tures within the backbone for the prediction task. The overall architecture
 223 is shown in Fig. 6.

224 *Projector.* The features, $\mathbf{z}_i^{\text{win}}$ and $\mathbf{z}_j^{\text{exp}}$, are embedded with a wide range of
 225 dataset-dependent attributes. We refine them to concentrate on the gene
 226 expression of interest by several multi-layer perceptrons (MLPs). Firstly,
 227 each window $\mathbf{z}_i^{\text{win}}$ is projected by $\text{MLP}_h^0(\cdot)$. Secondly, for each $\mathbf{z}_j^{\text{exp}}$, as its
 228 associated gene expression $\mathbf{y}_j^{\text{exp}}$, is available, we empower the refinement of
 229 $\mathbf{z}_j^{\text{exp}}$ by $\mathbf{y}_j^{\text{exp}}$. We concatenate $\mathbf{z}_j^{\text{exp}}$ and $\mathbf{y}_j^{\text{exp}}$ before feeding to $\text{MLP}_r^0(\cdot)$. We

230 have

$$\begin{aligned}\mathbf{h}_i^0 &= \text{MLP}_h^0(\mathbf{z}_i^{\text{win}}) , & \forall v_i \in \mathcal{V}^{\text{win}} , \\ \mathbf{r}_j^0 &= \text{MLP}_r^0([\mathbf{z}_j^{\text{exp}} \| \mathbf{y}_j^{\text{exp}}]) , & \forall v_j \in \mathcal{V}^{\text{exp}} ,\end{aligned}$$

231 where the superscripts of \mathbf{h}_i^0 and \mathbf{r}_j^0 denote that they are initial refined feature,
232 and $\cdot \| \cdot$ is a concatenation operator. $\text{MLP}_h^0(\cdot)$ and $\text{MLP}_r^0(\cdot)$ are three-layer
233 perceptrons with LeakyReLU activation functions. The bottom two layers
234 of $\text{MLP}_h^0(\cdot)$ and $\text{MLP}_r^0(\cdot)$ share the same parameter, restricting them in the
235 same feature space.

236 *Backbone.* We use a sequence of GraphSAGE layers as our backbone [9],
237 to allow interactions respectively among windows and exemplars. These in-
238 teractions evolve each window feature under their neighborhood structural
239 information and smooth each exemplar with neighbors to provide more ac-
240 curate features that are used in the following layers. Assuming there is L
241 layer and $l \in [1, \dots, L]$. At l^{th} layer, we mathematically define

$$\mathbf{h}_i^{l+1} = \mathbf{W}_h^l \left[\mathbf{h}_i^l \left\| \frac{1}{|\mathcal{N}_i^{\mathcal{E}_{\text{win} \rightarrow \text{win}}}|} \sum_{j \in \mathcal{N}_i^{\mathcal{E}_{\text{win} \rightarrow \text{win}}}} \mathbf{h}_j^l \right\| \right], \quad \forall v_i \in \mathcal{V}^{\text{win}} \quad (5)$$

$$\mathbf{r}_j^{l+1} = \mathbf{W}_r^l \left[\mathbf{r}_j^l \left\| \frac{1}{|\mathcal{N}_j^{\mathcal{E}_{\text{exp} \rightarrow \text{exp}}}|} \sum_{i \in \mathcal{N}_j^{\mathcal{E}_{\text{exp} \rightarrow \text{exp}}}} \mathbf{r}_i^l \right\| \right], \quad \forall v_j \in \mathcal{V}^{\text{exp}} \quad (6)$$

242 where \mathbf{W}_h^l and \mathbf{W}_r^l are linear weight matrices, and $\mathcal{N}_i^{\mathcal{E}_{\text{win} \rightarrow \text{win}}} = \{j \mid e_{jk} \in$
243 $\mathcal{E}_{\text{win} \rightarrow \text{win}} \wedge k = i\}$ and $\mathcal{N}_j^{\mathcal{E}_{\text{exp} \rightarrow \text{exp}}} = \{i \mid e_{ik} \in \mathcal{E}_{\text{exp} \rightarrow \text{exp}} \wedge k = j\}$ contain the
244 index of neighborhood nodes under the given edge set.

245 *GEB Block.* This block is concurrent with the GraphSAGE backbone, bring-
246 ing knowledge about gene expression $\mathbf{y}_j^{\text{exp}}$ of the exemplar to the window

247 feature \mathbf{h}_i^l . For brevity, we do not differentiate between the outputs of the
 248 GraphSAGE layer and the GEB block. At l^{th} layer, we project $\mathbf{y}_j^{\text{exp}}$ to \mathbf{s}_j^l ,
 249 and have interactions between \mathbf{h}_i^l and \mathbf{r}_j^l to obtain the revised window and
 250 exemplar features, i.e., \mathbf{h}_i^{l+1} and \mathbf{r}_j^{l+1} .

251 In detail, the difference between $\mathbf{h}_i^l - \mathbf{r}_j^l$ is passed to $\text{MLP}_m^t(\cdot)$, querying
 252 the potential knowledge of their gene expression difference. We chunk the
 253 outputs to $\mathbf{m}_{h,j,i}$ and $\mathbf{m}_{r,i,j}$. Then, they are used to retrieve gene expression
 254 from \mathbf{s}_j^l (i.e., a projection of $\mathbf{y}_j^{\text{exp}}$), and adjust feature significance of \mathbf{h}_i^l and \mathbf{r}_j^l .
 255 Semantically, $\mathbf{m}_{h,j,i}$ summarises ‘the existing knowledge of gene expression’,
 256 and $\mathbf{m}_{r,i,j}$ tells ‘the desired gene expression knowledge’. Mathematically, we
 257 have

$$\mathbf{h}_i^{l+1} = \mathbf{h}_i^l + \frac{1}{|\mathcal{N}_i^{\mathcal{E}_{\text{exp} \rightarrow \text{win}}}|} \sum_{j \in \mathcal{N}_i^{\mathcal{E}_{\text{exp} \rightarrow \text{win}}}} \text{MLP}_h^t([\mathbf{h}_i^l \| \mathbf{r}_j^l \| \mathbf{s}_j^l] \odot \mathbf{m}_{h,j,i}) , \quad \forall v_i \in \mathcal{V}^{\text{win}} , \quad (7)$$

$$\mathbf{r}_j^{l+1} = \mathbf{r}_j^l + \frac{1}{|\mathcal{N}_j^{\mathcal{E}_{\text{win} \rightarrow \text{exp}}}|} \sum_{i \in \mathcal{N}_j^{\mathcal{E}_{\text{win} \rightarrow \text{exp}}}} \text{MLP}_r^t([\mathbf{h}_i^l \| \mathbf{r}_j^l \| \mathbf{s}_j^l] \odot \mathbf{m}_{r,i,j}) , \quad \forall v_j \in \mathcal{V}^{\text{exp}} , \quad (8)$$

$$\mathbf{m}_{h,j,i}, \mathbf{m}_{r,i,j} = \text{Chunk}(\sigma(\text{MLP}_m^l(\mathbf{h}_i^l - \mathbf{r}_j^l))) , \quad \forall e_{ij} \in \mathcal{E}_{\text{win} \rightarrow \text{exp}} , \quad (9)$$

$$\mathbf{s}_j^l = \text{MLP}_s^l(\mathbf{y}_j^{\text{exp}}) , \quad \forall v_j \in \mathcal{V}^{\text{exp}} , \quad (10)$$

258 where $\mathcal{N}_i^{\mathcal{E}_{\text{exp} \rightarrow \text{win}}} := \{j \mid e_{jk} \in \mathcal{E}_{\text{exp} \rightarrow \text{win}} \wedge k = i\}$, $\mathcal{N}_j^{\mathcal{E}_{\text{win} \rightarrow \text{exp}}} := \{i \mid e_{ik} \in$
 259 $\mathcal{E}_{\text{win} \rightarrow \text{exp}} \wedge k = j\}$, $\text{MLP}_h^t(\cdot)$, $\text{MLP}_r^t(\cdot)$, and $\text{MLP}_m^l(\cdot)$ are multi-layer percep-
 260 trons, the \odot is element-wise multiplication, $\text{MLP}_s^l(\cdot)$ is a single layer per-
 261 ceptron, the $\text{Chunk}(\cdot)$ operator equally splits the input into two outputs,

and $\sigma(\cdot)$ is a Sigmoid function. By scaling the magnitudes of the features in Eq. 7, we directly inject the knowledge about the gene expression of the exemplars into the window feature, while we also update the exemplar feature to facilitate both interactions with the exemplars and the window feature in following layers. We finally apply a scatter mean operation to aggregate the updated features from the GEB block and the GraphSAGE layers.

Prediction Block. We extend the attention pooling from [42] to consider the contribution of exemplars to a window toward the gene expression prediction, while also encouraging the model to refine exemplars feature in a more window-dependent way. We measure the importance $\hat{\mathbf{a}}_{i,j}$ of each exemplar by using the difference between a window \mathbf{h}_i^L and a linearly projected exemplar \mathbf{r}_j^L . Then, the importance score of exemplars for a given window is normalized in the same vein with [43]. The normalized importance score $\mathbf{a}_{i,j}$ is used to weigh the exemplar and aggregate it into the window feature before making the gene expression prediction. Overall, our prediction block is defined as

$$\hat{\mathbf{a}}_{i,j} = \mathbf{h}_i^L - \mathbf{W}_a \mathbf{r}_j^L , \quad \forall e_{ij} \in \mathcal{E}_{\text{win} \rightarrow \text{exp}} , \quad (11)$$

$$\mathbf{a}_{i,j} = \frac{\exp \hat{\mathbf{a}}_{i,j}}{\sum_{k \in \mathcal{N}_i^{\text{exp} \rightarrow \text{win}}} \exp \hat{\mathbf{a}}_{i,k}} , \quad \forall e_{ij} \in \mathcal{E}_{\text{win} \rightarrow \text{exp}} , \quad (12)$$

$$\mathbf{y}_i = \mathbf{W}_y \left(\mathbf{h}_i^L + \sum_{j \in \mathcal{N}_i^{\text{exp} \rightarrow \text{win}}} \mathbf{a}_{i,j} \odot \mathbf{r}_j^L \right) , \quad \forall v_i \in \mathcal{V}^{\text{win}} \quad (13)$$

where \mathbf{W}_a and \mathbf{W}_y are linear weight matrix. Note that Eq. 12 performs computations element-wisely.

280 *Objective.* Our model is optimized with mean squared loss (i.e., the Euclidean
281 distance) \mathcal{L}_2 and batch-wise PCC \mathcal{L}_{pcc} . We have

$$\mathcal{L}_{\text{total}} = \mathcal{L}_2 + \mathcal{L}_{\text{pcc}}$$

282 **4. Experiments**

283 *Datasets.* We perform experiments in the publicly available STNet dataset
284 [3] and 10xProteomic datasets¹. The STNet dataset contains roughly 30,612
285 pairs of the slide image window and the gene expression. This dataset covers
286 68 slide images from 23 patients. Following [3], we target predicting expres-
287 sion of 250 gene types that have the largest mean across the dataset. The
288 10xProteomic dataset has 24,263 slide image windows and gene expression
289 pairs from 6 slide images. We select target gene types in the same way as the
290 STNet dataset. We apply log transformation and min-max normalization to
291 the target gene expression. Our normalization method is different from [3]
292 (they use log transformation and \mathcal{L}_1 normalization, i.e., log-transforming the
293 division of the expression of each gene type by the sum of expression of all
294 gene types, for each slide image window). Our normalization method allows
295 independent analysis of expression prediction of each gene type.

296 *4.1. Experimental Set-up*

297 *Baseline Methods.* We compare with extensive SOTA methods in domains
298 of gene expression prediction, ImageNet classification benchmarks, exemplar
299 learning, and GCNs.

¹<https://www.10xgenomics.com/resources/datasets>

- 300 • STNet [3], NSL [4], and EGN [11]. They are the SOTA methods in
301 gene expression prediction.
- 302 • ViT [43], MPViT [44] and CycleMLP [45]. We use the SOTA ImageNet
303 classification methods in our task. They are strong baselines in our
304 task. Specifically, we use ViT-B, MPViT-Base, and CycleMLP-B2.
305 Please refer to [43, 44, 45] for details.
- 306 • Retro [22] and ViTExp. We explore the SOTA exemplar learning meth-
307 ods, where the input window is represented in image \mathbf{w}_i . However,
308 Retro is originally developed for natural language processing. We adapt
309 it by providing the feature extractor output as the exemplar features.
310 ViTExp directly concatenates the exemplar features to the ViT patch
311 representation. The exemplar features are added with an embedding
312 to be differentiated from the patch representation of the slide image
313 window. Both Retro and ViTExp are based on the ViT-B architecture
314 [43].
- 315 • GraphSAGE [9], GATv2Conv [38], and TransformerConv [35]. In the
316 same vein as the last categories, we adapt the SOTA GCNs for exem-
317 plar learning, by allowing communications from windows to windows,
318 exemplars to exemplars, windows to exemplars, and exemplars to win-
319 dows.

320 *Evaluation Metrics.* We evaluate the proposed methods and alternative base-
321 lines with PCC, mean squared error (MSE), and mean absolute error (MAE).
322 We use PCC@F, PCC@S, and PCC@M to denote the first quantile, me-
323 dian, and mean of the PCC. The PCC@F verifies the least performed model
324 predictions. The PCC@S and PCC@M measure the median and mean of

325 correlations for each gene type, given predictions and ground truth for all
326 of the slide image windows. Meanwhile, the MSE and MAE measure the
327 sample-wise deviation between predictions and ground truth of each slide
328 image window for each gene type. For PCC@F, PCC@S, and PCC@M, the
329 higher value indicates better performance. In contrast, for MSE and MAE,
330 the lower value means better performance.

331 *Implementation Details.* We use the ResNet18 [10] that is officially pre-
332 trained on the ImageNet-1K dataset as our feature extractor. We implement
333 EGNN by using the *Pytorch* [46] and *PyTorch Geometric* [47] frameworks.
334 EGNN is respectively trained from scratch for 50 epochs and 300 epochs
335 on the STNet dataset and 10xProteomic dataset. We use batch size 1. All
336 windows of a slide image composite a batch with size 1 in our case. We
337 set the learning rate to 5×10^{-4} . Our weight decay is 1×10^{-4} . We use
338 a GraphSAGE backbone with hidden dimensions 512 and layers 4. All the
339 experiments are conducted with NVIDIA Tesla P100 GPUs.

Table 1: Quantitative gene expression prediction comparisons with SOTA methods on STNet dataset and 10xProteomic dataset. We bold the best results. We use ‘-’ to denote unavailable results. Models are evaluated by four-fold cross-validation and three-fold cross-validation on the above datasets. Our proposed EGNN framework consistently outperforms the SOTA methods in PCC@F $\times 10^1$, PCC@S $\times 10^1$ and PCC@M $\times 10^1$ for both datasets. GraphSAGE finds the best MSE $\times 10^2$ and MAE $\times 10^1$ on the STNet dataset.

Method	STNet Dataset					10xProteomic Dataset				
	MSE $\times 10^2$	MAE $\times 10^1$	PCC@F $\times 10^1$	PCC@S $\times 10^1$	PCC@M $\times 10^1$	MSE $\times 10^2$	MAE $\times 10^1$	PCC@F $\times 10^1$	PCC@S $\times 10^1$	PCC@M $\times 10^1$
STNet [3]	4.52	1.70	0.05	0.92	0.93	12.40	2.64	1.25	2.26	2.15
NSL [4]	-	-	-0.71	0.25	0.11	-	-	-3.73	1.84	0.25
ViT [43]	4.28	1.67	0.97	1.86	1.82	7.54	2.27	4.64	5.11	4.90
CycleMLP [45]	4.41	1.68	1.11	1.95	1.91	4.69	1.55	5.88	6.60	6.32
MPViT [44]	4.49	1.70	0.91	1.54	1.69	5.45	1.56	6.40	7.15	6.84
Retro [22]	4.53	1.71	0.99	1.74	1.79	5.25	1.65	5.46	6.35	6.04
ViTExp	4.46	1.69	0.87	1.72	1.74	5.04	1.66	5.59	6.36	6.00
EGN [11]	4.10	1.61	1.51	2.25	2.02	5.49	1.55	6.78	7.21	7.07
GraphSAGE [9]	3.79	1.57	1.98	2.90	2.83	4.98	1.65	6.82	7.42	7.25
GATv2Conv [38]	4.03	1.62	2.00	2.85	2.77	4.35	1.49	6.79	7.25	7.11
TransformerConv [35]	4.14	1.64	1.97	2.78	2.69	4.71	1.53	6.81	7.39	7.27
Ours	3.94	1.61	2.12	3.05	2.92	3.52	1.31	7.06	7.60	7.44

340 *Quantitative Evaluation.* We compare our EGNN framework with the base-
341 lines on the STNet dataset and the 10xProteomic dataset (Tab. 1). As the
342 gene expression prediction task emphasizes capturing the relativity varia-
343 tion, we bias on the PCC-related evaluation metrics, i.e., PCC@F, PCC@S,
344 and PCC@M. Our EGNN consistently achieves the best performance in
345 PCC@F, PCC@S, and PCC@M. Our findings are as follows: i) it’s worth
346 noting that GCN-based approaches achieve higher PCC-related evaluation
347 metrics than the approaches that predict gene expression individually from
348 each window, *e.g.*, EGN, the past SOTA method. GCN-based approaches
349 are capable of modeling spatial relations among windows, capturing rela-
350 tive changes among window features which improve the PCC-related eval-
351 uation metrics; ii) GraphSAGE, GATv2Conv, and TransformerConv, the
352 SOTA GCN methods, lead to the second-best performance in the STNet
353 dataset and 10xProteomic dataset in PCC-related evaluations. Meanwhile,
354 GraphSAGE finds the best MSE and MAE on the STNet dataset. These
355 GCN models all outperform their counterparts, the SOTA methods in the
356 ImageNet-1K classification task (*e.g.*, CycleMLP and MPViT), evidencing
357 our claims that reasoning spatial relations among windows out-weight fea-
358 ture learning in our gene expression prediction problem (Sec. 1); iii) the
359 PCC@F of our model significantly outperforms the baseline methods. This
360 metric evaluates the worst model capability by calculating the first quantile
361 of PCC across all gene types. The majority of gene types covered by the first
362 quantile have skewed expression distributions, which is the most challeng-
363 ing part of the prediction task. Our method has 0.012 - 0.024 higher than
364 the second-best performance from other methods in PCC@F; iv) STNet and

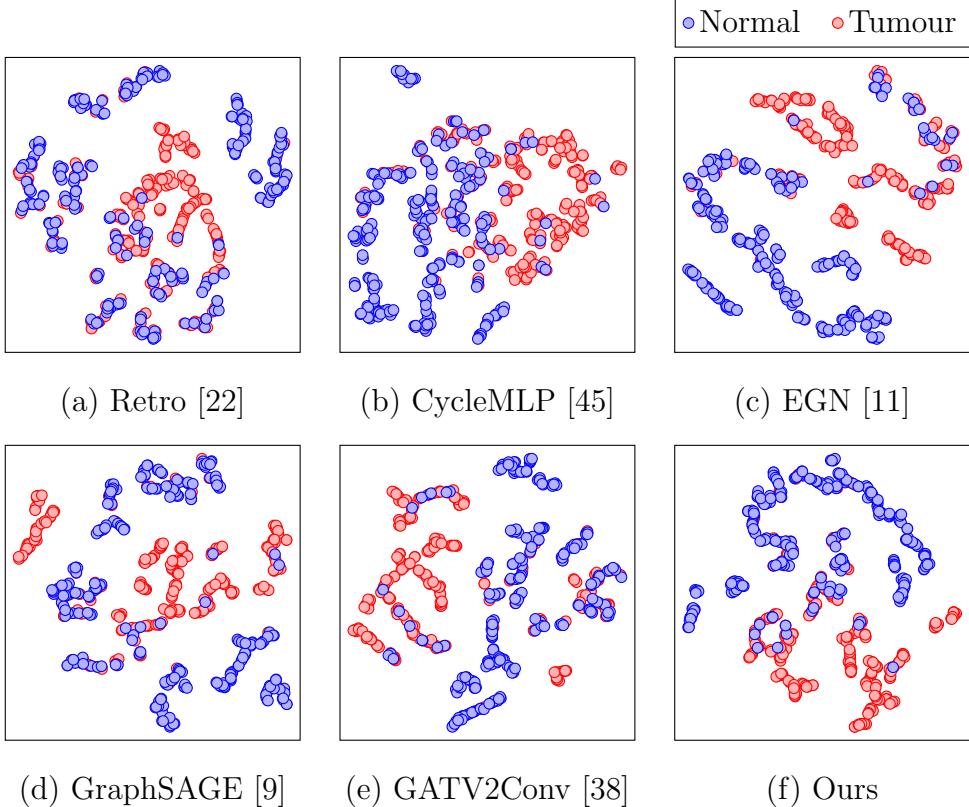


Figure 7: Quantitative evaluation of the top performed models from Tab. 1. We employ t-SNE [48] for feature dimension reduction. We use the extra labels (i.e., tumour and normal) from the STNet dataset for annotations.

365 NSL fail to achieve good performance. Again, the gene expression-related
 366 features are usually non-uniformly distributed across the slide image. They
 367 predict the gene expression of each window independently, failing to capture
 368 feature dependency among the slide image. Moreover, NSL shows a negative
 369 correlation with PCC@F. This validates our claims that predicting gene ex-
 370 pression directly from the color intensity is vulnerable, and it is only feasible
 371 in extreme cases, *e.g.*, the example of tumor-related gene expression in Sec. 1.

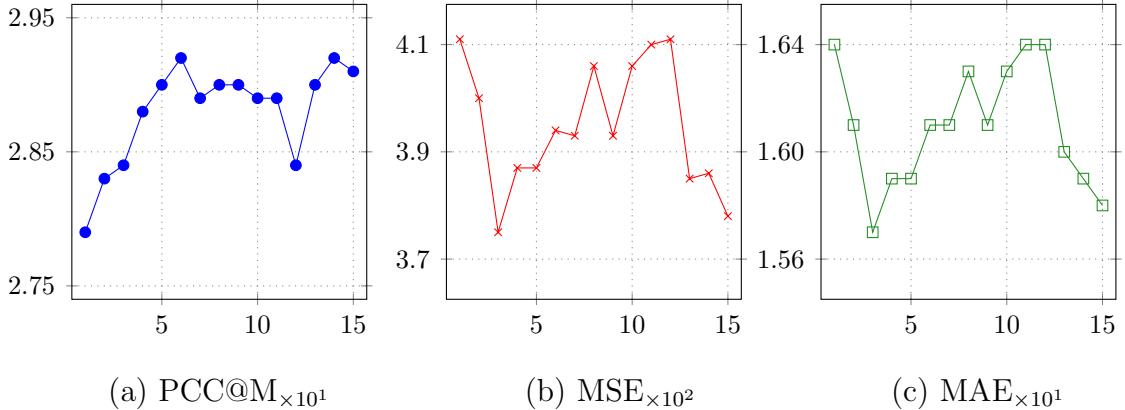


Figure 8: Ablation study on the number of exemplars used in our model. The number is varied from 1 to 15, and we respectively present PCC@M, MSE, and MAE in (a), (b), and (c).

372 *Quantitative Evaluation.* We present the latent space visualization (Fig. 7),
373 by considering the top performed models from Tab. 1. Fig. 7 (a, b, c) are
374 models that predict gene expression of each window individually, and Fig. 7
375 (d, e, f) are best performed GCN-based methods. To enable a clean visualiza-
376 tion, we randomly sample 256 slide image window features for each label, i.e.,
377 tumour and normal. Our method sufficiently separates the tumour features
378 from the normal features.

379 4.2. Ablation study

380 We study the capability of each model component by conducting a de-
381tailed ablation study on the STNet dataset.

382 *Number of Exemplars.* We present PCC@M (Fig. 8 (a)), MSE (Fig. 8 (b)),
383 and MAE (Fig. 8 (c)), by varying the number of exemplars used in our model
384 from 1 to 15. Having 6 exemplars finds the best PCC@M, and we have the
385 best MSE and MAE by using 3 exemplars. Again, our task emphasizes

Table 2: Ablation study on exemplar retrieval. We compare features from StyleEncoder, AlexNet and ResNet18 for exemplar retrieval. We use the ResNet18 with \mathcal{L}_1 for exemplar retrieval in our EGNN framework.

Feature Space	Distance Metrics	$MSE_{\times 10^2}$	$MAE_{\times 10^1}$	$PCC@M_{\times 10^1}$
StyleEncoder	<i>cosine</i>	3.87	1.59	2.61
StyleEncoder	\mathcal{L}_2	4.05	1.62	2.60
StyleEncoder	\mathcal{L}_1	4.02	1.62	2.62
AlexNet	LPIPS	4.21	1.65	2.64
ResNet18	<i>cosine</i>	4.00	1.62	2.88
ResNet18	\mathcal{L}_2	4.03	1.62	2.86
ResNet18	\mathcal{L}_1	3.94	1.61	2.92

386 capturing relative gene expression changes. Thus, our final recommendation
387 is to use 6 exemplars.

388 *Exemplar Retrieval.* We explore alternative approaches for retrieving exemplars (Tab. 2). We study exemplar retrieval strategy by considering StyleEncoder [11], AlexNet [49], and ResNet18 [10]. Among the three backbone networks, StyleEncoder is learned in an unsupervised manner with slide images, while AlexNet and ResNet18 are trained with supervision for ImageNet-1K (i.e., natural RGB images) classification task. We explore diverse distance matrices including LPIPS [50], \mathcal{L}_2 , \mathcal{L}_1 , and *cosine* (i.e., cosine similarity). We have the following findings: i) LPIPS distance retrieved exemplars lead to bad model performance. The distance is trained based on the human perception of natural RGB images. These natural RGB images are different from the slide images, leading to bad exemplar retrieves, and damaging the model performance; ii) the StyleEncoder is optimized for reconstructing slide images. However, as shown in Fig. 3, gene expression-related knowledge is not

Table 3: Ablation study on model architectures.

Settings	MSE _{$\times 10^2$}	MAE _{$\times 10^1$}	PCC@M _{$\times 10^1$}
Backbone only	4.22	1.66	2.76
W/o GEB	4.18	1.65	2.85
W/o Projector	4.10	1.63	2.84
Ours	3.94	1.61	2.92

401 balanced across the dataset, imposing biased knowledge to the StyleEncoder.
 402 When performing interactions among windows and exemplars in our model,
 403 such biases are potentially amplified and then decrease the model perfor-
 404 mance; iii) using ResNet18 with \mathcal{L}_1 distance achieves the best performance,
 405 though the pre-trained ResNet18 lacks gene expression-related knowledge.
 406 This retrieval method is an expert at encoding knowledge of image textures,
 407 benefiting from the large-scale training dataset, ImageNet-1K for our task.
 408 This allows modeling the relation among windows and exemplars in a more
 409 accurate way.

410 *Model Architectures.* We study the performance of three baseline settings,
 411 ‘Backbone only’, ‘w/o EB block’, and ‘w/o projector’, in Tab. 3. The ‘Back-
 412 bone only’ setting uses the GraphSAGE backbone architecture, considering
 413 the windows from the graph only. The ‘w/o GEB block’ setting removes the
 414 GEB block, and replaces our prediction block with a single linear layer for
 415 gene expression prediction. The ‘w/o projector’ setting replaces the projec-
 416 tor with a linear layer to unify the dimension. Our findings are as follows:
 417 i) the ‘Backbone only’ setting achieves the worst performance because of the
 418 absence of extra knowledge from the exemplar; ii) the ‘W/o GEB’ block has
 419 the second-worst performance because of a similar reason. However, it has a

Table 4: Ablation study on graph construction. The performance of our method in the ‘w/o exemplars’ setting is omitted, as our method is designed for exemplar learning only.

Settings	Method	$MSE_{\times 10^2}$	$MAE_{\times 10^1}$	$PCC@M_{\times 10^1}$
w/o exemplars	GraphSAGE	4.22	1.66	2.76
	GATv2Conv	4.34	1.68	2.75
	TransformerConv	4.45	1.70	2.62
w/ exemplars	GraphSAGE	3.79	1.57	2.83
	GATv2Conv	4.03	1.62	2.77
	TransformerConv	4.14	1.64	2.69
	Ours	3.94	1.61	2.92

420 projector to refine the ResNet18 feature to gene expression-related features;
 421 iii) with all proposed components, we have the best performance.

422 *Graph Construction.* We study if constructing the exemplar nodes in our
 423 graph can generally benefit other baseline GCN frameworks, i.e., Graph-
 424 SAGE, GATv2Conv, and TransformerConv (Tab. 4). As shown, the perfor-
 425 mance of these models has consistently been improved, by leveraging extra
 426 knowledge from the exemplars. With our GEB blocks, we have the best
 427 performance.

428 5. Limitation and Future Work

429 *Limitation.* Our method is capable to predict the gene expression of multiple
 430 windows in a single forward pass. However, computation costs have to be
 431 wasted, when a user is only interested in the gene expression of one window
 432 in a slide image. With a single window, our model is equivalent to sequential
 433 stacking linear layers, lacking interactions among windows and exemplars,
 434 and decreasing the model performance. Thus, dummy windows from the

435 slide image are needed to be sampled for forming a graph. The sampling and
436 prediction processes consume extra computation resources.

437 *Future Work.* Though exemplars could be used as anchors for allowing information propagation for windows distributed in a slide image, a more explicit
438 solution for facilitating information propagation needs to be explored. In our
439 future work, we will study downsampling and upsampling operations for our
440 exemplar and window-based graphs to build an encoder and decoder-based
441 architecture to allow broader information propagation in the slide image.
442

443 6. Conclusion and Broader Impact

444 This paper proposes an EGNN framework to accurately predict gene expression from each fine-grained area of tissue slide image, i.e., different windows. EGNN uses the GraphSAGE as a backbone while integrating with
445 exemplar learning concepts. We first have an extractor to retrieve the ex-
446emplars of the given tissue slide image window. Then, we construct a graph
447 to connect windows within the same slide image and their corresponding ex-
448emplars, and propose a GEB block to progressively revise the intermediate
449 GraphSAGE feature by reciprocating with the nearest exemplars. With ex-
450tensive experiments, we demonstrate the superiority of the EGNN framework
451 over the SOTA methods. EGNN is promising to facilitate studies on diseases
452 and novel treatments with accurate gene expression prediction.
453

455 *Acknowledgement.* The authors would like to thank Machine Learning &
456 Artificial Intelligence Future Science Platforms, CSIRO for computation re-
457 source funding.

458 **References**

- 459 [1] V. Marx, Method of the year: spatially resolved transcriptomics, *Nature
460 Methods* 18 (2021) 9–14. doi:10.1038/s41592-020-01033-y.
- 461 [2] M. Asp, J. Bergenstrhle, J. Lundeberg, Spatially resolved transcriptomesnext
462 generation tools for tissue exploration, *BioEssays* 42 (2020)
463 1900221. doi:10.1002/bies.201900221.
- 464 [3] B. He, L. Bergenstrhle, L. Stenbeck, A. Abid, A. Andersson, A. Borg,
465 J. Maaskola, J. Lundeberg, J. Zou, Integrating spatial gene expression
466 and breast tumour morphology via deep learning, *Nature Biomedical
467 Engineering* 4 (2020) 1–8. doi:10.1038/s41551-020-0578-x.
- 468 [4] M. Dawood, K. Branson, N. Rajpoot, F. u. A. A. Minhas, All you need
469 is color: Image based spatial gene expression prediction using neural
470 stain learning (08 2021).
- 471 [5] G. Huang, Z. Liu, L. van der Maaten, K. Weinberger, Densely connected
472 convolutional networks, 2017. doi:10.1109/CVPR.2017.243.
- 473 [6] S. Yang, H. Su, W. H. Hsu, W. Chen, Class-agnostic few-shot object
474 counting, in: IEEE Winter Conference on Applications of Computer
475 Vision, WACV 2021, Waikoloa, HI, USA, January 3-8, 2021, IEEE,
476 2021, pp. 869–877. doi:10.1109/WACV48630.2021.00091.
477 URL <https://doi.org/10.1109/WACV48630.2021.00091>
- 478 [7] H. Lin, Z. Ma, R. Ji, Y. Wang, X. Hong, Boosting crowd counting
479 via multifaceted attention, CoRR abs/2203.02636 (2022). arXiv:2203.

- 480 02636, doi:10.48550/arXiv.2203.02636.
- 481 URL <https://doi.org/10.48550/arXiv.2203.02636>
- 482 [8] T. N. Kipf, M. Welling, Semi-supervised classification with graph con-
483 convolutional networks, CoRR abs/1609.02907 (2016). arXiv:1609.02907.
484 URL <http://arxiv.org/abs/1609.02907>
- 485 [9] W. L. Hamilton, Z. Ying, J. Leskovec, Inductive representation learning
486 on large graphs, in: I. Guyon, U. von Luxburg, S. Bengio, H. M. Wal-
487 lach, R. Fergus, S. V. N. Vishwanathan, R. Garnett (Eds.), Advances
488 in Neural Information Processing Systems 30: Annual Conference on
489 Neural Information Processing Systems 2017, December 4-9, 2017,
490 Long Beach, CA, USA, 2017, pp. 1024–1034.
491 URL <https://proceedings.neurips.cc/paper/2017/hash/5dd9db5e033da9c6fb5ba83c7a7ebea9-Abstract.html>
492
- 493 [10] K. He, X. Zhang, S. Ren, J. Sun, Deep residual learning for image
494 recognition, 2016, pp. 770–778. doi:10.1109/CVPR.2016.90.
- 495 [11] Y. Yang, M. Hossain, E. Stone, S. Rahman, Exemplar guided deep neu-
496 ral network for spatial transcriptomics analysis of gene expression pre-
497 diction (10 2022).
- 498 [12] . Avsec, V. Agarwal, D. Visentin, J. Ledsam, A. Grabska-Barwinska,
499 K. Taylor, Y. Assael, J. Jumper, P. Kohli, D. Kelley, Effective
500 gene expression prediction from sequence by integrating long-range
501 interactions, Nature Methods 18 (2021) 1196–1203. doi:10.1038/
502 s41592-021-01252-x.

- 503 [13] B. Schmauch, A. Romagnoni, E. Pronier, C. Saillard, P. Maill, J. Calder-
504 aro, A. Kamoun, M. Sefta, S. Toldo, M. Zaslavskiy, T. Clozel, M. Moarii,
505 P. Courtiol, G. Wainrib, A deep learning model to predict rna-seq ex-
506 pression of tumours from whole slide images, Nature Communications
507 11 (08 2020). doi:10.1038/s41467-020-17678-4.
- 508 [14] X. Li, C.-Y. Wang, From bulk, single-cell to spatial rna sequencing,
509 International Journal of Oral Science 13 (12 2021). doi:10.1038/
510 s41368-021-00146-0.
- 511 [15] M. Gerlinger, A. Rowan, S. Horswell, J. Larkin, D. Endesfelder,
512 E. Gronroos, P. Martinez, N. Matthews, A. Stewart, P. Tarpey,
513 I. Varela, B. Phillimore, S. Begum, N. McDonald, A. Butler, D. Jones,
514 K. Raine, C. Latimer, C. Santos, C. Swanton, Intratumor heterogene-
515 ity and branched evolution revealed by multiregion sequencing, The
516 New England journal of medicine 366 (2012) 883–92. doi:10.1056/
517 NEJMoa1113205.
- 518 [16] V. Marx, Method of the year: spatially resolved transcriptomics, Nature
519 Methods 18 (2021) 9–14. doi:10.1038/s41592-020-01033-y.
- 520 [17] M. Bautista, A. Sanakoyeu, E. Sutter, B. Ommer, Cliquecnn: Deep
521 unsupervised exemplar learning (08 2016).
- 522 [18] B. Patro, V. Namboodiri, Deep exemplar networks for vqa and vqg (12
523 2019).
- 524 [19] Z. Chen, J. Chen, Y. Geng, J. Pan, Z. Yuan, H. Chen, Zero-Shot Visual

- 525 Question Answering Using Knowledge Graph, 2021, pp. 146–162. doi:
526 10.1007/978-3-030-88361-4_9.
- 527 [20] D. Teney, A. Hengel, Zero-shot visual question answering (11 2016).
- 528 [21] M. Farazi, S. Khan, N. Barnes, From known to the unknown: Trans-
529 ferring knowledge to answer questions about novel visual and seman-
530 tic concepts, Image and Vision Computing 103 (2020) 103985. doi:
531 10.1016/j.imavis.2020.103985.
- 532 [22] S. Borgeaud, A. Mensch, J. Hoffmann, T. Cai, E. Rutherford, K. Milli-
533 can, G. Driessche, J.-B. Lespiau, B. Damoc, A. Clark, D. Casas, A. Guy,
534 J. Menick, R. Ring, T. Hennigan, S. Huang, L. Maggiore, C. Jones,
535 A. Cassirer, L. Sifre, Improving language models by retrieving from tri-
536 lions of tokens (12 2021).
- 537 [23] Y. Wu, M. Rabe, D. Hutchins, C. Szegedy, Memorizing transformers (03
538 2022).
- 539 [24] P. Blatter, M. Kanakis, M. Danelljan, L. Gool, Efficient visual tracking
540 with exemplar transformers (12 2021).
- 541 [25] K. Guu, K. Lee, Z. Tung, P. Pasupat, M.-W. Chang, Realm: Retrieval-
542 augmented language model pre-training (02 2020).
- 543 [26] J. Thorne, A. Vlachos, C. Christodoulopoulos, A. Mittal, Fever: a large-
544 scale dataset for fact extraction and verification, 2018. doi:10.18653/
545 v1/N18-1074.

- 546 [27] P. Lewis, E. Perez, A. Piktus, F. Petroni, V. Karpukhin, N. Goyal,
547 H. Kttler, M. Lewis, W.-t. Yih, T. Rocktschel, S. Riedel, D. Kiela,
548 Retrieval-augmented generation for knowledge-intensive nlp tasks (05
549 2020).
- 550 [28] F. Petroni, P. Lewis, A. Piktus, T. Rocktschel, Y. Wu, A. Miller,
551 S. Riedel, How context affects language models' factual predictions,
552 2020.
- 553 [29] N. Moghe, S. Arora, S. Banerjee, M. Khapra, Towards exploiting back-
554 ground knowledge for building conversation systems (09 2018).
- 555 [30] M. Defferrard, X. Bresson, P. Vandergheynst, Convolutional neural
556 networks on graphs with fast localized spectral filtering, in: D. D. Lee,
557 M. Sugiyama, U. von Luxburg, I. Guyon, R. Garnett (Eds.), Advances
558 in Neural Information Processing Systems 29: Annual Conference on
559 Neural Information Processing Systems 2016, December 5-10, 2016,
560 Barcelona, Spain, 2016, pp. 3837–3845.
561 URL [https://proceedings.neurips.cc/paper/2016/hash/
562 04df4d434d481c5bb723be1b6df1ee65-Abstract.html](https://proceedings.neurips.cc/paper/2016/hash/04df4d434d481c5bb723be1b6df1ee65-Abstract.html)
- 563 [31] H. Zhu, P. Koniusz, Simple spectral graph convolution, in: 9th Inter-
564 national Conference on Learning Representations, ICLR 2021, Virtual
565 Event, Austria, May 3-7, 2021, OpenReview.net, 2021.
566 URL <https://openreview.net/forum?id=CY05T-YjWZV>
- 567 [32] P. Veličković, G. Cucurull, A. Casanova, A. Romero, P. Liò, Y. Bengio,
568 Graph attention networks, International Conference on Learning Rep-

- 569 resentations (2018).
- 570 URL <https://openreview.net/forum?id=rJXMpikCZ>
- 571 [33] F. Monti, D. Boscaini, J. Masci, E. Rodolà, J. Svoboda, M. M. Bron-
572 stein, Geometric deep learning on graphs and manifolds using mixture
573 model cnns, in: 2017 IEEE Conference on Computer Vision and Pattern
574 Recognition, CVPR 2017, Honolulu, HI, USA, July 21-26, 2017, IEEE
575 Computer Society, 2017, pp. 5425–5434. doi:10.1109/CVPR.2017.576.
576 URL <https://doi.org/10.1109/CVPR.2017.576>
- 577 [34] M. Fey, J. E. Lenssen, F. Weichert, H. Müller, Splinecnn: Fast
578 geometric deep learning with continuous b-spline kernels, in: 2018
579 IEEE Conference on Computer Vision and Pattern Recognition,
580 CVPR 2018, Salt Lake City, UT, USA, June 18-22, 2018, Computer
581 Vision Foundation / IEEE Computer Society, 2018, pp. 869–877.
582 doi:10.1109/CVPR.2018.00097.
583 URL http://openaccess.thecvf.com/content_cvpr_2018/html/Fey_SplineCNN_Fast_Geometric_CVPR_2018_paper.html
- 585 [35] Y. Shi, Z. Huang, S. Feng, H. Zhong, W. Wang, Y. Sun, Masked label
586 prediction: Unified message passing model for semi-supervised classi-
587 fication, in: Z. Zhou (Ed.), Proceedings of the Thirtieth International
588 Joint Conference on Artificial Intelligence, IJCAI 2021, Virtual Event /
589 Montreal, Canada, 19-27 August 2021, ijcai.org, 2021, pp. 1548–1554.
590 doi:10.24963/ijcai.2021/214.
591 URL <https://doi.org/10.24963/ijcai.2021/214>
- 592 [36] X. Fan, M. Gong, Y. Xie, F. Jiang, H. Li, Structured self-attention

- 593 architecture for graph-level representation learning, Pattern Recognit.
594 100 (2020) 107084. doi:10.1016/j.patcog.2019.107084.
595 URL <https://doi.org/10.1016/j.patcog.2019.107084>
- 596 [37] B. Yang, Y. Kang, L. Zhang, H. Li, GGAC: multi-relational image
597 gated GCN with attention convolutional binary neural tree for iden-
598 tifying disease with chest x-rays, Pattern Recognit. 120 (2021) 108113.
599 doi:10.1016/j.patcog.2021.108113.
600 URL <https://doi.org/10.1016/j.patcog.2021.108113>
- 601 [38] S. Brody, U. Alon, E. Yahav, How attentive are graph attention net-
602 works?, in: The Tenth International Conference on Learning Represen-
603 tations, ICLR 2022, Virtual Event, April 25-29, 2022, OpenReview.net,
604 2022.
605 URL <https://openreview.net/forum?id=F72ximsx7C1>
- 606 [39] T. Karras, S. Laine, M. Aittala, J. Hellsten, J. Lehtinen, T. Aila, Ana-
607 lyzing and improving the image quality of stylegan, in: 2020 IEEE/CVF
608 Conference on Computer Vision and Pattern Recognition, CVPR 2020,
609 Seattle, WA, USA, June 13-19, 2020, Computer Vision Foundation /
610 IEEE, 2020, pp. 8107–8116. doi:10.1109/CVPR42600.2020.00813.
611 URL https://openaccess.thecvf.com/content_CVPR_2020/html/Karras_Analyzing_and_Improving_the_Image_Quality_of_StyleGAN_CVPR_2020_paper.html
612
613
- 614 [40] Y. Yang, M. Z. Hossain, T. Gedeon, S. Rahman, S2FGAN: semanti-
615 cally aware interactive sketch-to-face translation, in: IEEE/CVF Win-
616 ter Conference on Applications of Computer Vision, WACV 2022,

617 Waikoloa, HI, USA, January 3-8, 2022, IEEE, 2022, pp. 3162–3171.

618 doi:10.1109/WACV51458.2022.00322.

619 URL <https://doi.org/10.1109/WACV51458.2022.00322>

[41] M. Fey, J. Lenssen, Fast graph representation learning with pytorch geometric (03 2019).

[42] K. Zhu, J. Wu, Residual attention: A simple but effective method for multi-label recognition (10 2021). doi:10.1109/ICCV48922.2021.00025.

[43] A. Dosovitskiy, L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai,
T. Unterthiner, M. Dehghani, M. Minderer, G. Heigold, S. Gelly,
J. Uszkoreit, N. Houlsby, An image is worth 16x16 words: Transform-
ers for image recognition at scale, in: 9th International Conference on
Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7,
2021, OpenReview.net, 2021.

631 URL <https://openreview.net/forum?id=YicbFdNTTy>

[44] Y. Lee, J. Kim, J. Willette, S. J. Hwang, Mpvit: Multi-path vision transformer for dense prediction, CoRR abs/2112.11010 (2021). arXiv: 2112.11010.

635 URL <https://arxiv.org/abs/2112.11010>

[45] S. Chen, E. Xie, C. Ge, D. Liang, P. Luo, Cyclemlp: A mlp-like architecture for dense prediction, CoRR abs/2107.10224 (2021). arXiv: 2107.10224.

639 URL <https://arxiv.org/abs/2107.10224>

- 640 [46] A. Paszke, S. Gross, F. Massa, A. Lerer, J. Bradbury, G. Chanan,
641 T. Killeen, Z. Lin, N. Gimelshein, L. Antiga, A. Desmaison, A. Köpf,
642 E. Z. Yang, Z. DeVito, M. Raison, A. Tejani, S. Chilamkurthy,
643 B. Steiner, L. Fang, J. Bai, S. Chintala, Pytorch: An imperative
644 style, high-performance deep learning library, in: H. M. Wallach,
645 H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. B. Fox, R. Garnett
646 (Eds.), Advances in Neural Information Processing Systems 32: Annual
647 Conference on Neural Information Processing Systems 2019, NeurIPS
648 2019, December 8-14, 2019, Vancouver, BC, Canada, 2019, pp. 8024–
649 8035.
- 650 URL [https://proceedings.neurips.cc/paper/2019/hash/
651 bdbca288fee7f92f2bfa9f7012727740-Abstract.html](https://proceedings.neurips.cc/paper/2019/hash/bdbca288fee7f92f2bfa9f7012727740-Abstract.html)
- 652 [47] M. Fey, J. E. Lenssen, Fast graph representation learning with pytorch
653 geometric, CoRR abs/1903.02428 (2019). [arXiv:1903.02428](https://arxiv.org/abs/1903.02428).
- 654 URL <http://arxiv.org/abs/1903.02428>
- 655 [48] P. E. Rauber, A. X. Falcão, A. C. Telea, Visualizing time-dependent data
656 using dynamic t-sne, in: E. Bertini, N. Elmquist, T. Wischgoll (Eds.),
657 18th Eurographics Conference on Visualization, EuroVis 2016 - Short
658 Papers, Groningen, The Netherlands, June 6-10, 2016, Eurographics
659 Association, 2016, pp. 73–77. doi:10.2312/eurovisshort.20161164.
- 660 URL <https://doi.org/10.2312/eurovisshort.20161164>
- 661 [49] A. Krizhevsky, I. Sutskever, G. E. Hinton, Imagenet classification with
662 deep convolutional neural networks, Commun. ACM 60 (6) (2017) 84–

663 90. doi:10.1145/3065386.
664 URL <http://doi.acm.org/10.1145/3065386>

665 [50] R. Zhang, P. Isola, A. A. Efros, E. Shechtman, O. Wang, The un-
666 reasonable effectiveness of deep features as a perceptual metric, in:
667 2018 IEEE Conference on Computer Vision and Pattern Recognition,
668 CVPR 2018, Salt Lake City, UT, USA, June 18-22, 2018, Computer
669 Vision Foundation / IEEE Computer Society, 2018, pp. 586–595.
670 doi:10.1109/CVPR.2018.00068.
671 URL http://openaccess.thecvf.com/content_cvpr_2018/html/
672 *Zhang_The_Unreasonable_Effectiveness_CVPR_2018_paper.html*