

REVIEW PAPER

## A Review of Location Encoding for GeoAI: Methods and Applications

Gengchen Mai<sup>a,b,c</sup>, Krzysztof Janowicz<sup>a,b</sup>, Yingjie Hu<sup>d</sup>, Song Gao<sup>e</sup>, Bo Yan<sup>a,b</sup>, Rui Zhu<sup>a,b</sup>, Ling Cai<sup>a,b</sup>, and Ni Lao<sup>f</sup>

<sup>a</sup>STKO Lab, Department of Geography, UC Santa Barbara, CA, USA, 93106;

<sup>b</sup>Center for Spatial Studies, UC Santa Barbara, CA, USA, 93106;

<sup>c</sup>Department of Computer Science, Stanford University, CA, USA, 94305;

<sup>d</sup>GeoAI Lab, Department of Geography, University at Buffalo, NY, USA, 14260;

<sup>e</sup>GeoDS Lab, Department of Geography, University of Wisconsin-Madison, WI, USA, 53706;

<sup>f</sup>Mosaix.ai, Palo Alto, CA, USA, 94303

### ARTICLE HISTORY

Compiled August 2, 2022

### ABSTRACT

A common need for artificial intelligence models in the broader geoscience is to represent and encode various types of spatial data, such as points (e.g., points of interest), polylines (e.g., trajectories), polygons (e.g., administrative regions), graphs (e.g., transportation networks), or rasters (e.g., remote sensing images), in a hidden embedding space so that they can be readily incorporated into deep learning models. One fundamental step is to encode a single point location into an embedding space, such that this embedding is learning-friendly for downstream machine learning models such as support vector machines and neural networks. We call this process *location encoding*. However, there lacks a systematic review on the concept of location encoding, its potential applications, and key challenges that need to be addressed. This paper aims to fill this gap. We first provide a formal definition of location encoding, and discuss the necessity of location encoding for GeoAI research from a machine learning perspective. Next, we provide a comprehensive survey and discussion about the current landscape of location encoding research. We classify location encoding models into different categories based on their inputs and encoding methods, and compare them based on whether they are parametric, multi-scale, distance preserving, and direction aware. We demonstrate that existing location encoding models can be *unified* under a shared formulation framework. We also discuss the application of location encoding for different types of spatial data. Finally, we point out several challenges in location encoding research that need to be solved in the future.

### KEYWORDS

Location encoding; GeoAI; representation learning; spatially-explicit machine learning

## 1. Introduction and Motivation

The rapid development of novel deep learning and representation learning techniques and the increasing availability of diverse, large-scale geospatial data have fueled substantial progress in geospatial artificial intelligence (GeoAI) research (Smith 1984, Couclelis

5 1986, Openshaw and Openshaw 1997, Janowicz *et al.* 2020). This includes progress on  
6 a wide spectrum of challenging tasks such as terrain feature detection and extraction  
7 (Li and Hsu 2020), land use classification (Zhong *et al.* 2019), navigation in the urban  
8 environment (Mirowski *et al.* 2018), image geolocalization (Weyand *et al.* 2016, Izbicki  
9 *et al.* 2019), toponym recognition and disambiguation (DeLozier *et al.* 2015, Wang *et al.*  
10 2020), geographic knowledge graph completion and summarization (Qiu *et al.* 2019,  
11 Yan *et al.* 2019), traffic forecasting (Li *et al.* 2018a), to name a few.

12 Despite the fact that these models are very different in design, they share a common  
13 characteristic - they need to *represent* (or *encode*) different types of spatial data, such  
14 as points (e.g., points of interest (POIs)), polylines (e.g., trajectories), polygons (e.g.,  
15 administrative regions), graphs/networks (e.g., transportation networks), or raster (e.g.,  
16 satellite images), in a hidden embedding space so that they can be utilized by machine  
17 learning models such as deep neural nets (NN). For raster data, this encoding process  
18 is straightforward since the regular grid structures can be directly handled by existing  
19 deep learning models such as convolutional neural networks (CNN) (Krizhevsky *et al.*  
20 2012). The representation problem gets more complicated for vector data such as point  
21 sets, polylines, polygons, and networks, which have more irregular spatial organization  
22 formats, because the concepts of location, distance, and direction among others do not  
23 have straightforward counterparts in existing NN and it is not trivial to design NN  
24 operations (e.g., convolution) for irregularly structured data (Valesia *et al.* 2019).

25 Early efforts perform data transformation operations to *convert* the underlying  
26 spatial data into a format which can be handled by existing NN modules (Wang *et al.*  
27 2019). However, this conversion process often leads to information loss. For example,  
28 many early research about point cloud classification and segmentation first *converted*  
29 3D point clouds into volumetric representations (e.g., voxelized shapes) (Maturana and  
30 Scherer 2015, Qi *et al.* 2016) or rendered them into 2D images (Su *et al.* 2015, Qi *et al.*  
31 2016). Then they applied 3D or 2D CNN on these converted data representations for the  
32 classification or segmentation tasks. These practices have a major limitation – choosing  
33 an appropriate spatial resolution for a volumetric representation is challenging (Qi *et al.*  
34 2017a). A finer spatial resolution leads to data sparsity and higher computation cost  
35 while a coarser spatial resolution provides poor prediction results.

36 The reason for performing such data conversions is a lack of means to directly handle  
37 vector data in deep neural nets. An alternative approach is to encode these spatial  
38 data models directly. The first step towards such goal is to encode a point location  
39 into an embedding space such that these location embeddings can be easily used in the  
40 downstream NN modules. This is the idea of *location encoding*.

41 Location encoding (Mac Aodha *et al.* 2019, Mai *et al.* 2020b, Zhong *et al.* 2020,  
42 Mai *et al.* 2020a, Gao *et al.* 2019, Xu *et al.* 2018, Chu *et al.* 2019) refers to a NN-  
43 based encoding process which represents a point/location into a high dimensional  
44 vector/embedding such that this embedding can preserve different spatial information  
45 (e.g., distance, direction) and, at the same time, be *learning-friendly*

46 for downstream machine learning (ML) models such as neural nets and support  
47 vector machines (SVM). By *learning friendly* we mean that the downstream model does  
48 not need to be very complex and does not require lots of training data to prevent model  
49 overfitting. The encoding results are called location embeddings. And the corresponding  
50 NN architecture is called *location encoder*, which is a general-purpose model that can  
51 be incorporated into different GeoAI models for different downstream tasks.

52 Figure 1 is an illustration of location encoding. Here, we use location-based species  
53 classification as an example of the downstream tasks which aims at predicting species  $y$   
54 based on a given location  $\mathbf{x}$ . The training objective is to learn the conditional distribution

55  $P(y|\mathbf{x})$ , i.e., the probability of observing  $y$  given  $\mathbf{x}$ , which is highly non-linear. The  
56 idea of location encoding can be understood as a feature decomposition process which  
57 decomposes location  $\mathbf{x}$  (e.g., a two-dimensional vector of latitude and longitude) into  
58 a learning-friendly high dimensional vector (e.g., a vector with 100 dimensions), such  
59 that the highly non-linear distribution  $P(y|\mathbf{x})$  can be learned with a relatively simple  
60 learner such as a linear SVM or a shallow NN model  $M()$ . The key benefits of such an  
61 architecture are to require less training data with simpler learners, and the possibility  
62 to leverage unsupervised training to better learn the location representations.

63 [Figure 1 about here.]

64 Recently, the effectiveness of location encoding has been demonstrated in multiple  
65 GeoAI tasks including geo-aware image classification (Yin *et al.* 2019, Chu *et al.* 2019,  
66 Mac Aodha *et al.* 2019, Mai *et al.* 2020b), POI classification (Mai *et al.* 2020b), place  
67 annotation (Yin *et al.* 2019), trajectory prediction (Xu *et al.* 2018, Yin *et al.* 2019),  
68 location privacy protection (Rao *et al.* 2020), geographic question answering (Mai  
69 *et al.* 2020a), 3D protein distribution reconstruction (Zhong *et al.* 2020), point cloud  
70 classification and segmentation (Qi *et al.* 2017a,b, Li *et al.* 2018b), and so on. Despite  
71 these successful stories, there is still a lack of a systematic review on such a topic. This  
72 paper fills this gap by providing a comparative survey on different location encoding  
73 models. We give a general conceptual formulation framework which unifies almost all  
74 existing location encoding methods.

75 It is worth mentioning that the location encoding discussed in this work is different  
76 from the traditional location encoding systems (i.e., geocoding systems)<sup>1</sup> which convert  
77 geographic coordinates into codes using an encoding scheme such as Geohash or codes  
78 for partition tiles such as Open Location Code and what3words. These traditional  
79 encoding systems are designed to support navigation and spatial indexing, while the  
80 neural location encoders we present here are used to support downstream ML models.

81 The contributions of our work are as follows:

- 82 (1) Although there are multiple existing works on location encoding, the necessity to  
83 design such a model is not clear. In this work, we formally define the location  
84 encoding problem and discuss the necessity from a machine learning perspective.
- 85 (2) We conduct a systematic review on existing location encoding research. A  
86 detailed classification system for location encoders is provided and all models  
87 are reformulated under a unified framework. This allows us to identify the  
88 commonalities and differences among different location encoding models. As  
89 far as we know, this is the first review on such a topic.
- 90 (3) We extend the idea of location encoding to the broader topic of encoding different  
91 types of spatial data (e.g., polylines, polygons, graphs, and rasters). The possible  
92 solutions and challenges are discussed.
- 93 (4) To emphasize the general applicability of location encoding, we discuss its potential  
94 applications in different geoscience domains. We hope these discussions can open  
95 up new areas of research.

96 The rest of this paper is structured as follows. We first introduce a formal definition  
97 of location encoding in Section 2. Then, in Section 3, we discuss the necessity of location  
98 encoding. Next, we provide a general framework for understanding the current landscape  
99 of location encoding research and survey a collection of representative work in Section  
100 4. In Section 5, we discuss how to apply location encoding for different types of spatial

---

<sup>1</sup><https://gogeomatics.ca/location-encoding-systems-could-geographic-coordinates-be-replaced-and-at-what-cost/>

101 data. Finally, we conclude our work and discuss future research directions in Section 6.

## 102 2. Definitions

103 **Definition 2.1** (Location Encoding). Given a set of points  $\mathcal{P} = \{p_i\}$ , e.g., the locations  
104 of sensors, species occurrences, and so on, where each point (e.g., an air quality sensor)  
105  $p_i = (\mathbf{x}_i, \mathbf{v}_i)$  is associated with a location  $\mathbf{x}_i \in \mathbb{R}^L$  (e.g., a sensor’s location) in  $L$ -D  
106 space ( $L = 2, 3$ ) and attributes  $\mathbf{v}_i \in \mathbb{R}^E$  (e.g., air quality measurements). We define the  
107 location encoder as a function  $Enc^{(\mathcal{P}, \theta)}(\mathbf{x}) : \mathbb{R}^L \rightarrow \mathbb{R}^d$  ( $L \ll d$ ), which is parameterized  
108 by  $\theta$  and maps any coordinate  $\mathbf{x}$  in space to a vector representation of  $d$  dimension.  
109 This process is called *location encoding* and the results are called *location embeddings*.

110 Here,  $Enc^{(\mathcal{P}, \theta)}(\mathbf{x})$  indicates that the encoding result of  $\mathbf{x}$  may also depend on other  
111 locations in  $\mathcal{P}$ . When  $Enc^{(\mathcal{P}, \theta)}(\mathbf{x})$  is independent of other points in  $\mathcal{P}$ , we can simplify  
112 it to  $Enc^{(\theta)}(\mathbf{x})$ . Note that sometimes, the input of the location encoder can be both  
113 locations and attributes, i.e.,  $Enc^{(\mathcal{P}, \theta)}(\mathbf{x}, \mathbf{v}) : \mathbb{R}^{L+E} \rightarrow \mathbb{R}^d$ .

114 Figure 1 illustrates the idea of location encoding in Definition 2.1. The 20 2D points  
115 serve as an example of  $\mathcal{P} = \{p_i\}$  with  $L = 2$  and  $n = |\mathcal{P}| = 20$ . Note that  $Enc^{(\mathcal{P}, \theta)}(\mathbf{x})$   
116 can not only be used to encode global location  $\mathbf{x}$ , but also be utilized to encode the  
117 spatial relation between two locations, i.e., the spatial affinity vector  $\Delta_{AB} = \mathbf{x}_A - \mathbf{x}_B$ .

118 One question we may ask is *whether a location encoder can preserve spatial*  
119 *information such as distance and direction information after the encoding process*. From  
120 a spatial information preservation perspective, there are two properties we expect a  
121 location encoder to have: distance preservation and direction awareness.

122 **Property 2.1** (Distance Preservation). The distance preservation property requires  
123 *two nearby locations to have similar location embeddings*. More concretely, given  
124 any pair of location  $(\mathbf{x}_A, \mathbf{x}_B)$ , the inner product/similarity between their resulting  
125 location embeddings, i.e.,  $\langle Enc^{(\mathcal{P}, \theta)}(\mathbf{x}_A), Enc^{(\mathcal{P}, \theta)}(\mathbf{x}_B) \rangle$  monotonically decreases when  
126 the distance<sup>2</sup> between  $\mathbf{x}_A$  and  $\mathbf{x}_B$ , i.e.,  $\|\mathbf{x}_A - \mathbf{x}_B\|$ , increases.

127 [Figure 2 about here.]

128 Property 2.1 can be seen as a reflection of Tobler’s First Law of Geography (TFL)  
129 (Tobler 1970) in location encoding. The requirement of distance preservation has  
130 been adopted by multiple existing location encoding works. For example, Gao *et al.*  
131 (2019) proposed a learnable location representation model  $v(\mathbf{x})$  which consists of  
132 three sub-models: vector matrix multiplication, magnified local isometry, and global  
133 adjacency kernel. The global adjacency kernel sub-model assumes  $\langle v(\mathbf{x}_a), v(\mathbf{x}_b) \rangle =$   
134  $(Kd)f(\|\mathbf{x}_a - \mathbf{x}_b\|)$ , where  $K, d$  are predefined constants and  $f(r)$  is the adjacency  
135 kernel that decreases monotonically as  $r$  increases. It can be seen that this sub-model  
136 directly satisfies Property 2.1. Mai *et al.* (2020b) also showed that their proposed  
137 multi-scale location encoder has a similar distance preservation property.

138 **Property 2.2** (Direction Awareness). *Locations that point into similar directions have*  
139 *more similar (relative) location embeddings than those who point into very different*  
140 *directions*. More concretely, as shown in Figure 2, given  $\mathbf{x}_O$  as the reference point and  
141  $y$  axis as the global north direction,  $\mathbf{x}_A, \mathbf{x}_B$ , and  $\mathbf{x}_C$  are on the same circle centered at  
142  $\mathbf{x}_O$  and therefore share the same distance to  $\mathbf{x}_O$ . The relative spatial relation between

---

<sup>2</sup>This can be Euclidean distance, manhattan distance, geodesic distance, great circle distance, and so on.

143  $\mathbf{x}_A$  and  $\mathbf{x}_O$  is defined as  $\Delta_{AO} = \mathbf{x}_A - \mathbf{x}_O$ . The direction of  $\Delta_{AO}$ , i.e.,  $\angle_{AO}$  is defined  
 144 as the clockwise angle between  $y$  axis and  $\Delta_{AO}$ . The same logic applies to  $\mathbf{x}_B$  and  $\mathbf{x}_C$ .  
 145 We say a location encoder is direction aware if it satisfies the following property: the  
 146 inner product  $\langle Enc^{(\mathcal{P},\theta)}(\Delta_{AO}), Enc^{(\mathcal{P},\theta)}(\Delta_{BO}) \rangle > \langle Enc^{(\mathcal{P},\theta)}(\Delta_{AO}), Enc^{(\mathcal{P},\theta)}(\Delta_{CO}) \rangle$  if  
 147  $|\angle_{AO} - \angle_{BO}| < |\angle_{AO} - \angle_{CO}|$ .

148 Property 2.2 is a reflection of the Generalized First Law of Geography (Zhu *et al.*  
 149 2019b) which includes direction into the consideration of similarities. In this paper,  
 150 we call a location encoder an *isotropic* location encoder if it only preserves the spatial  
 151 proximity but ignores the variance of location embeddings when direction changes.  
 152 We need to develop direction-aware, so-called anisotropic location encoders when the  
 153 isotropicity assumption can not be held anymore. In fact, in normal spatial analysis,  
 154 isotropicity is the “default” assumption in most of the time. Although anisotropic  
 155 versions of many geospatial analysis techniques have been developed such as directional  
 156 kriging (Te Stroet and Snepvangers 2005), anisotropic clustering (Mai *et al.* 2018),  
 157 direction remains on the level of an afterthought (Zhu *et al.* 2019b). A similar situation  
 158 can be seen in the current location encoding research, or GeoAI research in general.  
 159 Compared with studies on distance preserved location encoders, there has been much less  
 160 work on how to make a location encoder direction aware. Mai *et al.* (2020b) empirically  
 161 showed that their multi-scale location encoder as well as many baseline models are  
 162 direction-aware. However, this is just a by-product from their visualization analysis  
 163 of the response maps of these location encoders. No theoretical proof is provided. As  
 164 far as we know, there is no existing research that aims at developing a direction-aware  
 165 location encoder deliberately.

166 Besides the above two spatial information preservation properties, Location Encoder  
 167  $Enc^{(\mathcal{P},\theta)}()$  should also satisfy the following two properties to ensure its generalizability.

168 **Property 2.3** (Inductive Learning Method). Location Encoder  $Enc^{(\mathcal{P},\theta)}()$  is an  
 169 inductive learning method, i.e., the pretrained location encoder can be utilized to  
 170 encode any location without retraining even if it does not appear in the training set.

171 Property 2.3 makes  $Enc^{(\mathcal{P},\theta)}()$  differ from many existing transductive-learning-based  
 172 location representation learning methods such as Location2Vec (Zhu *et al.* 2019a),  
 173 POI2Vec (Feng *et al.* 2017), and Kejriwal and Szekely (2017). For example, Kejriwal  
 174 and Szekely (2017) converted a set of GeoNames locations into a  $k$ -th nearest neighbor  
 175 graph in which locations (i.e., nodes) are linked to nearby locations by distance-weighted  
 176 edges. A random-walk-based graph embedding method is used to learn an embedding  
 177 for each location. This method is essentially a transductive learning model: when new  
 178 locations are added to the training set, the graph is modified and the whole model has  
 179 to be retrained to obtain the embeddings of new locations.

180 **Property 2.4** (Task Independence). Location Encoder  $Enc^{(\mathcal{P},\theta)}()$  should be task-  
 181 independent or so called task-agnostic, i.e., the same model architecture can be used in  
 182 different downstream tasks without any modification.

183 Property 2.4 also differentiate  $Enc^{(\mathcal{P},\theta)}()$  from some existing task-dependent location  
 184 representation approaches. For instance, both Location2Vec (Zhu *et al.* 2019a) and  
 185 POI2Vec (Feng *et al.* 2017) learned the embeddings of locations (e.g., cell stations,  
 186 POIs) based on trajectories by adopting a Word2Vec-style training objective. This kind  
 187 of location representation cannot be easily transferred to other tasks beyond human  
 188 mobility. Similarly, Gao *et al.* (2020) discretized the study area into  $N \times N$  lattice and

189 learned the embedding of each location by letting a Long Short Term Memory (LSTM)  
190 based artificial agent navigate through the study area. The learned location embeddings  
191 are used for simulation purpose but not for other geospatial tasks.

192 **Property 2.5** (Parametric Model). A *parametric model* is a learning model with a  
193 finite set of parameters  $\theta$ .

194 A parametric model is not very flexible, but the model complexity is bounded. In  
195 contrast, a *non-parametric model* assumes that the data distribution cannot be defined  
196 with a finite set of parameters  $\theta$  and the size of parameters  $\theta$  can grow as the amount  
197 of data grows (Russell and Norvig 2015). So a non-parametric model is more flexible  
198 while the model complexity is unbounded.

199 Although Property 2.1, 2.2, and 2.5 are expected for location encoders, not all  
200 models we will discuss in Section 4 have these properties. See Table 1 for the detailed  
201 comparison. However, all location encoders discussed in Section 4 satisfy Property 2.3  
202 and 2.4. So we will not discuss these two properties separately for each model.

203 For simplicity, we will use  $Enc()$  and  $Enc^{(P)}()$  to indicate  $Enc^{(\theta)}()$  and  $Enc^{(P,\theta)}()$ .

### 204 3. The necessity of location encoding for GeoAI

205 In this section we motivate the need to embed a location  $\mathbf{x} \in \mathbb{R}^L$  ( $L = 2, 3$ ) into a high  
206 dimensional vector  $Enc^{(P,\theta)}(\mathbf{x}) \in \mathbb{R}^d$ , which may seem counter-intuitive at first. We  
207 mainly address this issue from a machine learning perspective.

208 A key concept in statistical machine learning is bias-variance trade-off (Hastie *et al.*  
209 2009). On the one hand, when a learning system is required to pick one hypothesis out  
210 of a large hypothesis space (e.g., deciding the parameters of a large 24 layer neural  
211 networks), it is flexible enough to approximate almost any non-linear distribution (low  
212 bias). However, it needs a lot of training data to prevent overfitting. This is called  
213 the low bias high variance situation. On the other hand, when the hypothesis space is  
214 restricted (e.g., linear regression or single layer neural nets) the system has little chance  
215 to over fit, but might be ill-suited to model the underlying distribution and result in low  
216 performance in both training and testing sets (high bias). This situation is called low  
217 variance high bias. For many applications the data distribution is complex and highly  
218 non-linear. We may not have enough domain knowledge to design good models with low  
219 variance (the effective model complexity) and low bias (the model data mismatch) at  
220 the same time. Moreover, we might want to avoid adopting too much domain knowledge  
221 into the model design which will make the resulting model task specific. For example,  
222 the distribution of plant species (such as  $P(y|\mathbf{x})$  in Figure 1) may be highly irregular  
223 influenced by several geospatial factors and interactions among species (Mac Aodha  
224 *et al.* 2019). Kernel (smoothing) methods (e.g., Radial Basis Function (RBF)) and  
225 neural networks (e.g., feed-forward nets) are two types of most successful models which  
226 require very little domain knowledge for model design. They both have well established  
227 ways of controlling the effective model complexity. The kernel methods are more suited  
228 to low dimensional input data – modeling highly non-linear distributions with little  
229 model complexity. However, they need to store the kernels during inference time which  
230 is not memory efficient. Neural networks have more representation power which means  
231 a deep network can approximate very complex functions with no bias, while requiring  
232 more domain knowledge for model design to achieve lower variance and bias.

233 From a statistical machine learning perspective, the main purposes of location

234 encoding is to produce *learning friendly representations of geographic locations* for  
 235 downstream models such as SVM and neural networks. By learning friendly we mean  
 236 that the downstream model does not need to be very complex and require large  
 237 training samples. For example, the location encoding process may perform a feature  
 238 decomposition ( $\mathbf{x} \in \mathbb{R}^L \rightarrow Enc^{(\mathcal{P}, \theta)}(\mathbf{x}) \in \mathbb{R}^d$ , where  $L < d$ ) so that the distribution we  
 239 want to model such as  $P(y|\mathbf{x})$  in Figure 1 becomes linear in the decomposed feature  
 240 space, and a simple linear model can be applied. Figure 3 illustrates this idea by using  
 241 a simple binary classification task. If we use the original geographic coordinates  $\mathbf{x}$  as  
 242 the input features to train the binary classifier, the resulting classifier  $M_1$  will be a  
 243 complex and nonlinear function which is prone to overfitting as shown in the left of  
 244 Figure 3. After the location encoding process, the geographic coordinates feature is  
 245 decomposed so that a simple linear model  $M_2$  can be used as the binary classifier<sup>3</sup>.

246 [Figure 3 about here.]

247 [Table 1 about here.]

## 248 4. A review of the current landscape of location encoding

249 In this section, we provide a comprehensive review of the existing location encoding  
 250 techniques. Instead of enumerating every existing location encoding approach we  
 251 organize our discussion in a top-down manner. We first classify location encoding  
 252 models into different groups according to the input of location encoders and how they  
 253 manipulate the spatial features. Firstly, according to the input, we can classify the  
 254 existing location encoders into two categories: *single point location encoder*  $Enc(\mathbf{x})$   
 255 and *aggregation location encoder*  $Enc^{(\mathcal{P})}(\mathbf{x})$ .  $Enc(\mathbf{x})$  only considers the current point’s  
 256 location while  $Enc^{(\mathcal{P})}(\mathbf{x})$  additionally considers points in its neighborhood  $\mathcal{N}(\mathbf{x}) \subseteq \mathcal{P}$ .  
 257 Then, in Section 4.1,  $Enc(\mathbf{x})$  is further classified into sub-categories based on the type  
 258 of positional encoder  $PE(\mathbf{x})$ . Next, in Section 4.2,  $Enc^{(\mathcal{P})}(\mathbf{x})$  is classified based on the  
 259 used neighborhood  $\mathcal{N}(\mathbf{x})$ . The survey result is summarized in Table 1. Finally, the  
 260 comparison among different models will be discussed in Section 4.3.

### 261 4.1. Single point location encoder $Enc(\mathbf{x})$

Interestingly, most existing single point location encoders, i.e.,  $Enc(\mathbf{x})$ , (Tang *et al.*  
 2015, Gao *et al.* 2019, Xu *et al.* 2018, Chu *et al.* 2019, Mac Aodha *et al.* 2019, Mai  
*et al.* 2020b, Zhong *et al.* 2020, Rao *et al.* 2020) share a similar structure:

$$Enc(\mathbf{x}) = \mathbf{NN}(PE(\mathbf{x})), \quad (1)$$

262 Here,  $\mathbf{NN}(\cdot) : \mathbb{R}^W \rightarrow \mathbb{R}^d$  is a learnable neural network component which maps the  
 263 input position embedding  $PE(\mathbf{x}) \in \mathbb{R}^W$  into the location embedding  $Enc(\mathbf{x}) \in \mathbb{R}^d$ . A  
 264 common practice is to define  $\mathbf{NN}(\cdot)$  as a multi-layer perceptron, while Mac Aodha *et al.*  
 265 (2019) adopted a more complex  $\mathbf{NN}(\cdot)$  which includes an initial fully connected layer,  
 266 followed by a series of residual blocks. The purpose of  $\mathbf{NN}(\cdot)$  is to provide a learnable  
 267 component for the location encoder, which captures the complex interaction between  
 268 input locations and target labels.

---

<sup>3</sup>The dimensionality of the location embedding space will be larger (e.g., 32 or 128); we use 3D here for ease of illustration.

269  $PE(\cdot)$  is the most important component which distinguishes different  $Enc(\mathbf{x})$ . Usually,  
 270  $PE(\cdot)$  is a *deterministic* function which transforms location  $\mathbf{x}$  into a  $W$ -dimension  
 271 vector, so-called position embedding. The purpose of  $PE(\cdot)$  is to do location feature  
 272 normalization (Chu *et al.* 2019, Mac Aodha *et al.* 2019, Rao *et al.* 2020) and/or feature  
 273 decomposition (Mai *et al.* 2020b, Zhong *et al.* 2020) so that the output  $PE(\mathbf{x})$  is  
 274 more learning-friendly for  $\mathbf{NN}(\cdot)$ . In Table 1 we further classify different  $Enc(\mathbf{x})$  into  
 275 four sub-categories based on their  $PE(\cdot)$ : discretization-based, direct, sinusoidal, and  
 276 sinusoidal multi-scale location encoder. Each of them will be discussed in detail below.

#### 277 4.1.1. Discretization-based location encoder

278 The early pioneers (Tang *et al.* 2015) argued that GPS coordinates are rather precise  
 279 location indicators which are difficult to use by a classifier. So instead of using the  
 280 coordinates, they discretized the whole study area into grid tiles and indicates each  
 281 point by the corresponding grid that it falls into.

282 **Definition 4.1** (Discretization-based Location Encoder). A discretization-based  
 283 location encoder divides the study area (e.g., the earth surface) into regular area  
 284 units such as grids, hexagons, or triangles -  $Enc_{discretize}(\mathbf{x}) = \mathbf{NN}(PE_{discretize}(\mathbf{x}))$   
 285 where  $PE_{discretize}(\cdot)$  is usually a tile lookup function which maps  $\mathbf{x}$  to a one hot vector  
 286 that indicates the corresponding tile id it falls into.

287 **onehot.** Early work in location encoding does not really have learnable component  
 288 specific to the location encoder. For example, Tang *et al.* (2015) divided the study area  
 289 (the contiguous United States) into  $M$  rectangle grids. Given a location  $\mathbf{x}$  (the geotag of  
 290 an image),  $PE_{onehot}(\mathbf{x}) \in \mathbb{R}^M$  is a one hot vector to indicate which grid  $\mathbf{x}$  falls into and  
 291  $\mathbf{NN}(\cdot)$  as an identity function, i.e.,  $Enc_{onehot}(\mathbf{x}) = \mathbf{NN}(PE_{onehot}(\mathbf{x})) = PE_{onehot}(\mathbf{x})$ .

292 **tile.** Later Mai *et al.* (2020b) introduced *tile* as one of  $Enc_{discretize}(\mathbf{x})$  which uses a  
 293 trainable location embedding matrix as  $\mathbf{NN}(\cdot)$ . This makes it possible for the model to  
 294 benefit from unsupervised training.

295 Although  $Enc_{discretize}(\mathbf{x})$  shows promising results on tasks such as geo-aware image  
 296 classification, it has several inherent limitations: 1) Each tile embedding are trained  
 297 separately and spatial dependency is ignored, i.e., they do not have the distance  
 298 preservation property; 2) They have only one fixed spatial scale which can not effectively  
 299 handle points with varied density; 3) Choosing the correct discretization is very  
 300 challenging (Openshaw 1981, Fotheringham and Wong 1991), and incorrect choices will  
 301 significantly affect the model’s performance and efficiency (Lechner *et al.* 2012).

302 There are some possible solutions for these problems. For problem 1 we can add a  
 303 regularization term in the loss function to make nearby tile embeddings have higher  
 304 cosine similarity. For problem 2 and 3 one can adopt an adaptive discretization (Weyand  
 305 *et al.* 2016) or multi-level discretization strategy as Kulkarni *et al.* (2020) did which uses  
 306 deeper levels (smaller tiles) for higher point density areas and shallower levels (larger  
 307 tiles) for sparse areas. However, finer spatial resolution or multi-level discretization  
 308 means more tiles and more learnable parameters which can easily lead to overfitting.

#### 309 4.1.2. Direct location encoder

310 Recently, researchers adopted a rather simple approach by directly applying neural  
 311 networks to (normalized) coordinates (Xu *et al.* 2018, Chu *et al.* 2019, Rao *et al.* 2020).

312 **Definition 4.2** (Direct Location Encoder). A direct location encoder is defined as



313  $Enc_{direct}(\mathbf{x}) = \mathbf{NN}(PE_{direct}(\mathbf{x}))$  where  $PE_{direct}(\mathbf{x})$  is usually a function to normalize  
 314 or standardize the input location feature  $\mathbf{x}$  and  $\mathbf{NN}(\cdot)$  is a multi-layer perceptron.

315 **direct.** There are many slight variations of  $Enc_{direct}(\mathbf{x})$  models. Chu *et al.* (2019)  
 316 took the input longitude and latitude (i.e.,  $\mathbf{x} = [\lambda, \phi]^T$  of a image) and normalized them  
 317 to range  $[-1, 1)$  by dividing them with constant values. Similarly, in trajectory synthesis  
 318 study, Rao *et al.* (2020) deployed  $PE_{direct}(\mathbf{x})$  which standardized each trajectory point  
 319  $\mathbf{x} = [\lambda, \phi]^T$  by using the centroid of all trajectory points. In order to perform pedestrian  
 320 trajectory prediction, Xu *et al.* (2018) also designed a simple location encoder whose  
 321  $PE_{direct}(\mathbf{x})$  normalizes  $\mathbf{x}$  to  $[0, 1]$ . Without a feature decomposition step, these models  
 322 often fail to capture the fine details of data distributions, and have worse prediction  
 323 accuracy than *tile* on specific tasks.

#### 324 4.1.3. Sinusoidal location encoder

325 **Definition 4.3** (Sinusoidal Location Encoder). A sinusoidal location encoder is defined  
 326 as  $Enc_{sinu}(\mathbf{x}) = \mathbf{NN}(PE_{sinu}(\mathbf{x}))$  where  $PE_{sinu}(\mathbf{x})$  is a deterministic function which  
 327 processes  $\mathbf{x}$  with sinusoidal functions, e.g.,  $\sin(\cdot)$ , after location feature normalization.

**wrap.** Mac Aodha *et al.* (2019) developed  $Enc_{wrap}(\mathbf{x}) = \mathbf{NN}(PE_{sinu}(\mathbf{x}))$  which  
 uses sinusoidal functions to wrap the geographic coordinates. In Equation 2, longitude  
 $\lambda$  and latitude  $\phi$  are first normalized into range  $[-1, 1]$  by dividing by  $180^\circ$  and  $90^\circ$   
 accordingly and then are fed into  $\sin(\pi x)$  and  $\cos(\pi x)$  functions.

$$PE_{wrap}(\mathbf{x}) = [\sin(\pi \frac{\lambda}{180^\circ}), \cos(\pi \frac{\lambda}{180^\circ}), \sin(\pi \frac{\phi}{90^\circ}), \cos(\pi \frac{\phi}{90^\circ})], \text{ where } \mathbf{x} = (\lambda, \phi) \quad (2)$$

328 The purpose to use sinusoidal functions is to wrap geographic coordinates around the  
 329 world Mac Aodha *et al.* (2019). This ensures that longitude  $\lambda_1 = -180^\circ$  and  $\lambda_2 = 180^\circ$   
 330 have the same encoding results. However, applying this encoding strategy to latitudes  
 331 is problematic. As for  $\phi_1 = -90^\circ$  and  $\phi_2 = 90^\circ$ , i.e., the South pole and North pole,  
 332 they would have identical encoding results which is problematic. Moreover, even if we  
 333 fix this problem, *wrap* is still not a spherical distance-preserved location encoder.

#### 334 4.1.4. Sinusoidal multi-scale location encoder

335 One limitation of all location encoders we discussed so far is that they can not handle  
 336 non-uniform point density (Qi *et al.* 2017a) or mixtures of distributions with very  
 337 different characteristics (Mai *et al.* 2020b). For example, given a set of POIs, some types  
 338 tend to cluster together such as night clubs, women’s clothing, restaurants while other  
 339 POI types are rather evenly distributed such as post offices, schools, and fire stations.  
 340 Similarly, as for species occurrences, some species herd together such as wildebeests  
 341 and zebras while the individuals of other species tend to walk alone and protect their  
 342 own territory such as tigers and bears. This will result in different spatial distribution  
 343 patterns of these species occurrences. In order to jointly model these spatial distributions,  
 344 we need an encoding method which supports multi-scale representations.

345 Inspired by the position encoding architecture in Transformer (Vaswani *et al.* 2017),  
 346 researchers developed multi-scale location encoders by using sinusoidal functions with  
 347 different frequencies (Mai *et al.* 2020b, Zhong *et al.* 2020).

**Definition 4.4** (Sinusoidal Multi-Scale Location Encoder). The sinusoidal multi-scale  
 location encoder is defined as  $Enc_{sinmul}(\mathbf{x}) = \mathbf{NN}(PE_{sinmul}(\mathbf{x}))$  where  $PE_{sinmul}(\mathbf{x})$

decomposes  $\mathbf{x}$  into a multi-scale representation based on different sinusoidal functions with different frequencies:

$$PE_{simul}(\mathbf{x}) = [PE_0^{(S)}(\mathbf{x}); \dots; PE_s^{(S)}(\mathbf{x}); \dots; PE_{S-1}^{(S)}(\mathbf{x})]. \quad (3)$$

348 Here  $S$  is the total number of scales.  $PE_s^{(S)}(\mathbf{x})$  processes the location features with  
349 different sinusoidal functions whose frequency is determined by the scale  $s$ .

350 **TF.** Zhong *et al.* (2020) slightly changed the position encoding architecture in  
351 Transformer (Vaswani *et al.* 2017) and applied them in high dimension data points  
352 such as 3D Cartesian coordinates.

**Definition 4.5** (Transformer-based Location Encoder). The transformer-based location encoder  $Enc_{TF}(\mathbf{x}) = \mathbf{NN}(PE_{TF}(\mathbf{x}))$  is following Definition 4.4. For each scale  $s \in \{0, 1, \dots, S-1\}$ ,  $PE_s^{(S)}(\mathbf{x}) = PE_s^{TF}(\mathbf{x})$  is defined by Equation 4. Here,  $\mathbf{x}^{[l]}$  indicates the  $l$ th dimension of  $\mathbf{x}$ .

$$PE_s^{TF}(\mathbf{x}) = [PE_{s,1}^{TF}(\mathbf{x}); \dots; PE_{s,l}^{TF}(\mathbf{x}); \dots; PE_{s,L}^{TF}(\mathbf{x})], \quad (4a)$$

$$\text{where } PE_{s,l}^{TF}(\mathbf{x}) = [\cos(\frac{2\pi S \mathbf{x}^{[l]}}{S^{(s+1)/S}}); \sin(\frac{2\pi S \mathbf{x}^{[l]}}{S^{(s+1)/S}})], \quad \forall l = 1, 2, \dots, L. \quad (4b)$$

353 Zhong *et al.* (2020) showed that  $Enc_{TF}(\mathbf{x})$  works well for noiseless data, but for noisy  
354 data they need to exclude the top 10% highest frequency components (the smallest  
355 several  $s$ ) in Equation 3. This indicates the necessity of another parameter to control  
356 the smallest scale we consider in sinusoidal functions. That is the usage of  $\lambda_{min}$  in  
357 *theory* and *grid* which we will discussed below.

358 **theory.** Space2Vec (Mai *et al.* 2020b) introduced *theory* as a 2D multi-scale location  
359 encoder by using sinusoidal functions with different frequencies.

**Definition 4.6** (Theory Location Encoder). Let  $\mathbf{a}_1 = [1, 0]^T$ ,  $\mathbf{a}_2 = [-1/2, \sqrt{3}/2]^T$ ,  $\mathbf{a}_3 = [-1/2, -\sqrt{3}/2]^T \in \mathbb{R}^2$  be three unit vectors which are oriented  $120^\circ$  apart from each other.  $\lambda_{min}, \lambda_{max}$  are the minimum and maximum grid scale, and  $g = \frac{\lambda_{max}}{\lambda_{min}}$ .  $Enc_{theory}(\mathbf{x})$  is following Definition 4.4 where in each scale  $s \in \{0, 1, \dots, S-1\}$ ,  $PE_s^{(S)}(\mathbf{x}) = PE_s^{theory}(\mathbf{x})$  is defined in Equation 5. Here,  $\langle \cdot, \cdot \rangle$  indicates vector inner product.

$$PE_s^{theory}(\mathbf{x}) = [PE_{s,1}^{theory}(\mathbf{x}); PE_{s,2}^{theory}(\mathbf{x}); PE_{s,3}^{theory}(\mathbf{x})], \quad (5a)$$

$$\text{where } PE_{s,j}^{theory}(\mathbf{x}) = [\cos(\frac{\langle \mathbf{x}, \mathbf{a}_j \rangle}{\lambda_{min} \cdot g^{s/(S-1)}}); \sin(\frac{\langle \mathbf{x}, \mathbf{a}_j \rangle}{\lambda_{min} \cdot g^{s/(S-1)}})], \quad \forall j = 1, 2, 3. \quad (5b)$$

360 Both  $Enc_{theory}(\cdot)$  and Gao *et al.* (2019) are inspired by the grid cell research from  
361 neuroscience field (Hafting *et al.* 2005, Blair *et al.* 2007, Killian *et al.* 2012, Agarwal  
362 *et al.* 2015). In fact, Gao *et al.* (2019) inspired and laid the theoretical foundation  
363 of  $Enc_{theory}(\mathbf{x})$ . As we discussed in Section 2, Gao *et al.* (2019) proposed a location  
364 representation model  $v(\mathbf{x})$  which consists of three sub-models. They also proposed a  
365 complex-value-based location encoder  $\Psi(\mathbf{x})$  as an analytical solution for  $v(\mathbf{x})$  which  
366 inspired  $Enc_{theory}(\cdot)$ . More specifically, given two location  $\mathbf{x}_a, \mathbf{x}_b$ , Gao *et al.* (2019)  
367 proved that  $\langle \Psi(\mathbf{x}_a), \Psi(\mathbf{x}_b) \rangle = 3(1 - \frac{\beta}{4} \|\mathbf{x}_b - \mathbf{x}_a\|^2)$  where  $\beta = \|\mathbf{a}_j\|_2^2 = 1$ . That means

368 the inner products between their location embeddings increase when  $\|\mathbf{x}_b - \mathbf{x}_a\|^2$   
 369 decrease, which satisfies Property 2.1. Mai *et al.* (2020b) showed that  $Enc_{theory}(\cdot)$  also  
 370 satisfies Property 2.1 both theoretically and empirically.

371 **grid.** *grid* is another type of  $Enc_{sinmul}(\mathbf{x})$  proposed by Space2Vec (Mai *et al.* 2020b).

**Definition 4.7** (Grid Location Encoder).  $Enc_{grid}(\mathbf{x})$  follows Definition 4.4 where at  
 each scale  $s \in \{0, 1, \dots, S-1\}$ ,  $PE_s^{(S)}(\mathbf{x}) = PE_s^{grid}(\mathbf{x})$  is defined by Equation 6. Here,  
 $\lambda_{min}$ ,  $\lambda_{max}$  and  $g$  follow the same definition as Definition 4.6.

$$PE_s^{grid}(\mathbf{x}) = [PE_{s,1}^{grid}(\mathbf{x}); \dots; PE_{s,l}^{grid}(\mathbf{x}); \dots; PE_{s,L}^{grid}(\mathbf{x})], \quad (6a)$$

$$\text{where } PE_{s,l}^{grid}(\mathbf{x}) = [\cos(\frac{\mathbf{x}^{[l]}}{\lambda_{min} \cdot g^{s/(S-1)}}); \sin(\frac{\mathbf{x}^{[l]}}{\lambda_{min} \cdot g^{s/(S-1)}})], \forall l = 1, 2, \dots, L. \quad (6b)$$

372 Mai *et al.* (2020b) shows that for both *theory* and *grid*,  $\lambda_{max}$  can be directly  
 373 determined based on the size of the study area while  $\lambda_{min}$  is the critical parameter  
 374 which decides the highest spatial resolution  $Enc(\mathbf{x})$  can handle.

#### 375 4.1.5. Comparison among different $Enc(\mathbf{x})$

376 Compared with *discretize* which yields identical embeddings for  $\mathbf{x}$  that fall into the  
 377 same tile, *direct* can distinguish nearby locations, i.e.,  $Enc_{direct}(\mathbf{x}_a) \neq Enc_{direct}(\mathbf{x}_b)$ ,  
 378 if  $\mathbf{x}_a \neq \mathbf{x}_b$ . Mai *et al.* (2020b) compared them in different tasks and found out that  
 379 without an appropriate location feature normalization  $PE_{direct}(\mathbf{x})$ , *direct* will show  
 380 lower performance than *tile*. This indicates the importance of  $PE_{direct}(\mathbf{x})$ .

381 One advantage of *direct* is its simple architecture with fewer hyperparameters.  
 382 However, compared with  $PE_{sinu}(\mathbf{x})$  and  $PE_{sinmul}(\mathbf{x})$ ,  $PE_{direct}(\mathbf{x})$  is rather hard for  
 383  $NN(\cdot)$  to learn from and may produce over-generalized distributions.

384 Compared with *TF* and *grid*, *theory* has a theoretical foundation to ensure Property  
 385 2.1. However, *theory* can only be applied to point sets in 2D space. In contrast, *TF*  
 386 and *grid* lack a theoretical guarantee for Property 2.1 while they can be utilized for  
 387 points in any  $L$ -D space. *TF* and *grid* follow similar idea while *grid* has an additional  
 388 parameter  $\lambda_{min}$ , which is more flexible for data sets with different characteristics.

## 389 4.2. Aggregation location encoder $Enc^{(P)}(\mathbf{x})$

**Definition 4.8** (Aggregation Location Encoder). The aggregation location encoder  
 $Enc^{(P)}(\mathbf{x})$  jointly considers the location feature  $\mathbf{x}$  and the aggregated features from the  
 neighborhood of  $\mathbf{x}$ , denoted as  $\mathcal{N}(\mathbf{x})$ . Inspired by the Graph Neural Network (GNN)  
 framework (Xu *et al.* 2019), a generic model setup of  $Enc^{(P)}(\mathbf{x})$  can be defined as

$$\mathbf{h}_x^{(0)} = Enc(\mathbf{x}), \quad (7a)$$

$$\mathbf{g}_x^{(m)} = Agg_{\mathbf{x}_i \in \mathcal{N}(\mathbf{x})} \{\mathbf{h}_{\mathbf{x}_i}^{(m-1)}\}, \quad (7b)$$

$$\mathbf{h}_x^{(m)} = Cmb(\mathbf{h}_x^{(m-1)}, \mathbf{g}_x^{(m)}), \quad (7c)$$

$$Enc^{(P)}(\mathbf{x}) = Rdt(\mathbf{h}_x^{(M)}). \quad (7d)$$

390 It consists of  $M$  Location Encoding Aggregation (LEA) layers, which iteratively  
 391 update the location representations. In Equation 7a, the initial location embedding

392  $\mathbf{h}_{\mathbf{x}}^{(0)}$  can be computed based on any single point location encoder discussed in  
 393 Section 4.1. Each LEA layer constitutes one neighborhood aggregation operation  
 394  $Agg_{\mathbf{x}_i \in \mathcal{N}(\mathbf{x})}\{\cdot\}$  (Equation 7b) and a feature combination operation  $Cmb(\cdot, \cdot)$  (Equation  
 395 7c).  $Agg_{\mathbf{x}_i \in \mathcal{N}(\mathbf{x})}\{\cdot\}$  aggregates the feature  $\mathbf{h}_{\mathbf{x}_i}^{(m-1)}$  of  $\mathbf{x}_i$  in  $\mathcal{N}(\mathbf{x})$  from the previous LEA  
 396 layer which can be seen as an analogy of the convolution operation of CNN on point sets.  
 397  $Agg_{\mathbf{x}_i \in \mathcal{N}(\mathbf{x})}\{\cdot\}$  can be element-wise max/min/mean pooling, sum, or any permutation  
 398 invariant architectures (Zaheer *et al.* 2017, Qi *et al.* 2017a, Veličković *et al.* 2018, Mai  
 399 *et al.* 2019a).  $Cmb(\cdot, \cdot)$  combines the point-wise feature  $\mathbf{h}_{\mathbf{x}}^{(m-1)}$  with the aggregated  
 400 features  $\mathbf{g}_{\mathbf{x}}^{(m)}$  which can be vector concatenation, element-wise max, min or mean. In the  
 401 last layer, a readout function  $Rdt(\cdot)$ , which can be an identity function or a multi-layer  
 402 perceptron, produces the final aggregated location embedding for  $\mathbf{x}$  (See Equation 7d).

403 Equation 7 can be treated as an analogy of the GNN framework (Battaglia *et al.*  
 404 2018, Xu *et al.* 2019, Wu *et al.* 2020b) which is a general framework for different  
 405 neural network architectures applied on graphs such as GCN (Kipf and Welling 2017),  
 406 GraphSAGE (Hamilton *et al.* 2017), GG-NN (Li *et al.* 2016), GAN (Veličković *et al.*  
 407 2018), MPNN (Gilmer *et al.* 2017), R-GCN (Schlichtkrull *et al.* 2018), CGA (Mai *et al.*  
 408 2019a), TransGCN (Cai *et al.* 2019), and so on.

409  $\mathcal{N}(\mathbf{x})$  can be defined in different ways such as the top  $k$ th nearest locations (Appleby  
 410 *et al.* 2020, Mai *et al.* 2020b), locations within a buffer radius (Qi *et al.* 2017b), or  
 411 locations within the same voxel as  $\mathbf{x}$  (Zhou and Tuzel 2018). According to the definition of  
 412  $\mathcal{N}(\mathbf{x})$ , we classify  $Enc^{(\mathcal{P})}(\mathbf{x})$  into different categories: kernel, global, local neighborhood,  
 413 and hierarchical neighborhood aggregation location encoder as summarized in Table 1.  
 414 We will discuss each in detail in the following section.

#### 415 4.2.1. Kernel-based location encoder

416 A kernel-based location encoder needs two components: a predefined kernel function  
 417  $k(\cdot, \cdot)$  and a set of kernel center points  $\mathcal{Q} = \{p_j\}$ . The selection of  $k(\cdot, \cdot)$  depends on  
 418 the nature of the dataset. The popular options are RBF kernels and Mercer kernels  
 419 (Vapnik 2013).  $\mathcal{Q}$  can be equal to or a subset of the training point set, i.e.,  $\mathcal{Q} \subseteq \mathcal{P}$  or  
 420 can be a predefined point set such as the centers of regular grids (Yin *et al.* 2019).

**Definition 4.9** (Kernel-Based Location Encoder). Given a kernel function  $k(\cdot, \cdot)$  and  
 the kernel center point set  $\mathcal{Q} = \{p_j\}$ , we define the kernel-based location encoder as  
 Equation 8 by following Definition 4.8 where  $M = 1$ .

$$\mathbf{h}_{\mathbf{x}}^{(0)} = \mathbf{x}, \quad (8a)$$

$$\mathbf{g}_{\mathbf{x}}^{(1)} = Agg_{p_i \in \mathcal{N}(\mathbf{x})}^{(kernel)}\{\mathbf{h}_{\mathbf{x}_i}^{(0)}\} = \Psi(\mathbf{x}) = [k(\mathbf{x}, \mathbf{x}_1); \dots; k(\mathbf{x}, \mathbf{x}_{|\mathcal{Q}|})], \quad (8b)$$

$$\mathbf{h}_{\mathbf{x}}^{(1)} = Cmb(\mathbf{h}_{\mathbf{x}}^{(0)}, \mathbf{g}_{\mathbf{x}}^{(1)}) = \mathbf{g}_{\mathbf{x}}^{(1)}, \quad (8c)$$

$$Enc_{kernel}^{(\mathcal{P})}(\mathbf{x}) = Rdt(\mathbf{h}_{\mathbf{x}}^{(1)}) = \mathbf{NN}(\mathbf{h}_{\mathbf{x}}^{(1)}). \quad (8d)$$

421 Here, each  $\mathbf{x}$  has a fixed “neighborhood” -  $\mathcal{Q}$ , i.e.,  $\mathcal{N}(\mathbf{x}) = \mathcal{Q} = \{p_j = (\mathbf{x}_i)\}$ .  $[\cdot; \cdot]$   
 422 indicates vector concatenation.  $|\mathcal{Q}|$  indicates the total number of kernels and  $\mathbf{NN}(\cdot)$  is  
 423 a multi-layer perceptron.

424 From a location feature decomposition perspective,  $Enc_{kernel}^{(\mathcal{P})}(\mathbf{x})$  can be seen as using  
 425 kernel trick to decompose  $\mathbf{x}$  into multiple kernel features (see Section 3). From a location  
 426 encoding aggregation perspective, it can be seen as aggregating the kernel features

427 derived from the interaction between  $\mathbf{x}$  and each kernel. Since each location  $\mathbf{x} \in \mathcal{P}$   
 428 shares the same  $\mathcal{N}(\mathbf{x}) = \mathcal{Q}$ ,  $Agg_{p_i \in \mathcal{N}(\mathbf{x})}^{(kernel)}\{\cdot\}$  does not need to be a permutation invariant  
 429 function. Here we use a concatenation of the  $|\mathcal{Q}|$  kernel features of  $\mathbf{x}$ . Examples of  
 430  $Enc_{kernel}^{(\mathcal{P})}(\mathbf{x})$  include GPS2Vec (Yin *et al.* 2019) and the *rbf* baseline used by Space2Vec  
 431 (Mai *et al.* 2020b). The main difference among them are the definitions of  $\mathcal{Q}$  and  $k(\cdot, \cdot)$ .

432 **GPS2Vec.** GPS2Vec (Yin *et al.* 2019) divided the Earth into multiple UTM zones  
 433 and trained different  $Enc_{kernel}^{(\mathcal{P})}(\mathbf{x})$  separately. Each zone is further divided into  $q$  grids  
 434 whose centers are used as  $\mathcal{Q}$  and  $k(\mathbf{x}, \mathbf{x}_j) = \exp(-\frac{\|\mathbf{x} - \mathbf{x}_j\|_2}{\zeta})$ . Here  $\|\cdot\|_2$  indicates L2

435 norm and  $\zeta$  is a constant attenuation coefficient. We denote this model as  $Enc_{GPS2Vec}^{(\mathcal{P})}(\cdot)$ .

436 **rbf.** Mai *et al.* (2020b) proposed a RBF kernel based location encoder as one of their  
 437 baselines, called *rbf*.  $\mathcal{Q}$  are randomly sampled  $q$  points from  $\mathcal{P}$ , i.e.,  $\mathcal{Q} \subseteq \mathcal{P}$ . The kernel  
 438 function  $k(\mathbf{x}, \mathbf{x}_j) = \exp(-\frac{\|\mathbf{x} - \mathbf{x}_j\|_2^2}{2\sigma^2})$  is a RBF kernel. We denote it as  $Enc_{rbf}^{(\mathcal{P})}(\cdot)$ .

439 **Adapted kernel\*** Instead of using constant kernel bandwidth, Berg *et al.* (2014)  
 440 also proposed an idea of adaptive kernels to model the spatio-temporal distribution  
 441 prior of bird species. They precomputed this prior as a fixed scale kernel density  
 442 estimation (KDE) map. Here, instead of making a precomputed KDE map, we adopt  
 443 this idea to design an adaptive kernel based location encoder  $Enc_{ak}(\mathbf{x})$  in which  
 444  $k(\mathbf{x}, \mathbf{x}_j) = \exp(-\frac{\|\mathbf{x} - \mathbf{x}_j\|_2^2}{2h_\sigma(\mathbf{x})^2})$ .  $h_\sigma(\mathbf{x})$  is a kernel bandwidth function of location  $\mathbf{x}$ , e.g.,

445 defining  $h_\sigma(\mathbf{x})$  as the half of the distance from  $\mathbf{x}$  to its  $K$ th nearest neighbor. We use \*  
 446 to differentiate  $Enc_{ak}(\mathbf{x})$  from the original model proposed by Berg *et al.* (2014).

447 Despite its simple design and the ability to handle non-linear distributions,  
 448  $Enc_{kernel}(\mathbf{x})$  has several shortcomings. Firstly,  $Enc_{kernel}(\mathbf{x})$  has to memorize  $\mathcal{Q}$  at  
 449 testing time which will affect memory efficiency. Secondly, since  $\mathbf{h}_\mathbf{x}^{(1)} \in \mathbb{R}^{|\mathcal{Q}|}$ , the size  
 450 of  $\mathcal{Q}$  directly decides the number of learnable parameters in  $\mathbf{NN}(\cdot)$ . This creates a  
 451 performance-efficiency trade-off problem. When  $|\mathcal{Q}|$  is small,  $Enc_{kernel}(\mathbf{x})$  has less  
 452 memory burden.  $\mathbf{NN}(\cdot)$  also has fewer learnable parameters which need less training  
 453 data and is less likely to over fit. However,  $\mathcal{Q}$  is rather distributed sparsely over the  
 454 study area which affects the quality of the encoding results. When  $|\mathcal{Q}|$  is rather large,  
 455 the encoding results are more accurate, but we need a huge amount of memory to store  
 456  $\mathcal{Q}$  and  $\mathbf{NN}(\cdot)$  needs more learning parameters. Moreover, the prediction of  $Enc_{kernel}(\mathbf{x})$   
 457 also depends on the distribution of  $\mathcal{Q}$  since it performs poorly on data sparse regions.

#### 458 4.2.2. Global aggregation location encoder

459 The global aggregation location encoder  $Enc_{global}^{(\mathcal{P})}(\mathbf{x})$  defines  $\mathcal{N}(\mathbf{x})$  as all locations in  
 460  $\mathcal{P}$ , i.e.,  $\mathcal{N}_{global}(\mathbf{x}) = \mathcal{P}$ .

461 **PointNet.** PointNet (Qi *et al.* 2017a) was originally proposed for 3D point cloud  
 462 classification and segmentation. Its point cloud segmentation architecture can be  
 463 formulated as a global aggregation location encoder as shown in Equation 9. Equation 9a  
 464 first embeds each point into an initial embedding with a *direct*-like single point location  
 465 encoder  $Enc_{pnet}(\mathbf{x})$ . It normalizes the input 3D point feature  $\mathbf{x} = [x, y, z]^T \in \mathbb{R}^3$  on to  
 466 a unit sphere, denoted as  $PE_{pnet}(\cdot)$ , before feeding them into  $\mathbf{NN}(\cdot)$ .  $\mathbf{NN}(\cdot)$  consists  
 467 of two affine transformations (parameterized as t-nets)  $TN_1(\cdot)$  and  $TN_2(\cdot)$ , which are  
 468 separated by a multi-layer perceptron  $MLP_1(\cdot)$ . These affine transformations help to  
 469 make the semantic labeling of a point cloud invariant to geometric transformations (Qi  
 470 *et al.* 2017a). Note that PointNet only has one LEA layer, i.e.,  $M = 1$ .  $Agg_{\mathbf{x}_i \in \mathcal{N}(\mathbf{x})}^{(pnet)}\{\cdot\}$

471 is a multi-layer perceptron  $MLP_2(\cdot)$  for each  $\mathbf{x}_i \in \mathcal{N}_{pnet}(\mathbf{x}) = \mathcal{P}$  followed by an  
 472 element-wise max pooling  $MaxPool_{\mathbf{x}_i \in \mathcal{N}_{pnet}(\mathbf{x})}\{\cdot\}$  (See Equation 9b).  $Cmb(\cdot, \cdot)$  is a  
 473 vector concatenation operation followed by a multi-layer perceptron  $MLP_3(\cdot)$  and the  
 474 readout function  $Rdt(\cdot)$  is an identity function as shown in Equation 9c, 9d.

$$\mathbf{h}_{\mathbf{x}}^{(0)} = Enc_{pnet}(\mathbf{x}) = \mathbf{NN}(PE_{pnet}(\mathbf{x})) = TN_2(MLP_1(TN_1(PE_{pnet}(\mathbf{x}))), \quad (9a)$$

$$\mathbf{g}_{\mathbf{x}}^{(1)} = Agg_{\mathbf{x}_i \in \mathcal{N}_{pnet}(\mathbf{x})}^{(pnet)}\{\mathbf{h}_{\mathbf{x}_i}^{(0)}\} = MaxPool_{\mathbf{x}_i \in \mathcal{N}_{pnet}(\mathbf{x})}\{MLP_2(\mathbf{h}_{\mathbf{x}_i}^{(0)})\}, \quad (9b)$$

$$\mathbf{h}_{\mathbf{x}}^{(1)} = Cmb^{(pnet)}(\mathbf{h}_{\mathbf{x}}^{(0)}, \mathbf{g}_{\mathbf{x}}^{(1)}) = MLP_3([\mathbf{h}_{\mathbf{x}}^{(0)}; \mathbf{g}_{\mathbf{x}}^{(1)}]), \quad (9c)$$

$$Enc_{pnet}^{(\mathcal{P})}(\mathbf{x}) = Rdt^{(pnet)}(\mathbf{h}_{\mathbf{x}}^{(1)}) = \mathbf{h}_{\mathbf{x}}^{(1)}. \quad (9d)$$

#### 475 4.2.3. Local aggregation location encoder

476 The local aggregation location encoder  $Enc_{local}^{(\mathcal{P})}(\mathbf{x})$  considers a local neighborhood  
 477  $\mathcal{N}(\mathbf{x})$  such as all locations within a buffer of  $\mathbf{x}$  with radius  $r$ , i.e.,  $\mathcal{N}_{buffer}(\mathbf{x}) = \{\mathbf{x}_i \mid$   
 478  $\mathbf{x} - \mathbf{x}_i \parallel_2 \leq r, \forall \mathbf{x}_i \in \mathcal{P} \wedge \mathbf{x}_i \neq \mathbf{x}\}$ .

**VoxelNet.** Zhou and Tuzel (2018) discretizes the 3D space into unit voxels. For any location  $\mathbf{x}$ , its neighborhood is defined as  $\mathcal{N}_{vox}(\mathbf{x}) = \{\mathbf{x}_i \mid \iota(\mathbf{x}) = \iota(\mathbf{x}_i), \forall \mathbf{x}_i \in \mathcal{P}\}$ , where  $\iota(\mathbf{x})$  is a voxel lookup function which returns the ID of the voxel in which  $\mathbf{x}$  falls into.

$$\mathbf{h}_{\mathbf{x}}^{(0)} = Enc_{vox}(\mathbf{x}) = \mathbf{NN}(PE_{vox}(\mathbf{x})) = PE_{vox}(\mathbf{x}), \quad (10a)$$

$$\mathbf{g}_{\mathbf{x}}^{(m)} = Agg_{\mathbf{x}_i \in \mathcal{N}_{vox}(\mathbf{x})}^{(vox)}\{\mathbf{h}_{\mathbf{x}_i}^{(m-1)}\} = MaxPool_{\mathbf{x}_i \in \mathcal{N}_{vox}(\mathbf{x})}\{FCN^{(m)}(\mathbf{h}_{\mathbf{x}_i}^{(m-1)})\}, \quad (10b)$$

$$\mathbf{h}_{\mathbf{x}}^{(m)} = Cmb^{(vox)}(\mathbf{h}_{\mathbf{x}}^{(m-1)}, \mathbf{g}_{\mathbf{x}}^{(m)}) = [\mathbf{h}_{\mathbf{x}}^{(m-1)}; \mathbf{g}_{\mathbf{x}}^{(m)}], \quad (10c)$$

$$Enc_{vox}^{(\mathcal{P})}(\mathbf{x}) = Rdt^{(vox)}(\mathbf{h}_{\mathbf{x}}^{(M)}) = \mathbf{h}_{\mathbf{x}}^{(M)}. \quad (10d)$$

479 First, a *direct*-like location encoder  $Enc_{vox}(\cdot)$  (Equation 10a) encodes  $\mathbf{x} = [x, y, z]^T$   
 480 into  $PE_{vox}(\mathbf{x}) = [x, y, z, x - \bar{x}, y - \bar{y}, z - \bar{z}]^T$ . Here  $\mathbf{NN}(\cdot)$  is an identity function<sup>4</sup>, and  
 481  $(\bar{x}, \bar{y}, \bar{z})$  is the centroid of all points in  $\mathcal{N}_{vox}(\mathbf{x})$ . Then aggregation (Equation 10b) is  
 482 done by a pointwise fully connected layer  $FCN^{(m)}(\cdot)$  followed by an element-wise max  
 483 pooling  $MaxPool_{\mathbf{x}_i \in \mathcal{N}_{vox}(\mathbf{x})}\{\cdot\}$ .  $\mathbf{g}_{\mathbf{x}}^{(m)}$  encodes the shape contained within the current  
 484 voxel  $\mathcal{N}_{vox}(\mathbf{x})$ .  $Cmb(\cdot, \cdot)$  is simply a vector concatenation operation (See Equation 10c),  
 485 and  $Rdt_{vox}(\cdot)$  is a identity function (See Equation 10d).

**SAGAT.** Mai *et al.* (2020b) proposed a modified graph attention network (GAT) (Veličković *et al.* 2018) to model the spatial interactions among nearby locations (e.g., POIs).  $\mathcal{N}_{knn}(\mathbf{x})$  is defined as all  $k$  nearest neighbors of location  $\mathbf{x}$ . The original model focuses on encoding feature information  $\mathbf{v}$  for each location such as POI type, but here we generalize it as a generic local aggregation location encoder (Spatial-Aware Graph

<sup>4</sup>In Zhou and Tuzel (2018), they encode each 3D point as  $[x, y, z, v, x - \bar{x}, y - \bar{y}, z - \bar{z}]^T$  where  $v$  is an attribute of each point, e.g., received reflectance. Here we skip  $v$  to focus on the location information. However, as we said in Definition 2.1, it is very easy to extend  $Enc^{(\mathcal{P}, \theta)}(\mathbf{x})$  to  $Enc^{(\mathcal{P}, \theta)}(\mathbf{x}, \mathbf{v})$ .

Attention Network, in short, **SAGAT**):

$$\mathbf{h}_{\mathbf{x}}^{(0)} = Enc_{sagat}(\mathbf{x}), \quad (11a)$$

$$\mathbf{g}_{\mathbf{x}}^{(m)} = Agg_{\mathbf{x}_i \in \mathcal{N}_{knn}(\mathbf{x})}^{(sagat)}\{\mathbf{h}_{\mathbf{x}_i}^{(m-1)}\} = \sigma\left(\frac{1}{U} \sum_{u=1}^U \sum_{\mathbf{x}_i \in \mathcal{N}_{knn}(\mathbf{x})} \alpha_{iu}^{(m)} \mathbf{h}_{\mathbf{x}_i}^{(m-1)}\right), \quad (11b)$$

$$\mathbf{h}_{\mathbf{x}}^{(m)} = Cmb^{(sagat)}(\mathbf{h}_{\mathbf{x}}^{(m-1)}, \mathbf{g}_{\mathbf{x}}^{(m)}) = \mathbf{g}_{\mathbf{x}}^{(m)}, \quad (11c)$$

$$Enc_{sagat}^{(\mathcal{P})}(\mathbf{x}) = Rdt^{(sagat)}(\mathbf{h}_{\mathbf{x}}^{(M)}) = \mathbf{h}_{\mathbf{x}}^{(M)}, \quad (11d)$$

$$\text{where } \alpha_{iu}^{(m)} = \frac{\exp(\sigma_{iu}^{(m)})}{\sum_{\mathbf{x}_o \in \mathcal{N}_{knn}(\mathbf{x})} \exp(\sigma_{ou}^{(m)})}, \quad (11e)$$

$$\sigma_{iu}^{(m)} = LeakyReLU((\mathbf{a}_u^{(m)})^T [\mathbf{h}_{\mathbf{x}}^{(m-1)}; \mathbf{h}_{\mathbf{x}_i}^{(m-1)}; Enc(\mathbf{x} - \mathbf{x}_i)]). \quad (11f)$$

486 In Equation 11a, SAGAT first uses  $Enc_{sagat}(\mathbf{x})$  to lift each location into the embedding  
 487 space. Originally, Mai *et al.* (2020b) used a feature encoder as  $Enc_{sagat}(\mathbf{x})$  to represent  
 488  $\mathbf{v}$  into a feature embedding. However, here, we define  $Enc_{sagat}(\mathbf{x})$  to be any single  
 489 point location encoder discussed in Section 4.1. Equation 11b is the most important  
 490 neighborhood aggregation step. Multi-head attention is adopted to aggregate the  
 491 neighboring location embedding  $\mathbf{h}_{\mathbf{x}_i}^{(m-1)}$  of  $\mathbf{x}_i \in \mathcal{N}_{knn}(\mathbf{x})$  from the previous layer.  $U$  is  
 492 the total number of attention heads.  $\alpha_{iu}^{(m)}$  is the attention coefficient of  $\mathbf{h}_{\mathbf{x}_i}^{(m-1)}$  in the  
 493  $u$ th attention head within the  $m$ th layer which is normalized across  $\mathcal{N}_{knn}(\mathbf{x})$  based  
 494 on Equation 11e.  $\sigma_{iu}^{(m)}$  is its non-normalized counterpart. In Equation 11f,  $\mathbf{h}_{\mathbf{x}}^{(m-1)}$  and  
 495  $\mathbf{h}_{\mathbf{x}_i}^{(m-1)}$  are the representations of the center location  $\mathbf{x}$  and its neighbor  $\mathbf{x}_i$  in the  $m-1$ -th  
 496 layer. Compared with GAT (Veličković *et al.* 2018),  $Enc(\mathbf{x} - \mathbf{x}_i)$  is added in the attention  
 497 score computation so the spatial affinity  $\Delta \mathbf{x}_i = \mathbf{x} - \mathbf{x}_i$  is considered in the aggregation  
 498 step. This practice makes SAGAT “spatial-aware”.  $\mathbf{a}_u^{(m)}$  is a learnable attention vector  
 499 for the  $u$ th attention head and  $LeakyReLU(\cdot)$  is the LeakyReLU activation function.  
 500  $Enc(\cdot)$  is an encoder for spatial relation, which can be any single-point location encoder  
 501 such as *tile*, *rbf*, *direct*, *theory*, *grid*, *TF*, and so on. Mai *et al.* (2020b) only utilized  
 502 one aggregation layer, i.e.,  $M = 1$ . Here, we generalize it into multiple layers, i.e.,  
 503  $M \geq 1$ . The combination and readout function are both identity functions.

504

[Figure 4 about here.]

505

Figure 4 illustrates the  $m$ th aggregation layer of SAGAT. In this example, we consider  
 506 three neighbors. Three yellow vectors  $\mathbf{h}_{\mathbf{x}_i}^{(m-1)}$  ( $i = 1, 2, 3$ ) are the hidden embeddings  
 507 of three neighbors. Three green vectors  $Enc(\mathbf{x} - \mathbf{x}_i)$  are the spatial displacement  
 508 embeddings which are used in the spatial-aware graph attention. The red  
 509 vector are the hidden embedding of center location  $\mathbf{x}$  in the next layer.

510

**DGCNN.** Dynamic Graph CNNs (in short, DGCNN) (Wang *et al.* 2019) is proposed  
 511 to jointly consider global shape structure and local neighborhood information when  
 512 learning on point clouds. DGCNN also has  $M$  LEA layers. Compared with other  
 513 location encoders, DGCNN has two crucial distinctions: 1) dynamic graph neighborhood  
 514  $\mathcal{N}_{dgcnn}^{(m)}(\mathbf{x})$  and 2) an edge convolution module  $EdgeConv_{\mathbf{x}_i \in \mathcal{N}_{dgcnn}^{(m)}(\mathbf{x})}^{(m)}\{\cdot\}$ .

$$\mathbf{g}_{\mathbf{x}}^{(m)} = \text{Agg}_{\mathbf{x}_i \in \mathcal{N}_{dgcnn}^{(m)}(\mathbf{x})}^{(dgcnn)} \{\mathbf{h}_{\mathbf{x}_i}^{(m-1)}\} = \text{EdgeConv}_{\mathbf{x}_i \in \mathcal{N}_{dgcnn}^{(m)}(\mathbf{x})}^{(m)} \{\mathbf{h}_{\mathbf{x}_i}^{(m-1)}\} \quad (12a)$$

$$= \text{MaxPool}_{\mathbf{x}_i \in \mathcal{N}_{dgcnn}^{(m)}(\mathbf{x})} \{MLP_e^{(m)}([\mathbf{h}_{\mathbf{x}}^{(m-1)}; \mathbf{h}_{\mathbf{x}_i}^{(m-1)} - \mathbf{h}_{\mathbf{x}}^{(m-1)}])\}, \quad (12b)$$

$$\mathbf{h}_{\mathbf{x}}^{(m)} = \text{Cmb}^{(dgcnn)}(\mathbf{h}_{\mathbf{x}}^{(m-1)}, \mathbf{g}_{\mathbf{x}}^{(m)}) = \mathbf{g}_{\mathbf{x}}^{(m)}. \quad (12c)$$

515 The former means that instead of using a fixed neighborhood as VoxelNet and  
 516 SAGAT does, DGCNN recomputes the neighborhood of  $\mathbf{x}$  for each LEA layer.  
 517 In the  $m$ th layer, given the embeddings of all locations from the previous layer  
 518  $\mathcal{H}^{(m-1)} = \{\mathbf{h}_{\mathbf{x}_j}^{(m-1)} | \forall \mathbf{x}_j \in \mathcal{P}\}$ ,  $\mathcal{N}_{dgcnn}^{(m)}(\mathbf{x})$  is defined as the top K nearest neighbors  
 519 of  $\mathbf{h}_{\mathbf{x}}^{(m-1)}$  from  $\mathcal{H}^{(m-1)}$  in the embedding space (not Euclidean space). Since  $\mathcal{H}^{(m-1)}$   
 520 is updated after each LEA layer,  $\mathcal{N}_{dgcnn}^{(m)}(\mathbf{x})$  is called a dynamic neighborhood. As  
 521 for the edge convolution,  $\text{EdgeConv}_{\mathbf{x}_i \in \mathcal{N}_{dgcnn}^{(m)}(\mathbf{x})}^{(m)}\{\cdot\}$  is an aggregation operator which  
 522 captures local geometric structure while maintaining permutation invariance. Equation  
 523 12 mainly shows how the edge convolution works while skipping other steps such as  
 524 location embedding initialization and the readout function.  $\text{EdgeConv}_{\mathbf{x}_i \in \mathcal{N}_{dgcnn}^{(m)}(\mathbf{x})}^{(m)}\{\cdot\}$   
 525 concatenates  $\mathbf{h}_{\mathbf{x}}^{(m-1)}$  and its affinity to its neighbor  $\mathbf{h}_{\mathbf{x}_i}^{(m-1)} - \mathbf{h}_{\mathbf{x}}^{(m-1)}$ . The result is fed  
 526 into a multi-layer perceptron  $MLP_e^{(m)}(\cdot)$  followed by a pointwise max pooling. The  
 527 combination operator is an identity function (See Equation 12c). Although  $\mathcal{N}_{dgcnn}^{(m)}(\mathbf{x})$   
 528 is dynamic, we only consider single-scale neighborhoods for each  $\mathbf{x}$ . So we classify  
 529 DGCNN can a local aggregation location encoder.

530 [Figure 5 about here.]

#### 531 4.2.4. Hierarchical aggregation location encoder

532 Instead of defining one neighborhood per location, the hierarchical aggregation  
 533 location encoder  $\text{Enc}_{hieagg}^{(\mathcal{P})}(\mathbf{x})$  defines a multi-scale neighborhood for location  $\mathbf{x}$ ,  
 534 i.e.,  $\mathcal{N}_{hieagg}(\mathbf{x}) = \{\mathcal{N}_{hieagg}^{(1)}(\mathbf{x}), \mathcal{N}_{hieagg}^{(2)}(\mathbf{x}), \dots, \mathcal{N}_{hieagg}^{(m)}(\mathbf{x}), \dots, \mathcal{N}_{hieagg}^{(M)}(\mathbf{x})\}$  which can be  
 535 aggregated in a hierarchical manner as illustrated in Figure 5.

536 **Definition 4.10** (Hierarchical Aggregation Location Encoder).  $\text{Enc}_{hieagg}^{(\mathcal{P})}(\mathbf{x})$  first  
 537 hierarchically aggregates location features from these neighborhoods and then  
 538 propagates the aggregated features back to each location.  $\text{Enc}_{hieagg}^{(\mathcal{P})}(\mathbf{x})$  usually adopts  
 539 a Conv-DeConv (Noh *et al.* 2015) like architecture which consists of two modules: 1)  
 540 **Point Set Encoder**  $\text{PTConv}_{hieagg}^{(\mathcal{P})}(\cdot)$ : a stack of  $M$  neighborhood aggregation layers  
 541 (or called Point Conv layers (Li *et al.* 2018b)) in which each location embedding in the  
 542  $m$  layer is aggregated from its neighborhood  $\mathcal{N}_{hieagg}^{(m-1)}(\mathbf{x})$  in the previous layer; 2) **Point**  
 543 **Set Decoder**  $\text{PTDeConv}_{hieagg}^{(\mathcal{P})}(\cdot)$ : a stack of  $M$  point feature propagation layers (or  
 544 called Point DeConv layers) in which each DeConv-like architecture propagates the  
 545 location embeddings from the previous layers to a denser set of locations. This whole  
 546 architecture follows the U-Net design (Ronneberger *et al.* 2015):



$$Enc_{hieagg}^{(\mathcal{P})}(\mathbf{x}) = PTDeConv_{hieagg}^{(\mathcal{P})}(PTConv_{hieagg}^{(\mathcal{P})}(\mathbf{x})) \quad (13)$$

547 Here,  $PTConv_{hieagg}^{(\mathcal{P})}(\cdot)$  follows the model setup in Equation 7. Sometimes,  
 548  $PTDeConv_{hieagg}^{(\mathcal{P})}(\cdot)$  also follows Equation 7 such as PointCNN (Li *et al.* 2018b).

549 Figure 5 shows an illustration of this Conv-DeConv architecture of  $Enc_{hieagg}^{(\mathcal{P})}(\mathbf{x})$ .

550 **PointNet++**. As an extension of PointNet, PointNet++ (Qi *et al.* 2017b) follows  
 551 the Conv-DeConv like architecture in Figure 5 to achieve hierarchical feature aggregation  
 552 and propagation. The point set encoder  $PTConv_{pnet+}^{(\mathcal{P})}(\mathbf{x})$  consists of  $M$  Point Conv  
 553 Layers (so-called set abstraction levels (SAL) in Qi *et al.* (2017b)) each of which is  
 554 composed of three key layers:

- 555 (1) **The sampling layer** at the  $m$ th SAL samples a point subset  $\mathcal{P}^{(m)}$  from the  
 556 previous SAL -  $\mathcal{P}^{(m-1)}$ , and, therefore,  $\mathcal{P}^{(m)} \subset \mathcal{P}^{(m-1)}$  and  $\mathcal{P}^{(0)} = \mathcal{P}$ .
- 557 (2) **The grouping layer** at the  $m$ th SAL groups points in  $\mathcal{P}^{(m-1)}$  into the  
 558 neighborhood  $\mathcal{N}_{pnet+}^{(m)}(\mathbf{x}) \subseteq \mathcal{P}^{(m-1)}$  of each point  $\mathbf{x} \in \mathcal{P}^{(m)}$ .  $\mathcal{N}_{pnet+}^{(m)}(\mathbf{x})$  is defined  
 559 as all points  $\mathbf{x}_i \in \mathcal{P}^{(m-1)}$  within radius  $r^{(m)}$ .
- 560 (3) **The aggregation layer** at the  $m$ th SAL produces new location embedding  $\mathbf{h}_{\mathbf{x}}^{(m)}$   
 561 for each  $\mathbf{x} \in \mathcal{P}^{(m)}$  by aggregating  $\mathbf{h}_{\mathbf{x}_i}^{(m-1)}$  of  $\mathbf{x}_i$  in  $\mathbf{x}$ 's neighborhood  $\mathcal{N}_{pnet+}^{(m)}(\mathbf{x})$ .

$$\mathbf{h}_{\mathbf{x}}^{(0)} = Enc_{direct}(\mathbf{x}), \quad (14a)$$

$$\mathbf{g}_{\mathbf{x}}^{(m)} = Agg_{\mathbf{x}_i \in \mathcal{N}_{pnet+}^{(m)}(\mathbf{x})}^{(ssg)}\{\mathbf{h}_{\mathbf{x}_i}^{(m-1)}\} \quad (14b)$$

$$= MaxPool_{\mathbf{x}_i \in \mathcal{N}_{pnet+}^{(m)}(\mathbf{x})}\{MLP^{(m)}(\mathbf{h}_{\mathbf{x}_i}^{(m-1)})\}, \quad (14c)$$

$$\mathbf{h}_{\mathbf{x}}^{(m)} = Cmb^{(ssg)}(\mathbf{h}_{\mathbf{x}}^{(m-1)}, \mathbf{g}_{\mathbf{x}}^{(m)}) = \mathbf{g}_{\mathbf{x}}^{(m)}, \quad (14d)$$

$$PTConv_{ssg}^{(\mathcal{P})}(\mathbf{x}) = Rdt^{(ssg)}(\mathbf{h}_{\mathbf{x}}^{(M)}) = \mathbf{h}_{\mathbf{x}}^{(M)}. \quad (14e)$$

562 Qi *et al.* (2017b) proposed several versions of PointNet++: **SSG** (ablated PointNet++  
 563 with single scale grouping in each level), **MSG** (multi-scale grouping PointNet++),  
 564 and **MRG** (multi resolution grouping PointNet++). All of them are  $Enc_{hieagg}^{(\mathcal{P})}(\mathbf{x})$   
 565 according to Definition 4.10. Equation 14 shows how the point set encoder of **SSG**  
 566 works. In the  $m$ th SAL, the PointNet layer aggregates features in  $\mathcal{N}_{pnet+}^{(m)}(\mathbf{x})$  by using a  
 567 PointNet-like aggregation function  $Agg_{\mathbf{x}_i \in \mathcal{N}_{pnet+}^{(m)}(\mathbf{x})}^{(ssg)}\{\cdot\}$  -  $MLP^{(m)}(\cdot)$  followed by a max  
 568 pooling function (See Equation 14c).

569 Because  $\mathcal{P}^{(M)} \subset \mathcal{P}^{(M-1)} \subset \dots \subset \mathcal{P}^{(0)} = \mathcal{P}$ ,  $\mathcal{P}^{(M)}$  is a small subset of  $\mathcal{P}$  which can  
 570 be seen as the skeleton points of  $\mathcal{P}$ .  $PTConv_{ssg}^{(\mathcal{P})}(\cdot)$  can only produce embeddings for  
 571  $\mathbf{x} \in \mathcal{P}^{(M)}$ . To obtain embeddings for each  $\mathbf{x} \in \mathcal{P}$ ,  $PTDeConv_{ssg}^{(\mathcal{P})}(\cdot)$  is used to propagate  
 572 location features back to each  $\mathbf{x} \in \mathcal{P}$  by using an inverse distance weighted interpolation  
 573 method. This can be seen as a reverse process of  $PTConv_{ssg}^{(\mathcal{P})}(\cdot)$ . The sampled point sets  
 574  $\{\mathcal{P}^{(0)} = \mathcal{P}, \mathcal{P}^{(1)}, \mathcal{P}^{(2)}, \dots, \mathcal{P}^{(m)}, \dots, \mathcal{P}^{(M)}\}$  are used in a reverse manner to progressively  
 575 interpolate the location features back to a larger sampled point set until we get location  
 576 embeddings for all  $\mathbf{x} \in \mathcal{P}$ . This idea follows the Conv-DeConv idea in Equation 13.

577 Compared with **SSG**, **MSG** changes the point set encoder by concatenating location  
578 embeddings of the same points obtained from neighborhoods with different spatial  
579 scales. **MSG** is rather computationally expensive since it aggregates features in larger  
580 scale neighborhoods for each centroid point. **MRG** solves this by concatenating location  
581 embeddings obtained from different **SSG** point set encoders with varied numbers of  
582 SAL layers. Please refer to Qi *et al.* (2017b) for detailed description.

583 **PointCNN**. Similar to PointNet++, PointCNN (Li *et al.* 2018b) also utilizes a point  
584 set encoder with  $M$  Point Conv layers (which we call PointCNN Layers here). Each  
585 layer also consists of three key steps: sampling layer, grouping layer, aggregation layer  
586 (PointNet layer). The differences from **SSG** are mainly in grouping and aggregation  
587 layer. The  $m$ th PointCNN grouping layer defines the neighborhood  $\mathcal{N}_{ptcnn}^{(m)}(\mathbf{x}) \subset \mathcal{P}^{(m-1)}$   
588 as a set of  $K$  points uniformly sampled from the  $K\mathcal{D}$  nearest neighbors of  $\mathbf{x}$  obtained  
589 from  $\mathcal{P}^{(m-1)}$ . This can be seen as an analogy of the dilated convolution idea from the  
590 traditional CNN models and  $\mathcal{D}$  is the dilation rate. The PointCNN aggregation layer  
591 defines a convolution operation,  $Conv^{(ptcnn)}(\kappa^{(m)}, \cdot)$ , over  $\mathcal{N}_{ptcnn}^{(m)}(\mathbf{x})$  as an analogy to  
592 the convolution operation over images.  $\kappa^{(m)}$  is the convolution kernels in the  $m$ th layer.

$$\mathbf{h}_{\mathbf{x}_i, \delta}^{(m)} = MLP_{\delta}^{(m)}(\mathbf{x}_i - \mathbf{x}), \forall \mathbf{x}_i \in \mathcal{N}_{ptcnn}^{(m)}(\mathbf{x}), \quad (15a)$$

$$\mathcal{X}^{(m)} = MLP^{(m)}(\Gamma_{\mathcal{N}_{ptcnn}^{(m)}(\mathbf{x})}(\mathbf{x}_i - \mathbf{x})), \quad (15b)$$

$$\mathbf{g}_{\mathbf{x}}^{(m)} = Agg_{\mathbf{x}_i \in \mathcal{N}_{ptcnn}^{(m)}(\mathbf{x})}^{(ptcnn)}\{\mathbf{h}_{\mathbf{x}_i}^{(m-1)}\} \quad (15c)$$

$$= Conv^{(ptcnn)}(\kappa^{(m)}, \mathcal{X}^{(m)} \times \Gamma_{\mathcal{N}_{ptcnn}^{(m)}(\mathbf{x})}([\mathbf{h}_{\mathbf{x}_i, \delta}^{(m)}; \mathbf{h}_{\mathbf{x}_i}^{(m-1)}])), \quad (15d)$$

$$\mathbf{h}_{\mathbf{x}}^{(m)} = Cmb^{(ptcnn)}(\mathbf{h}_{\mathbf{x}}^{(m-1)}, \mathbf{g}_{\mathbf{x}}^{(m)}) = \mathbf{g}_{\mathbf{x}}^{(m)}, \quad (15e)$$

$$PTConv_{ptcnn}^{(\mathcal{P})}(\mathbf{x}) = Rdt^{(ptcnn)}(\mathbf{h}_{\mathbf{x}}^{(M)}) = [\mathbf{h}_{\mathbf{x}}^{(M)}; MLP_g(\mathbf{x})], \quad (15f)$$

593 Equation 15 describes how the point set encoder of PointCNN works. In Equation  
594 15a, a multiple-layer perceptron  $MLP_{\delta}^{(m)}(\cdot)$  individually lifts the spatial affinity  $\mathbf{x}_i - \mathbf{x}$   
595 for each neighbor  $\mathbf{x}_i \in \mathcal{N}_{ptcnn}^{(m)}(\mathbf{x})$  into  $\mathbf{h}_{\mathbf{x}_i, \delta}^{(m)} \in \mathbb{R}^{C_{\delta}}$ , a  $C_{\delta}$  dimensional embedding. Then  
596 in Equation 15b,  $\Gamma_{\mathcal{N}_{ptcnn}^{(m)}(\mathbf{x})}(\mathbf{x}_i - \mathbf{x}) \in \mathbb{R}^{K \times L}$  represents a stack of the spatial affinity  
597 vector  $\mathbf{x}_i - \mathbf{x} \in \mathbb{R}^L$  of all  $\mathbf{x}_i \in \mathcal{N}_{ptcnn}^{(m)}(\mathbf{x})$  which results in a  $K \times L$  matrix. The  
598 multi-layer perceptron  $MLP^{(m)}(\cdot)$  converts this matrix into a  $K \times K$   $\mathcal{X}$ -transformation  
599 matrix -  $\mathcal{X}^{(m)}$ . Next, in Equation 15d, a point convolution operator  $Conv^{(ptcnn)}(\kappa^{(m)}, \cdot)$   
600 aggregates the concatenation  $[\mathbf{h}_{\mathbf{x}_i, \delta}^{(m)}; \mathbf{h}_{\mathbf{x}_i}^{(m-1)}] \in \mathbb{R}^{C_{\delta} + d^{(m-1)}}$ .  $\mathcal{X}^{(m)}$  is used here to permute  
601 this  $K \times (C_{\delta} + d^{(m-1)})$  matrix and has to be aware of the order of all  $\mathbf{x}_i \in \mathcal{N}_{ptcnn}^{(m)}(\mathbf{x})$ .  
602 In the final readout function (Equation 15f) another  $MLP_g(\cdot)$  is used to directly lift  
603  $\mathbf{x}$  into a high dimensional embedding which is concatenated with  $\mathbf{h}_{\mathbf{x}}^{(M)} \in \mathbb{R}^{d^{(M)}}$ . This  
604 can be seen as an analogy of the skip-connection in traditional CNN. In Equation 15d,  
605 when  $m = 1$ ,  $\mathbf{h}_{\mathbf{x}_i}^{(m-1)} = \mathbf{h}_{\mathbf{x}_i}^{(0)}$  needs to be initialized. Given  $p_i = (\mathbf{x}_i, \mathbf{v}_i)$  (See Definition  
606 2.1), Li *et al.* (2018b) made  $\mathbf{h}_{\mathbf{x}_i}^{(0)} = \mathbf{v}_i$ . We also can choose  $\mathbf{h}_{\mathbf{x}_i}^{(0)} = \mathbf{x}_i$ .

607 Similar PointCNN layers are deployed in the point set decoder  $PTDeConv_{ptcnn}^{(\mathcal{P})}(\mathbf{x})$   
608 to form a Conv-DeConv like architecture. Similar  $Conv^{(ptcnn)}(\kappa^{(m)}, \cdot)$  operator is used  
609 while the only difference is  $PTDeConv_{ptcnn}^{(\mathcal{P})}(\mathbf{x})$  has more points but less feature channels  
610 in its output vs. its input, and  $\{\mathcal{P}^{(m)}\}$  are forwarded from  $PTConv_{ptcnn}^{(\mathcal{P})}(\mathbf{x})$ .

**Graph-Conv GAN**. Valsesia *et al.* (2019) presented a graph convolution based

$Enc_{hieagg}^{(\mathcal{P})}(\mathbf{x})$  for 3D point cloud generation based on Generative Adversarial Network (GAN) (Goodfellow *et al.* 2014). We denote the location encoder of the model as  $Enc_{gcgan}^{(\mathcal{P})}(\mathbf{x})$ . The same Conv-DeConv idea is used here (See Definition 4.10) and  $PTConv_{gcgan}^{(\mathcal{P})}(\mathbf{x})$  consists of  $M$  Point Conv layers each of which is composed of a sampling, a grouping, and an aggregation layer:

$$\mathbf{h}_{\mathbf{x}}^{(0)} = Enc(\mathbf{x}), \quad (16a)$$

$$\mathbf{g}_{\mathbf{x}}^{(m)} = Agg_{\mathbf{x}_i \in \mathcal{N}_{knn}^{(m)}(\mathbf{x})}^{(gcgan)} \{ \mathbf{h}_{\mathbf{x}_i}^{(m-1)} \} = \sum_{\mathbf{x}_i \in \mathcal{N}_{knn}^{(m)}(\mathbf{x})} \frac{F^{(m)}(\mathbf{h}_{\mathbf{x}_i}^{(m-1)} - \mathbf{h}_{\mathbf{x}}^{(m-1)}) \mathbf{h}_{\mathbf{x}_i}^{(m-1)}}{|\mathcal{N}_{knn}^{(m)}(\mathbf{x})|}, \quad (16b)$$

$$\mathbf{h}_{\mathbf{x}}^{(m)} = Cmb^{(gcgan)}(\mathbf{h}_{\mathbf{x}}^{(m-1)}, \mathbf{g}_{\mathbf{x}}^{(m)}) = \sigma(\mathbf{g}_{\mathbf{x}}^{(m)} + \mathbf{W}^{(m)} \mathbf{h}_{\mathbf{x}}^{(m-1)} + \mathbf{b}^{(m)}), \quad (16c)$$

$$PTConv_{gcgan}^{(\mathcal{P})}(\mathbf{x}) = Rdt^{(gcgan)}(\mathbf{h}_{\mathbf{x}}^{(M)}) = \mathbf{h}_{\mathbf{x}}^{(M)}. \quad (16d)$$

611 We denote the location embedding of  $\mathbf{x}$  at the  $m$ th layer as  $\mathbf{h}_{\mathbf{x}}^{(m)} \in \mathbb{R}^{d^{(m)}}$ .  $Enc(\mathbf{x})$  in  
 612 Equation 16a can be any single point location encoder.  $Agg_{\mathbf{x}_i \in \mathcal{N}_{knn}^{(m)}(\mathbf{x})}^{(gcgan)} \{ \mathbf{h}_{\mathbf{x}_i}^{(m-1)} \} \in \mathbb{R}^{d^{(m)}}$   
 613 in Equation 16b uses a graph convolution operator to aggregate the neighborhood  
 614  $\mathcal{N}_{knn}^{(m)}(\mathbf{x})$  at the  $m$ th layer.  $\mathcal{N}_{knn}^{(m)}(\mathbf{x})$  is defined as  $\mathbf{x}$ 's  $K$ -th nearest neighbors -  $\mathcal{N}_{knn}^{(m)}(\mathbf{x}) =$   
 615  $KNN(\mathbf{x}, K, \mathcal{P}^{(m-1)})$ .  $\mathbf{h}_{\mathbf{x}_i}^{(m-1)} - \mathbf{h}_{\mathbf{x}}^{(m-1)} \in \mathbb{R}^{d^{(m-1)}}$  capture the spatial affinity between  
 616 neighboring location  $\mathbf{x}_i$  and  $\mathbf{x}$ .  $F^{(m)}(\cdot)$  denotes a fully-connected network which regresses  
 617  $\mathbf{h}_{\mathbf{x}_i}^{(m-1)} - \mathbf{h}_{\mathbf{x}}^{(m-1)}$  into a matrix  $F^{(m)}(\mathbf{h}_{\mathbf{x}_i}^{(m-1)} - \mathbf{h}_{\mathbf{x}}^{(m-1)}) \in \mathbb{R}^{d^{(m)} \times d^{(m-1)}}$ . Equation 16c  
 618 shows how to combine the neighborhood feature  $\mathbf{g}_{\mathbf{x}}^{(m)} \in \mathbb{R}^{d^{(m)}}$  with  $\mathbf{x}$ 's own feature  
 619 from the previous layer  $\mathbf{h}_{\mathbf{x}}^{(m-1)} \in \mathbb{R}^{d^{(m-1)}}$ .  $\mathbf{W}^{(m)} \in \mathbb{R}^{d^{(m)} \times d^{(m-1)}}$  and  $\mathbf{b}^{(m)} \in \mathbb{R}^{d^{(m)}}$  are a  
 620 learnable matrix and a bias vector respectively.

621 It is worth mentioning that the main difference among SSG, PointCNN, and Graph-  
 622 Conv GAN is different aggregation layers used in their point set encoders (See Equation  
 623 14, 15, and 16). Figure 5 shows their shared Conv-DeConv architecture.

### 624 4.3. Comparison among different models

625 After introducing  $Enc(\mathbf{x})$  and  $Enc^{(\mathcal{P})}(\mathbf{x})$ , it is worth to compare them from different  
 626 aspects. First, we provide a general comparison between  $Enc(\mathbf{x})$  and  $Enc^{(\mathcal{P})}(\mathbf{x})$ :

- 627 (1)  $Enc(\mathbf{x})$  encodes  $\mathbf{x}$  independently without considering its spatial context. In  
 628 contrast,  $Enc^{(\mathcal{P})}(\mathbf{x})$  jointly consider  $\mathbf{x}$  and its neighbor  $\mathcal{N}(\mathbf{x})$ .  $Enc^{(\mathcal{P})}(\mathbf{x})$  can be  
 629 seen as a generalization of  $Enc(\mathbf{x})$  in which any  $Enc(\mathbf{x})$  can be used to compute  
 630  $\mathbf{h}_{\mathbf{x}}^{(0)}$  (See Equation 7).
- 631 (2) When a new location  $\mathbf{x}'$  is added to  $\mathcal{P}$ ,  $Enc(\mathbf{x}_i)$  is unaffected for all  $\mathbf{x}_i \in \mathcal{P}$ . In  
 632 contrast, as for  $Enc^{(\mathcal{P})}(\mathbf{x})$ , for all  $\mathbf{x}_i \in \{ \mathbf{x}_i | \mathbf{x}_i \in \mathcal{P} \wedge \mathbf{x}' \in \mathcal{N}(\mathbf{x}_i) \}$ ,  $Enc^{(\mathcal{P})}(\mathbf{x}_i)$  will  
 633 be updated since  $\mathcal{N}(\mathbf{x}_i)$  is modified. However,  $Enc_{kernel}^{(\mathcal{P})}(\mathbf{x})$  is unaffected since  
 634  $\mathcal{N}(\mathbf{x}) = \mathcal{Q}$  is the same for all  $\mathbf{x} \in \mathcal{P}$  which is unchanged.
- 635 (3)  $Enc(\mathbf{x})$  has a rather high inference speed, while the aggregation operator in  
 636  $Enc^{(\mathcal{P})}(\mathbf{x})$  is time-consuming.
- 637 (4) Since  $Enc^{(\mathcal{P})}(\mathbf{x})$  additionally considers  $\mathcal{N}(\mathbf{x})$ , it has richer features for model  
 638 prediction and has a potential higher performance compared to  $Enc(\mathbf{x})$ .

639 We can see that both  $Enc(\mathbf{x})$  and  $Enc^{(\mathcal{P})}(\mathbf{x})$  have advantages and disadvantages.  
 640 Although both of them are task-agnostic, they excel at different tasks and the model

641 selection should depend on the current task:

- 642 (1) The first criterion is the input for a single model prediction - one location  $\mathbf{x}$   
643 (e.g., geo-aware image classification) or a whole point set  $\mathcal{P}$  (e.g., point cloud  
644 segmentation). The former setup prefers  $Enc(\mathbf{x})$  given its fast inference speed.  
645 The latter one prefers  $Enc^{(\mathcal{P})}(\mathbf{x})$  since it can additionally captures spatial context  
646 information. Moreover, some tasks require producing a single embedding for the  
647 whole point set  $\mathcal{P}$  or parts of it such as point cloud classification and objection  
648 recognition. For these tasks,  $Enc^{(\mathcal{P})}(\mathbf{x})$  is the only choice.
- 649 (2) Another criterion is the preference between faster inference speed or higher  
650 prediction accuracy. A mobile application may prefer a faster inference speed in  
651 which  $Enc(\mathbf{x})$  excels. An enterprise application might prefer higher prediction  
652 accuracy where  $Enc^{(\mathcal{P})}(\mathbf{x})$  is preferred.

653 Next, we compare different sub-categories of location encoders from five different  
654 perspectives. The results shown in Table 1 will be discussed in detail as follows.

- 655 (1)  $L$ : In terms of the spatial dimension of  $\mathbf{x}$  a location encoder can handle, almost all  
656 models can handle different  $L$ , e.g.,  $L = 2, 3$ . Exceptions are *GPS2Vec*, *wrap*, and  
657 *theory*. *GPS2Vec* and *wrap* are specifically designed for GPS coordinates which  
658 can be uniquely identified by  $\phi$  and  $\lambda$ . *theory* is designed only for 2D coordinates  
659 since it is inspired by neuroscience research about grid cells which are critical for  
660 self-motion integration and navigation of mammals in a 2D space.
- 661 (2) Parametric: As shown in Table 1, all models except *rbf* and adaptive kernel\* are  
662 parametric models, which means their learnable parameters have a fixed size. *rbf*  
663 and adaptive kernel\* can be either parametric or non-parametric models. Since  
664  $\mathbf{h}_{\mathbf{x}}^{(1)} \in \mathbb{R}^{|\mathcal{Q}|}$  (See Equation 8b, 8c), the size of the kernel center set  $\mathcal{Q}$  decides the  
665 number of learnable parameters in  $\mathbf{NN}(\cdot)$  of  $Enc_{kernel}^{(\mathcal{P})}(\mathbf{x})$ . If  $\mathcal{Q} = \mathcal{P}$ , then both  
666 *rbf* and adaptive kernel\* become non-parametric models. Otherwise, if  $\mathcal{Q}$  has a  
667 fixed size, both of them are parametric model.
- 668 (3) Multi-scale: Both  $Enc_{sinmul}(\mathbf{x})$  and  $Enc_{hieagg}^{(\mathcal{P})}(\mathbf{x})$  utilize multi-scale approaches.  
669 So they are better at capturing locations with non-uniform density or a mixture  
670 of distributions with different characteristics. However, they adopt different  
671 multi-scale approaches. The former designs multi-scale representations based  
672 on  $PE(\mathbf{x})$  which uses sinusoidal functions with different frequencies, while the  
673 latter aggregates the neighborhood of  $\mathbf{x}$  in a hierarchical manner. From a practical  
674 perspective, many studies have shown that multi-scale location encoders can  
675 outperform single-scale models and models without scale related parameters  
676 on various tasks. For example, Qi *et al.* (2017b) showed that PointNet++ can  
677 outperform PointNet on both point cloud classification and segmentation task. Mai  
678 *et al.* (2020b) showed that multi-scale models (*theory* and *grid*) can outperform  
679 *tile*, *direct*, *wrap*, and *rbf* on both POI type classification task and geo-aware  
680 image classification task. Mai *et al.* (2020a) showed that *theory* can outperform  
681 *direct* on geographic question answering task.
- 682 (4) Distance Preservation: In terms of the question whether a location encoder is  
683 distance preserved (Property 2.1), we mainly consider them from an empirical  
684 perspective, e.g., whether the response map of a location encoder shows a spatial  
685 continuity pattern or a spatial heterogeneity pattern. The former implies that this  
686 model is distance preserved, while the latter indicates otherwise. For example, Mai  
687 *et al.* (2020b) systematically compared the response maps of different location  
688 encoders such as *tile*, *direct*, *wrap*, *theory*, *grid*, and *rbf*, after training on the

689 POI type classification task (See Figure 2 in Mai *et al.* (2020b)). The results  
690 show that *direct*, *wrap*, *theory*, and *rbf* are distance preserved while *tile* and  
691 *direct* are not. Moreover, as we discussed in Section 4.1.4, *theory* also have a  
692 theoretical proof for distance preservation. So we put “Yes+” in Table 1. In terms  
693 of other location encoders, since no experiment or theoretical proof has been done,  
694 whether they have this property is unknown.

695 (5) Direction Awareness: A similar logic is used here. The response maps produced by  
696 Mai *et al.* (2020b) showed that *direct* and *theory* is aware of direction information  
697 while *rbf* is not. No conclusion can be drawn for other models.

698 (6) SAGAT is a bit different. Its properties also depend on the property of the used  
699  $Enc(\cdot)$  in Equation 11f. SAGAT can be a parametric or non-parametric (e.g.,  
700 use *rbf* as  $Enc(\cdot)$  ) model. It becomes a multi-scale approach if  $Enc(\cdot)$  uses a  
701 multi-scale representation. Whether SAGAT has the distance preservation and  
702 direction awareness property also depends on the used  $Enc(\cdot)$ .

## 703 5. Applying location encoding to different types of spatial data

704 Location encoders can be directly utilized on multiple point set-based GeoAI tasks such  
705 as geo-aware image classification, POI type classification, and point cloud segmentation.  
706 However, there are many other tasks that are defined on other types of spatial data  
707 such as polylines, polygons, and graphs (networks). This section discusses the potential  
708 of location encoders to model these types of spatial data.

### 709 5.1. Polyline

710 The location-to-polyline relation can be seen as an analogy of the word-to-sentence  
711 relation. In NLP, a sentence, as an ordered sequence of words, can be encoded by  
712 different sequential neural nets such as different recurrent neural networks (RNN)  
713 (Hochreiter and Schmidhuber 1997, Cho *et al.* 2014) and Transformer (Vaswani *et al.*  
714 2017). Their idea is to feed the embedding of each word token into a sequential model  
715 at each time step to encode the whole word sequence as one single hidden state or a  
716 sequence of hidden states.

717 Similarly, we can encode a polyline as an ordered sequence of locations, by using  
718 these sequential neural network models. At each time step, we will encode the current  
719 location into a location embedding and feed it into the sequential model. In fact, several  
720 recent work about human mobility directly follow this idea. Xu *et al.* (2018) utilized a  
721 *direct* location encoder to represent each trajectory point into a location embedding.  
722 Then a trajectory, represented as a sequence of location embeddings, are encoded by  
723 an LSTM for pedestrian trajectory prediction. Similarly, Rao *et al.* (2020) proposed an  
724 LSTM-TrajGAN framework to generate privacy-preserving synthetic trajectory data in  
725 which each trajectory point was encoded by a *direct* location encoder.

### 726 5.2. Polygon

727 Encoding polygon geometries into the embedding space is a logical next step. It is  
728 very useful for several geospatial tasks which require comparing polygon geometries  
729 such as geographic entity alignment (Trisedya *et al.* 2019), spatial topological reasoning  
730 (Regalia *et al.* 2019), and geographic question answering (Mai *et al.* 2019b, 2020a, 2021).

731 However, unlike a polyline, which can be represented by an ordered sequence of  
732 locations, a polygon should be represented by all locations within it. The topological  
733 relationships between any location  $\mathbf{x}$  and a polygon should be preserved after the  
734 polygon encoding process. In other words, a polygon encoder should be topology aware.  
735 As far as we know, polygon encoding is still an ongoing research problem that does not  
736 have satisfactory solutions. Mai *et al.* (2020a) presented a geographic entity bounding  
737 box encoding model as the first step towards polygon encoding by uniformly sampling a  
738 location from within the bounding box of a geographic entity and feeding it to a location  
739 encoder. Despite its innovativeness, this model still cannot handle fine-grained polygon  
740 geometries. Yan *et al.* (2021) proposed a graph convolutional autoencoder (GCAE)  
741 which can encode simple polygons, i.e., polygons that do not intersect themselves and  
742 have no holes. GCAE converts the exterior of a simple polygon into a graph and then  
743 encodes this graph into an embedding space. The shortcoming of GCAE is that it  
744 cannot handle polygons with holes and multipolygons. Moreover, it cannot preserve  
745 the topology information. So one interesting future research direction is developing a  
746 topology-aware polygon encoder which can handle both simple and complex polygons.

### 747 5.3. Graph

748 Graph (or network) is also an important spatial data format used in multiple geospatial  
749 data sets such as transportation networks (Li *et al.* 2018a, Cai *et al.* 2020), spatial  
750 social networks (Andris 2016), and geographic knowledge graph (GeoKG) (Mai *et al.*  
751 2020a). A graph can be defined as  $\mathcal{G} = (\mathcal{V}, \mathcal{E})$  where  $\mathcal{V}$  and  $\mathcal{E}$  are the set of nodes and  
752 edges in this graph. In the geospatial domain, each node  $e \in \mathcal{V}$  or a subset of nodes in  $\mathcal{V}$   
753 is associated with a location  $\mathbf{x}$ <sup>5</sup> such as the sensor locations in a sensor network, users’  
754 locations in a spatial social network, or locations of geographic entities in a GeoKG.  
755 We further call this kind of graph a *spatially embedded graph*.

756 The early practice to encode spatially embedded graphs is to treat them as normal  
757 non-spatial graphs and use some existing GNN models or (knowledge) graph embedding  
758 models (Grover and Leskovec 2016, Bordes *et al.* 2013, Trouillon *et al.* 2017). In order  
759 to add the spatial information as additional features without significantly modifying the  
760 existing architectures, we can modify the node encoder by using one  $Enc(\mathbf{x})$  in Table 1  
761 as the node encoder or one component of it while keeping other components unchanged.  
762 The model can be trained with the same loss function. Mai *et al.* (2020a) adopted  
763 exactly this practice and developed a spatially-explicit knowledge graph embedding  
764 model. Similar ideas can be applied to other spatially embedded graphs.

765 Interestingly, other than the normal graph data, many pioneer research applied GNN  
766 models to a point set through a *point-set-to-graph* conversion. They first converted  
767 point set  $\mathcal{P}$  into a graph based on spatial relations, e.g., a k-th nearest neighbor spatial  
768 graph, in which nodes indicate points while edges are associated with pairwise distance  
769 based weights. After this conversion, a GNN model is applied on this graph so that node  
770 attribute prediction can be done based on not only the nodes’ own features but also  
771 their spatial context. Many GeoAI research has adopted this practice to tackle different  
772 tasks including air quality forecasting (Lin *et al.* 2018), place characteristics prediction  
773 (Zhu *et al.* 2020), GeoNames entity embedding learning (Kejriwal and Szekely 2017),  
774 and different spatial interpolation problems (Appleby *et al.* 2020, Wu *et al.* 2020a).  
775 We argue that *this kind of distance weighted graph method is insufficient to capture the*  
776 *relative spatial relations*, since it necessarily forfeits information about the spatial layout

---

<sup>5</sup>Each node can associate with more complex geometries, and one single location is a simplification.

777 of points. Some important spatial information is lost such as the direction relations  
778 which are important for certain tasks when isotropic assumption is not held. Instead,  
779 we advocate the idea of using any  $Enc^{(\mathcal{P})}(\mathbf{x})$  discussed in Section 4.2 for these tasks  
780 since  $Enc^{(\mathcal{P})}(\mathbf{x})$  is better at capturing spatial relations among locations.

#### 781 5.4. Raster

782 Convolutional Neural Networks (CNNs) (Lecun and Bengio 1995) are at the core of  
783 many highly successful models in manipulating raster data such as image classification,  
784 image generation, and image understanding. This great success is due to the ability  
785 of the convolution operation to exploit the principles of locality, stationarity, and  
786 compositionality. Locality is due to the local connectivity, stationarity is owed to  
787 shift-invariance, and compositionality stems from the multi-resolution structure of the  
788 raster data (Bronstein *et al.* 2017). The number of learnable parameters is greatly  
789 reduced because of its feature locality and weight sharing across the data domain  
790 (Valsesia *et al.* 2019). Due to the success of location encoding on vector data, it is  
791 particularly interesting to think about the questions *how we can apply location encoding*  
792 *techniques on rasters and what the benefits are.*

793 Interestingly, with increasing popularity of the Transformer (Vaswani *et al.* 2017)  
794 architecture, several efforts have been made to replace CNN with a Transformer-like  
795 architecture for raster-based tasks. The idea is that instead of using CNN kernels, we  
796 first encode the pixel features as well as pixel locations into the embedding space with  
797 a location-encoder-like architecture, so-called *pixel position encoding*, and then a self-  
798 attention is applied on top of these pixel embeddings for different vision tasks. However,  
799 one problem with this approach is that this per-pixel based self-attention has a very high  
800 computational cost. One solution, among others, proposed by Vision Transformer (ViT)  
801 (Dosovitskiy *et al.* 2021) uses a per-image-patch (instead of per-pixel) self-attention  
802 which significantly lowers the computational cost. Dosovitskiy *et al.* (2021) showed  
803 that ViT can outperform traditional CNN-based models on several image classification  
804 benchmarks. However, the patch size becomes an important hyperparameter which will  
805 significantly affect model performance. Applying location encoders on raster data is a  
806 very new research direction. Existing work mainly focuses on encoding the positions  
807 of pixels on an image. When a pixel represents an area on the earth’s surface (e.g.,  
808 pixels in a satellite image), it is potentially very beneficial to encode pixel’s geographic  
809 locations rather than its image positions. The geo-locations can serve as a channel,  
810 which transfers knowledge learnt from large quantities of unlabeled data (geographic  
811 data, geo-tagged image, or geo-tagged text) to the supervised learning tasks.

## 812 6. Conclusion and Vision for Future Work

813 In this work, we formulate location encoding as an inductive learning based, task  
814 agnostic encoding technique for geographic locations. A formal definition of location  
815 encoding is provided, and two expected properties – distance preservation and direction  
816 awareness – are discussed from the perspective of GIScience. We illustrate the necessity  
817 of location encoding for GeoAI from a statistical machine learning perspective. A general  
818 classification framework has been provided to understand the current landscape of  
819 location encoding research (See Table 1). We classify the existing location encoders into  
820 two categories: single point location encoder  $Enc(\mathbf{x})$  and aggregation location encoder  
821  $Enc^{(\mathcal{P})}(\mathbf{x})$ . For each category, we unify the location encoders into the same formulation

822 framework (See Equations 1 and 7). Different location encoders are also compared  
823 based on various characteristics. Finally, we demonstrate the possible usage of location  
824 encoding for different types of spatial data.

825 There are several interesting future research directions of location encoding:

- 826 (1) **Region representation learning:** As we discussed in Section 5.2, there is  
827 no satisfactory solution for polygon encoding (so called region representation  
828 learning), which will be very useful in various tasks such as geographic entity  
829 alignment and topological relation reasoning. How to design a topology-aware  
830 polygon encoder which can handle simple polygons, polygon with holes, and  
831 multipolygons simultaneously is an interesting future research direction.
- 832 (2) **Spatiotemporal point encoding:** All the methods discussed so far are focused  
833 on location information whereas the temporal aspect of geospatial data is also very  
834 important. Several important related questions are: 1) How to utilize temporal  
835 information in GeoAI models? 2) Can we encode temporal information in a similar  
836 manner as spatial information? 3) What are the important properties we need to  
837 preserve when doing temporal encoding? 4) How to combine temporal encoding  
838 and location encoding in a single framework? As for event sequences that happen  
839 synchronously (Kazemi *et al.* 2019), i.e., sampled at regular intervals, the temporal  
840 information can be modeled implicitly by RNNs, or fed in RNNs as another input  
841 dimension after transforming time into handcrafted features (Du *et al.* 2016, Li  
842 *et al.* 2017, Rao *et al.* 2020). Instead of using handcrafted temporal features,  
843 recent work proposed to encode time as learnable vector representations such as  
844 Time2Vec (Kazemi *et al.* 2019) and Cai *et al.* (2020). These temporal encoders  
845 are expected to preserve important properties such as periodicity, temporal  
846 continuity, invariance to time rescaling, and so on. However, there are no systematic  
847 comparison studies among these temporal encoding approaches. As for combining  
848 location and temporal encoding, one obvious way is to add temporal information  
849 as an additional dimension of the location features. Mac Aodha *et al.* (2019)  
850 adopted this practice by adding time as an additional feature of  $PE_{wrap}(\mathbf{x})$  in  
851 Equation 2. This led to a small performance improvement (0.25%-1.37%).  
852 However, they failed to consider those important properties of time mentioned  
853 above. Future research is needed to study the pros and cons of different temporal  
854 encoding approaches and how to combine it with location encoding.
- 855 (3) **Spherical location encoding:** As we discussed in Section 5.4, currently, there  
856 are no existing location encoders which can preserve spherical surface distance.  
857 When we are dealing with large-scale geospatial data sets (e.g., global SST data,  
858 species occurrences all over the world) in which the map distortion problem is no  
859 longer negligible, a spherical-aware location encoder is required which enable us  
860 to directly *calculate on a round planet* (Chrisman 2017).
- 861 (4) **Unsupervised learning for location encoding:** Most of the location encoders  
862 listed in Table 1 are trained in a supervised learning fashion which prohibits  
863 the application of the trained location embedding on other tasks. In contrast,  
864 text encoding methods, e.g., BERT, are trained in an unsupervised manner from  
865 numerous unlabeled data, and the pretrained model can be utilized in different  
866 downstream tasks (Devlin *et al.* 2018). How to design an unsupervised learning  
867 framework for location encoding is a very attractive research direction. Recently,  
868 multiple point cloud generative models have been proposed such as r-GAN/l-GAN  
869 (Achlioptas *et al.* 2018), Graph-Conv GAN (Valsesia *et al.* 2019), tree-GAN (Shu  
870 *et al.* 2019), PointFlow (Yang *et al.* 2019), and Generative PointNet (Xie *et al.*



871 2021). Their objective is to reconstruct given point clouds. This presents one  
872 possible unsupervised learning framework of location encoding for unmarked points  
873 (points without attributes). Another interesting idea is unsupervised learning of  
874 the spatial distribution of marked points (points with attributes).

#### 875 **Data and Codes Availability Statement**

876 There is no code implementation for this review paper.

#### 877 **Acknowledgements**

878 We would like to thank Dr. Fei Du for his comments on the differences between location  
879 encoding and geohash. We also want to Thank Prof. Stefano Ermon for his suggestions  
880 on unsupervised learning for location encoding. We would like to thank the three  
881 anonymous reviewers for their thoughtful comments and suggestions.

882 This work is mainly funded by the National Science Foundation under Grant No.  
883 2033521 A1 – KnowWhereGraph: Enriching and Linking Cross-Domain Knowledge  
884 Graphs using Spatially-Explicit AI Technologies. Gengchen Mai acknowledges the  
885 support by UCSB Schmidt Summer Research Accelerator Award, Microsoft AI for  
886 Earth Grant, and IARPA SMART 2020-0072. Yingjie Hu acknowledges support by  
887 the National Science Foundation under Grant No. BCS-2117771 – Geospatial Artificial  
888 Intelligence Approaches for Understanding Location Descriptions in Natural Disasters  
889 and Their Spatial Biases. Any opinions, findings, and conclusions or recommendations  
890 expressed in this material are those of the authors and do not necessarily reflect the  
891 views of the National Science Foundation.

#### 892 **Notes on contributors**

893 **Gengchen Mai** is a Postdoctoral Scholar at the Stanford AI Lab, Department  
894 of Computer Science, Stanford University. He is also affiliated with the Stanford  
895 Sustainability and AI Lab. He obtained his Ph.D. degree in Geography from the  
896 Department of Geography, University of California, Santa Barbara in 2021 and  
897 B.S. degree in GIS from Wuhan University in 2015. His research mainly focuses on  
898 Spatially-Explicit Machine Learning, Geospatial Knowledge Graph, and Geographic  
899 Question Answering. The first draft of this manuscript was submitted when Gengchen  
900 was a Ph.D. student at the Space and Time for Knowledge Organization Lab,  
901 Department of Geography, UC Santa Barbara.

902 **Krzysztof Janowicz** is a Professor for Geoinformatics at the University of California,  
903 Santa Barbara and director of the Center for Spatial Studies. His research focuses on  
904 how humans conceptualize the space around them based on their behavior, focusing  
905 particularly on regional and cultural differences with the ultimate goal of assisting  
906 machines to better understand the information needs of an increasingly diverse user  
907 base. Janowicz’s expertise is in knowledge representation and reasoning as they apply  
908 to spatial and geographic data, e.g., in the form of knowledge graphs.

909 **Yingjie Hu** is an Assistant Professor in GIScience at the Department of Geography,  
910 University at Buffalo, The State University of New York. He holds a Ph.D. in Geography  
911 from the University of California, Santa Barbara. His main research interests include  
912 Geospatial Artificial Intelligence and Spatial Data Science.

913 **Song Gao** is an Assistant Professor in GIScience at the Department of Geography,  
914 University of Wisconsin-Madison. He holds a Ph.D. in Geography at the University

915 of California, Santa Barbara. His main research interests include Place-Based GIS,  
916 Geospatial Data Science, and GeoAI approaches to Human Mobility and Social Sensing.

917 **Bo Yan** received his bachelor’s degree in GIS from Wuhan University and his PhD  
918 in Geography from the University of California Santa Barbara. Dr Yan’s research area  
919 includes Geospatial Knowledge Graphs and machine learning in the geospatial domain.

920 **Rui Zhu** is a Postdoctoral Scholar at the Center for Spatial Studies, University of  
921 California, Santa Barbara. He obtained his Ph.D. in Geography from the same university  
922 in 2020; his M.S. degree was in Information Science from the University of Pittsburgh.  
923 Zhu’s research interests include geospatial semantics, spatial statistics, as well as their  
924 broader interactions in geospatial artificial intelligence (GeoAI).

925 **Ling Cai** is a fourth-year PhD student at the Space and Time for Knowledge  
926 Organization Lab, Department of Geography, University of California, Santa  
927 Barbara. She obtained her M.S. degree in Geographical Information Science from  
928 Chinese Academy of Sciences and B.S. degree from Wuhan University. Her research  
929 interests include qualitative spatial temporal reasoning, temporal knowledge graph,  
930 neuro-symbolic AI, and urban computing.

931 **Ni Lao** is currently a research scientist at Google. He holds a Ph.D. degree in  
932 Computer Science from the Language Technologies Institute, School of Computer  
933 Science at Carnegie Mellon University. He was the Chief scientist and co-founder of  
934 Mosaix.ai, a voice search AI startup. He is an expert in Machine Learning (ML),  
935 Knowledge Graph (KG) and Natural Language Understanding (NLU). He is known for  
936 his work on large scale inference and learning on KGs. This paper was done when Ni  
937 worked at Mosaix.ai.

## 938 References

- 939 Achlioptas, P., *et al.*, 2018. Learning representations and generative models for 3d point clouds.  
940 *In: International Conference on Machine Learning.* 40–49.
- 941 Agarwal, A., *et al.*, 2015. A boon in studies of cognitive functions: brain and grid cells. *Current*  
942 *Science*, 108 (12), 2142–2144.
- 943 Andris, C., 2016. Integrating social network data into gisystems. *International Journal of*  
944 *Geographical Information Science*, 30 (10), 2009–2031.
- 945 Appleby, G., Liu, L., and Liu, L.P., 2020. Kriging convolutional networks. *In: AAAI 2020.*
- 946 Battaglia, P.W., *et al.*, 2018. Relational inductive biases, deep learning, and graph networks.  
947 *arXiv preprint arXiv:1806.01261.*
- 948 Berg, T., *et al.*, 2014. Birdsnap: Large-scale fine-grained visual categorization of birds. *In:*  
949 *CVPR 2014.* 2011–2018.
- 950 Blair, H.T., Welday, A.C., and Zhang, K., 2007. Scale-invariant memory representations emerge  
951 from moire interference between grid fields that produce theta oscillations: a computational  
952 model. *Journal of Neuroscience*, 27 (12), 3211–3229.
- 953 Bordes, A., *et al.*, 2013. Translating embeddings for modeling multi-relational data. *In: Advances*  
954 *in neural information processing systems.* 2787–2795.
- 955 Bronstein, M.M., *et al.*, 2017. Geometric deep learning: going beyond euclidean data. *IEEE*  
956 *Signal Processing Magazine*, 34 (4), 18–42.
- 957 Cai, L., *et al.*, 2020. Traffic transformer: Capturing the continuity and periodicity of time series  
958 for traffic forecasting. *Transactions in GIS.*
- 959 Cai, L., *et al.*, 2019. TransGCN: Coupling transformation assumptions with graph convolutional  
960 networks for link prediction. *In: ACM K-CAP 2019.* 131–138.
- 961 Cho, K., *et al.*, 2014. Learning phrase representations using rnn encoder–decoder for statistical  
962 machine translation. *In: EMNLP 2014.* 1724–1734.

- 963 Chrisman, N.R., 2017. Calculating on a round planet. *International Journal of Geographical*  
964 *Information Science*, 31 (4), 637–657.
- 965 Chu, G., *et al.*, 2019. Geo-aware networks for fine-grained recognition. *In: Proceedings of the*  
966 *IEEE International Conference on Computer Vision Workshops*. 0–0.
- 967 Couclelis, H., 1986. Artificial intelligence in geography: Conjectures on the shape of things to  
968 come. *The professional geographer*, 38 (1), 1–11.
- 969 DeLozier, G., Baldrige, J., and London, L., 2015. Gazetteer-independent toponym resolution  
970 using geographic word profiles. *In: AAAI 2015*. 2382–2388.
- 971 Devlin, J., *et al.*, 2018. Bert: Pre-training of deep bidirectional transformers for language  
972 understanding. *arXiv preprint arXiv:1810.04805*.
- 973 Dosovitskiy, A., *et al.*, 2021. An image is worth 16x16 words: Transformers for image recognition  
974 at scale. *In: ICLR 2021*.
- 975 Du, N., *et al.*, 2016. Recurrent marked temporal point processes: Embedding event history to  
976 vector. *In: ACM SIGKDD 2016*. 1555–1564.
- 977 Feng, S., *et al.*, 2017. POI2Vec: Geographical latent representation for predicting future visitors.  
978 *In: AAAI 2017*.
- 979 Fotheringham, A.S. and Wong, D.W., 1991. The modifiable areal unit problem in multivariate  
980 statistical analysis. *Environment and planning A*, 23 (7), 1025–1044.
- 981 Gao, R., *et al.*, 2019. Learning grid cells as vector representation of self-position coupled with  
982 matrix representation of self-motion. *In: ICLR 2019*.
- 983 Gao, R., *et al.*, 2020. A representational model of grid cells based on matrix lie algebras. *arXiv*  
984 *preprint arXiv:2006.10259*.
- 985 Gilmer, J., *et al.*, 2017. Neural message passing for quantum chemistry. *In: ICML 2017*.
- 986 Goodfellow, I., *et al.*, 2014. Generative adversarial nets. *In: Advances in Neural Information*  
987 *Processing Systems*. 2672–2680.
- 988 Grover, A. and Leskovec, J., 2016. node2vec: Scalable feature learning for networks. *In: ACM*  
989 *SIGKDD 2016*. 855–864.
- 990 Hafting, T., *et al.*, 2005. Microstructure of a spatial map in the entorhinal cortex. *Nature*, 436  
991 (7052), 801–806.
- 992 Hamilton, W., Ying, Z., and Leskovec, J., 2017. Inductive representation learning on large  
993 graphs. *In: Advances in neural information processing systems*. 1024–1034.
- 994 Hastie, T., Tibshirani, R., and Friedman, J., 2009. *The elements of statistical learning: Data*  
995 *mining, inference, and prediction*. Springer.
- 996 Hochreiter, S. and Schmidhuber, J., 1997. Long short-term memory. *Neural computation*, 9 (8),  
997 1735–1780.
- 998 Izbicki, M., Papalexakis, E.E., and Tsotras, V.J., 2019. Exploiting the earth’s spherical geometry  
999 to geolocate images. *In: Joint European Conference on Machine Learning and Knowledge*  
1000 *Discovery in Databases*. Springer, 3–19.
- 1001 Janowicz, K., *et al.*, 2020. GeoAI: Spatially explicit artificial intelligence techniques for  
1002 geographic knowledge discovery and beyond. *International Journal of Geographic Information*  
1003 *Science*.
- 1004 Kazemi, S.M., *et al.*, 2019. Time2vec: Learning a vector representation of time. *arXiv preprint*  
1005 *arXiv:1907.05321*.
- 1006 Kejriwal, M. and Szekely, P., 2017. Neural embeddings for populated geonames locations. *In:*  
1007 *International Semantic Web Conference*. Springer, 139–146.
- 1008 Killian, N.J., Jutras, M.J., and Buffalo, E.A., 2012. A map of visual space in the primate  
1009 entorhinal cortex. *Nature*, 491 (7426), 761–764.
- 1010 Kipf, T.N. and Welling, M., 2017. Semi-supervised classification with graph convolutional  
1011 networks. *In: ICLR 2017*.
- 1012 Krizhevsky, A., Sutskever, I., and Hinton, G.E., 2012. Imagenet classification with deep  
1013 convolutional neural networks. *In: Advances in neural information processing systems*.  
1014 1097–1105.
- 1015 Kulkarni, S., *et al.*, 2020. Spatial language representation with multi-level geocoding. *arXiv*  
1016 *preprint arXiv:2008.09236*.

- 1017 Lechner, A.M., *et al.*, 2012. Investigating species–environment relationships at multiple scales:  
1018 Differentiating between intrinsic scale and the modifiable areal unit problem. *Ecological*  
1019 *Complexity*, 11, 91–102.
- 1020 Lecun, Y. and Bengio, Y., 1995. *Convolutional networks for images, speech and time series*.  
1021 The MIT Press, 255–258.
- 1022 Li, W. and Hsu, C.Y., 2020. Automated terrain feature identification from remote sensing  
1023 imagery: a deep learning approach. *International Journal of Geographical Information Science*,  
1024 34 (4), 637–660.
- 1025 Li, Y., *et al.*, 2018a. Diffusion convolutional recurrent neural network: Data-driven traffic  
1026 forecasting. *In: ICLR 2018*.
- 1027 Li, Y., Du, N., and Bengio, S., 2017. Time-dependent representation for neural event sequence  
1028 prediction. *arXiv preprint arXiv:1708.00065*.
- 1029 Li, Y., *et al.*, 2018b. PointCNN: Convolution on x-transformed points. *In: Advances in neural*  
1030 *information processing systems*. 820–830.
- 1031 Li, Y., *et al.*, 2016. Gated graph sequence neural networks. *In: ICLR 2016*.
- 1032 Lin, Y., *et al.*, 2018. Exploiting spatiotemporal patterns for accurate air quality forecasting  
1033 using deep learning. *In: ACM SIGSPATIAL 2018*. 359–368.
- 1034 Mac Aodha, O., Cole, E., and Perona, P., 2019. Presence-only geographical priors for fine-grained  
1035 image classification. *In: ICCV 2019*. 9596–9606.
- 1036 Mai, G., *et al.*, 2020a. SE-KGE: A location-aware knowledge graph embedding model for  
1037 geographic question answering and spatial semantic lifting. *Transactions in GIS*, 24, 623–655.
- 1038 Mai, G., *et al.*, 2018. ADCN: An anisotropic density-based clustering algorithm for discovering  
1039 spatial point patterns with noise. *Transactions in GIS*, 22 (1), 348–369.
- 1040 Mai, G., *et al.*, 2019a. Contextual graph attention for answering logical queries over incomplete  
1041 knowledge graphs. *In: ACM K-CAP 2019*. 171–178.
- 1042 Mai, G., *et al.*, 2020b. Multi-scale representation learning for spatial feature distributions using  
1043 grid cells. *In: ICLR 2020*.
- 1044 Mai, G., *et al.*, 2021. Geographic Question Answering: Challenges, Uniqueness, Classification,  
1045 and Future Directions. *AGILE 2021: GIScience Series*, 2, 1–21.
- 1046 Mai, G., *et al.*, 2019b. Relaxing unanswerable geographic questions using a spatially explicit  
1047 knowledge graph embedding model. *In: AGILE 2019*. Springer, 21–39.
- 1048 Maturana, D. and Scherer, S., 2015. Voxnet: A 3d convolutional neural network for real-time  
1049 object recognition. *In: 2015 IEEE/RSJ International Conference on Intelligent Robots and*  
1050 *Systems (IROS)*. IEEE, 922–928.
- 1051 Mirowski, P., *et al.*, 2018. Learning to navigate in cities without a map. *In: Advances in Neural*  
1052 *Information Processing Systems*. 2419–2430.
- 1053 Noh, H., Hong, S., and Han, B., 2015. Learning deconvolution network for semantic segmentation.  
1054 *In: Proceedings of the IEEE international conference on computer vision*. 1520–1528.
- 1055 Openshaw, S., 1981. The modifiable areal unit problem. *Quantitative geography: A British view*,  
1056 60–69.
- 1057 Openshaw, S. and Openshaw, C., 1997. *Artificial intelligence in geography*. John Wiley & Sons,  
1058 Inc.
- 1059 Qi, C.R., *et al.*, 2017a. PointNet: Deep learning on point sets for 3d classification and  
1060 segmentation. *In: CVPR 2017*. 652–660.
- 1061 Qi, C.R., *et al.*, 2016. Volumetric and multi-view cnns for object classification on 3d data. *In:*  
1062 *CVPR 2016*. 5648–5656.
- 1063 Qi, C.R., *et al.*, 2017b. PointNet++: Deep hierarchical feature learning on point sets in a metric  
1064 space. *In: Advances in neural information processing systems*. 5099–5108.
- 1065 Qiu, P., *et al.*, 2019. Knowledge embedding with geospatial distance restriction for geographic  
1066 knowledge graph completion. *ISPRS International Journal of Geo-Information*, 8 (6), 254.
- 1067 Rao, J., *et al.*, 2020. LSTM-TrajGAN: A deep learning approach to trajectory privacy protection.  
1068 *In: GIScience 2020*. 12:1–12:17.
- 1069 Regalia, B., Janowicz, K., and McKenzie, G., 2019. Computing and querying strict, approximate,  
1070 and metrically refined topological relations in linked geographic data. *Transactions in GIS*,

1071 23 (3), 601–619.

1072 Ronneberger, O., Fischer, P., and Brox, T., 2015. U-net: Convolutional networks for biomedical  
1073 image segmentation. *In: International Conference on Medical image computing and computer-*  
1074 *assisted intervention*. Springer, 234–241.

1075 Russell, S. and Norvig, P., 2015. *Artificial Intelligence: a modern approach*. Pearson.

1076 Schlichtkrull, M., *et al.*, 2018. Modeling relational data with graph convolutional networks. *In:*  
1077 *European Semantic Web Conference*. Springer, 593–607.

1078 Shu, D.W., Park, S.W., and Kwon, J., 2019. 3d point cloud generative adversarial network  
1079 based on tree structured graph convolutions. *In: ICCV 2019*. 3859–3868.

1080 Smith, T.R., 1984. Artificial intelligence and its applicability to geographical problem solving.  
1081 *The Professional Geographer*, 36 (2), 147–158.

1082 Su, H., *et al.*, 2015. Multi-view convolutional neural networks for 3d shape recognition. *In:*  
1083 *ICCV 2015*. 945–953.

1084 Tang, K., *et al.*, 2015. Improving image classification with location context. *In: ICCV 2015*.  
1085 1008–1016.

1086 Te Stroet, C.B. and Snepvangers, J.J., 2005. Mapping curvilinear structures with local anisotropy  
1087 kriging. *Mathematical geology*, 37 (6), 635–649.

1088 Tobler, W.R., 1970. A computer movie simulating urban growth in the detroit region. *Economic*  
1089 *geography*, 46 (sup1), 234–240.

1090 Trisedya, B.D., Qi, J., and Zhang, R., 2019. Entity alignment between knowledge graphs using  
1091 attribute embeddings. *In: AAAI 2019*. vol. 33, 297–304.

1092 Trouillon, T., *et al.*, 2017. Knowledge graph completion via complex tensor factorization. *The*  
1093 *Journal of Machine Learning Research*, 18 (1), 4735–4772.

1094 Valsesia, D., Fracastoro, G., and Magli, E., 2019. Learning localized generative models for 3d  
1095 point clouds via graph convolution. *In: ICLR 2019*.

1096 Vapnik, V., 2013. *The nature of statistical learning theory*. Springer science & business media.

1097 Vaswani, A., *et al.*, 2017. Attention is all you need. *In: Advances in neural information processing*  
1098 *systems*. 5998–6008.

1099 Veličković, P., *et al.*, 2018. Graph attention networks. *In: ICLR 2018*.

1100 Wang, J., Hu, Y., and Joseph, K., 2020. Neurotrp: A neuro-net toponym recognition model for  
1101 extracting locations from social media messages. *Transactions in GIS*.

1102 Wang, Y., *et al.*, 2019. Dynamic Graph CNN for learning on point clouds. *Acm Transactions*  
1103 *On Graphics (tog)*, 38 (5), 1–12.

1104 Weyand, T., Kostrikov, I., and Philbin, J., 2016. Planet-photo geolocation with convolutional  
1105 neural networks. *In: European Conference on Computer Vision*. Springer, 37–55.

1106 Wu, Y., *et al.*, 2020a. Inductive graph neural networks for spatiotemporal kriging. *arXiv preprint*  
1107 *arXiv:2006.07527*.

1108 Wu, Z., *et al.*, 2020b. A comprehensive survey on graph neural networks. *IEEE Transactions*  
1109 *on Neural Networks and Learning Systems*.

1110 Xie, J., *et al.*, 2021. Generative pointnet: Deep energy-based learning on unordered point sets  
1111 for 3d generation, reconstruction and classification. *In: CVPR 2021*. 14976–14985.

1112 Xu, K., *et al.*, 2019. How powerful are graph neural networks? *In: ICLR 2019*.

1113 Xu, Y., Piao, Z., and Gao, S., 2018. Encoding crowd interaction with deep neural network for  
1114 pedestrian trajectory prediction. *In: CVPR 2018*. 5275–5284.

1115 Yan, B., *et al.*, 2019. A spatially explicit reinforcement learning model for geographic knowledge  
1116 graph summarization. *Transactions in GIS*, 23 (3), 620–640.

1117 Yan, X., *et al.*, 2021. Graph convolutional autoencoder model for the shape coding and cognition  
1118 of buildings in maps. *International Journal of Geographical Information Science*, 35 (3),  
1119 490–512.

1120 Yang, G., *et al.*, 2019. Pointflow: 3d point cloud generation with continuous normalizing flows.  
1121 *In: ICCV 2019*. 4541–4550.

1122 Yin, Y., *et al.*, 2019. GPS2Vec: Towards generating worldwide gps embeddings. *In: ACM*  
1123 *SIGSPATIAL 2019*. 416–419.

1124 Zaheer, M., *et al.*, 2017. Deep sets. *In: Advances in neural information processing systems*.

1125 3391–3401.  
1126 Zhong, E.D., *et al.*, 2020. Reconstructing continuous distributions of 3d protein structure from  
1127 cryo-em images. *In: ICLR 2020*.  
1128 Zhong, L., Hu, L., and Zhou, H., 2019. Deep learning based multi-temporal crop classification.  
1129 *Remote sensing of environment*, 221, 430–443.  
1130 Zhou, Y. and Tuzel, O., 2018. Voxelnet: End-to-end learning for point cloud based 3d object  
1131 detection. *In: CVPR 2018*. 4490–4499.  
1132 Zhu, D., *et al.*, 2020. Understanding place characteristics in geographic contexts through graph  
1133 convolutional neural networks. *Annals of the American Association of Geographers*, 110 (2),  
1134 408–420.  
1135 Zhu, M., *et al.*, 2019a. Location2Vec: a situation-aware representation for visual exploration  
1136 of urban locations. *IEEE Transactions on Intelligent Transportation Systems*, 20 (10),  
1137 3981–3990.  
1138 Zhu, R., Janowicz, K., and Mai, G., 2019b. Making direction a first-class citizen of tobler’s first  
1139 law of geography. *Transactions in GIS*, 23 (3), 398–416.



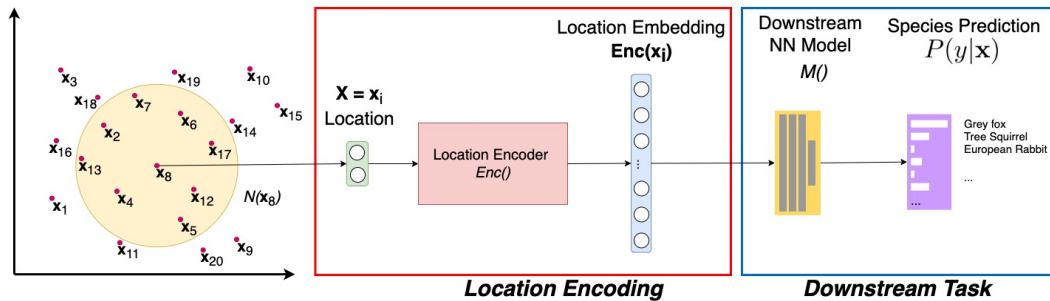


Figure 1.: An illustration of location encoding. Here, we use *location-based species classification* as an example of the downstream tasks. Those 20 points in 2D space represent species occurrence records. Each occurrence can be written as  $p_i = (\mathbf{x}_i, y_i)$  where  $\mathbf{x}_i$  indicates the 2D locations and  $y_i$  indicates the corresponding species type, i.e., the ground truth label.  $\mathcal{N}(\mathbf{x}_i)$  indicates the spatial neighborhood of  $\mathbf{x}_i$ . A location encoder  $Enc(\cdot)$  takes 2D location  $\mathbf{x}_i$  as its input and outputs a location embedding as a high dimensional vector. This embedding is further fed into a downstream NN model  $M()$  for species prediction. The whole model architecture can be trained end-to-end in a supervised learning manner.



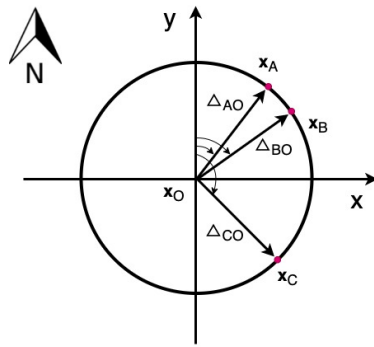


Figure 2.: An illustration of the direction preservation property of location encoding.

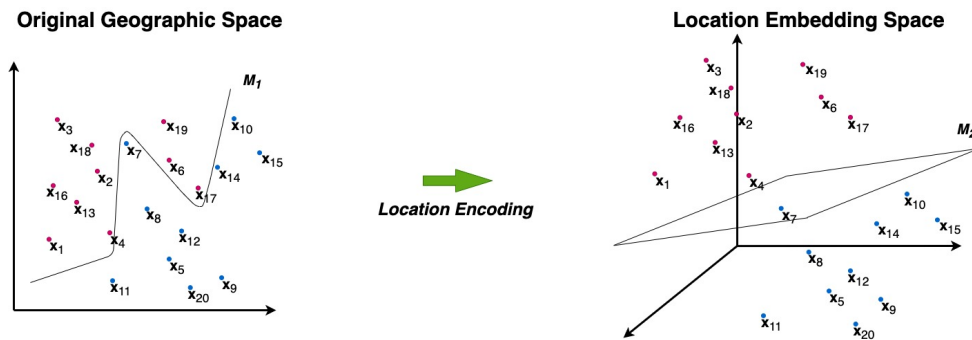


Figure 3.: An illustration of how location encoding can help to produce learning-friendly representations of geographic locations for downstream models. We use the same 20 points in Figure 1 as an example of  $\mathcal{P} = \{p_i\}$ . The red and blue points indicate they belong to two different classes.  $M_1$  and  $M_2$  are the illustrations of the trained binary classifiers in the original geographic space and the location embedding space.

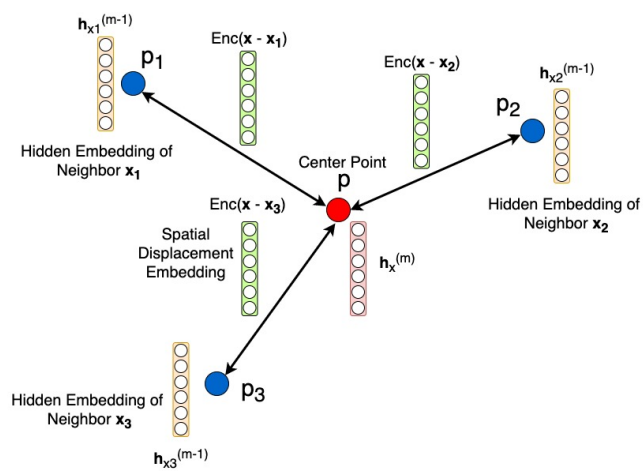


Figure 4.: An illustration of the  $m$ th aggregation layer of SAGAT.

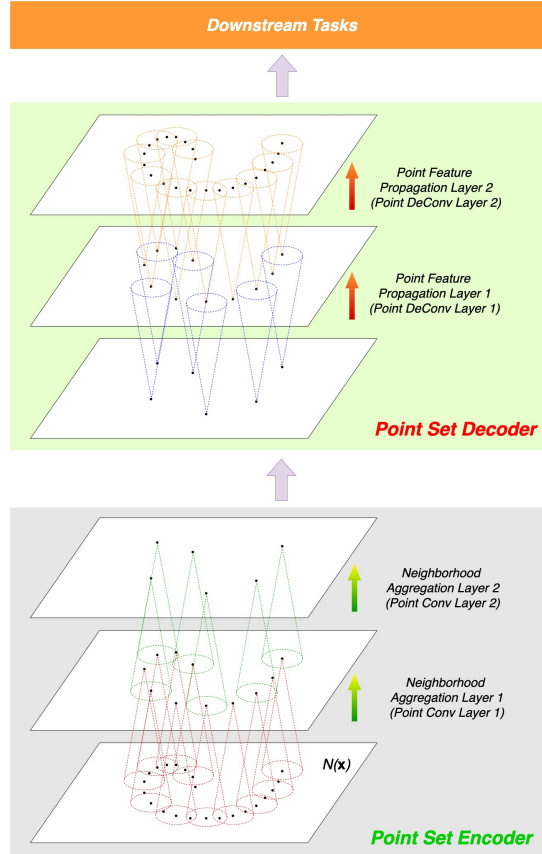


Figure 5.: An illustration of the Conv-DeConv architecture (Noh *et al.* 2015) for  $Enc_{hierarchical}^{(\mathcal{P})}(\mathbf{x})$  such as PointNet++ (Qi *et al.* 2017b), PointCNN (Li *et al.* 2018b), and Graph-Conv GAN (Valsesia *et al.* 2019).

1141 **List of Tables**

1142 1 Overview of location encoding approaches. Single point location encoders  
1143  $Enc(\mathbf{x})$  and aggregation location encoders  $Enc^{(\mathcal{P})}(\mathbf{x})$  are further  
1144 classified based on either  $PE(\mathbf{x})$  or  $\mathcal{N}(\mathbf{x})$  (see Figure 1). (M) indicates  
1145 multi-scale representation. \* indicates a generalized version of the original  
1146 model cited. We consider multiple criteria of location encoders: 1)  $L$ :  
1147 The spatial dimension of  $\mathcal{P}$ ; 2) Parametric: Is the location encoder a  
1148 parametric model (Yes) or non-parametric model (No)? 3) Mul.S.: Does  
1149 the location encoding adopt a multi-scale approach? 4) Dist.P.: Does  
1150 this location encoder preserve distance (Property 2.1)? 5) Dir.A.: Is this  
1151 location encoder aware of direction (Property 2.2)? For Dist.P. and Dir.A.  
1152 “Yes” or “No” indicates whether the property can be proved empirically  
1153 (for example by using the response maps of trained location encoders  
1154 Mai *et al.* (2020b)). “-” indicates that the property is unknown. “Yes+”  
1155 indicates that the property was shown both theoretically and empirically. 38

Table 1.: Overview of location encoding approaches. Single point location encoders  $Enc(\mathbf{x})$  and aggregation location encoders  $Enc^{(\mathcal{P})}(\mathbf{x})$  are further classified based on either  $PE(\mathbf{x})$  or  $\mathcal{N}(\mathbf{x})$  (see Figure 1). (M) indicates multi-scale representation. \* indicates a generalized version of the original model cited. We consider multiple criteria of location encoders: 1)  $L$ : The spatial dimension of  $\mathcal{P}$ ; 2) Parametric: Is the location encoder a parametric model (Yes) or non-parametric model (No)? 3) Mul.S.: Does the location encoding adopt a multi-scale approach? 4) Dist.P.: Does this location encoder preserve distance (Property 2.1)? 5) Dir.A.: Is this location encoder aware of direction (Property 2.2)? For Dist.P. and Dir.A. “Yes” or “No” indicates whether the property can be proved empirically (for example by using the response maps of trained location encoders Mai *et al.* (2020b)). “-” indicates that the property is unknown. “Yes+” indicates that the property was shown both theoretically and empirically.

Encoder	$PE(\mathbf{x})$	Model	$L$	Parametric	Mul.S.	Dist.P.	Dir.A.
$Enc(\mathbf{x})$	Discretization	onehot(Tang <i>et al.</i> 2015)	2,3	Yes	No	No	-
		tile(Mai <i>et al.</i> 2020b)	2,3	Yes	No	No	-
	Direct	direct(Xu <i>et al.</i> 2018, Chu <i>et al.</i> 2019, Rao <i>et al.</i> 2020)	2,3	Yes	No	Yes	Yes
	Sinusoidal	wrap(Mac Aodha <i>et al.</i> 2019)	2	Yes	No	Yes	-
	Sinusoidal (M)	TF(Zhong <i>et al.</i> 2020)	2,3	Yes	Yes	-	-
theory/Space2Vec (Mai <i>et al.</i> 2020b) grid/Space2Vec (Mai <i>et al.</i> 2020b)		2 2,3	Yes Yes	Yes Yes	Yes+ Yes	Yes -	
Encoder	$\mathcal{N}(\mathbf{x})$	Model	$L$	Parametric	Mul.S.	Dist.Pr.	Dir.A.
$Enc^{(\mathcal{P})}(\mathbf{x})$	Kernel	GPS2Vec(Yin <i>et al.</i> 2019)	2	Yes	No	-	-
		rbf(Mai <i>et al.</i> 2020b)	2,3	Yes/No	No	Yes	No
		Adapted kernel* (Berg <i>et al.</i> 2014)	2,3	Yes/No	No	-	-
	Global	PointNet (Qi <i>et al.</i> 2017a)	2,3	Yes	No	-	-
	Local	VoxelNet (Zhou and Tuzel 2018)	2,3	Yes	No	-	-
		SAGAT (Mai <i>et al.</i> 2020b)	2,3	Yes/No	Yes/No	Yes/No	Yes/No
		DGCNN (Wang <i>et al.</i> 2019)	2,3	Yes	No	-	-
	Hierarchical	PointNet++ (Qi <i>et al.</i> 2017b)	2,3	Yes	Yes	-	-
		PointCNN (Li <i>et al.</i> 2018b)	2,3	Yes	Yes	-	-
GraphCNN (Valsesia <i>et al.</i> 2019)		2,3	Yes	Yes	-	-	