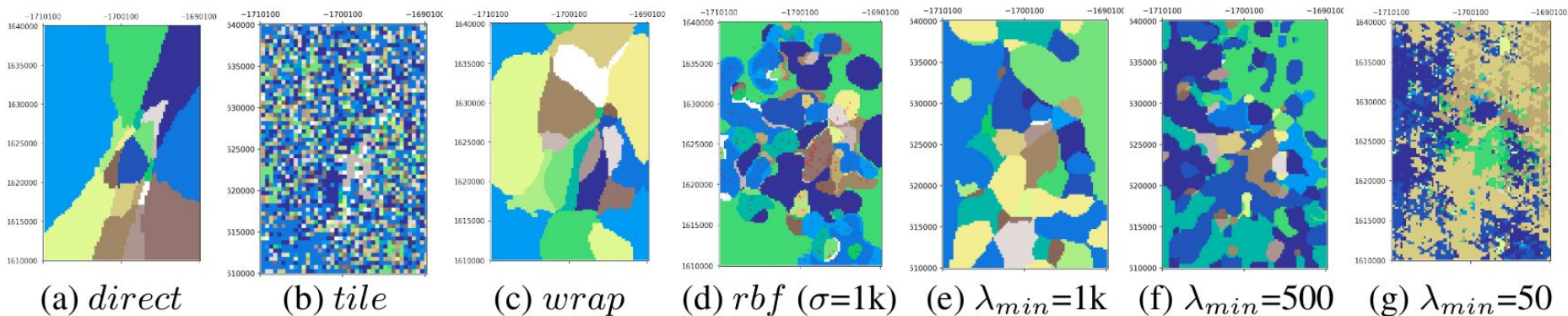


# Multi-Scale Representation Learning for Spatial Feature Distributions using Grid Cells

Gengchen Mai<sup>1</sup>, Krzysztof Janowicz<sup>1</sup>, Bo Yan<sup>2</sup>, Rui Zhu<sup>1</sup>, Ling Cai<sup>1</sup>, Ni Lao<sup>3</sup>  
<sup>1</sup>STKO Lab, UC Santa Barbara; <sup>2</sup> LinkedIn Corporation; <sup>3</sup> SayMosaic Inc.

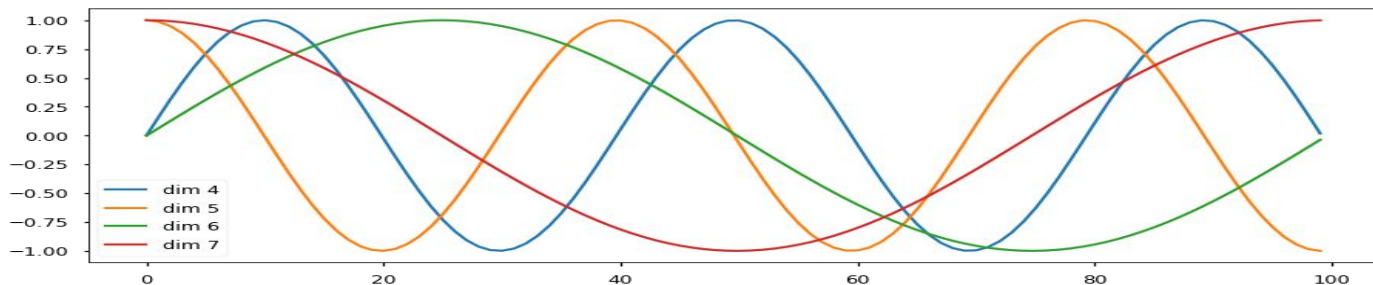
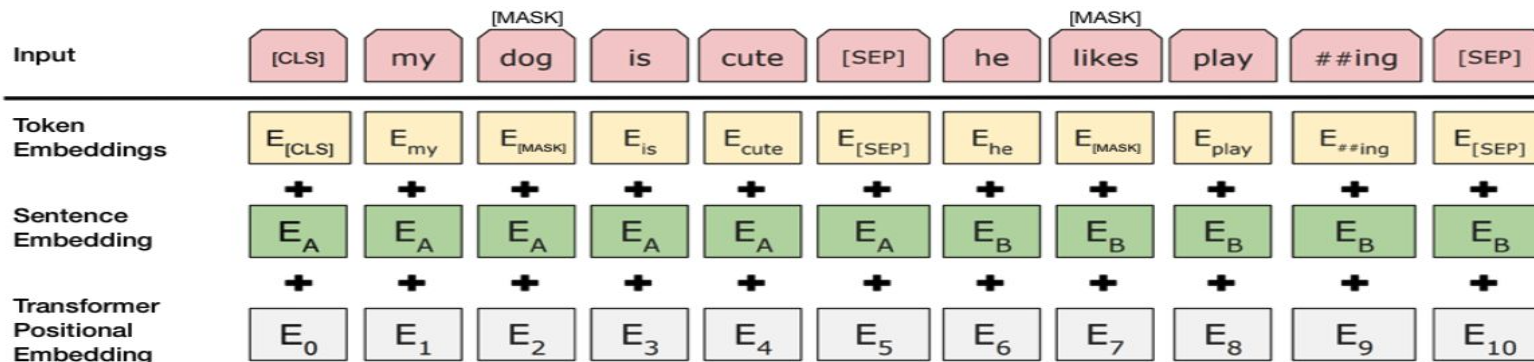


Embedding clustering of different location encoding models:

(a)-(d) baselines (e)-(f) **Space2Vec**

# Unsupervised Text Encoding

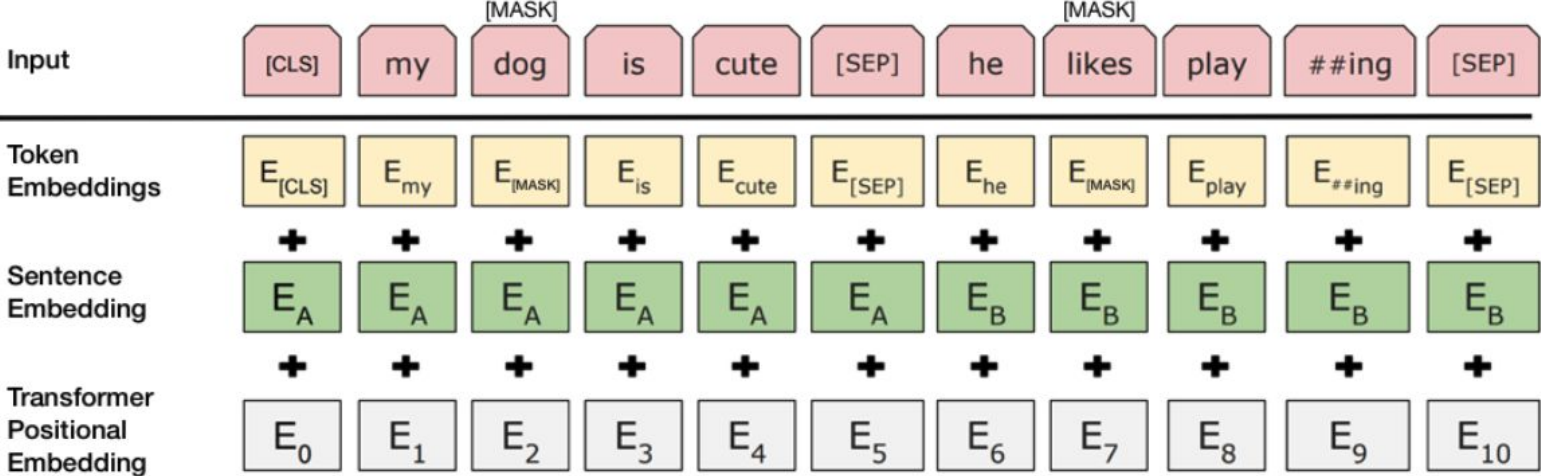
Position Encoding: encode word positions with sinusoid functions of different frequencies



Transformer (Vaswani et al., 2017) BERT (Devlin et al., 2019)

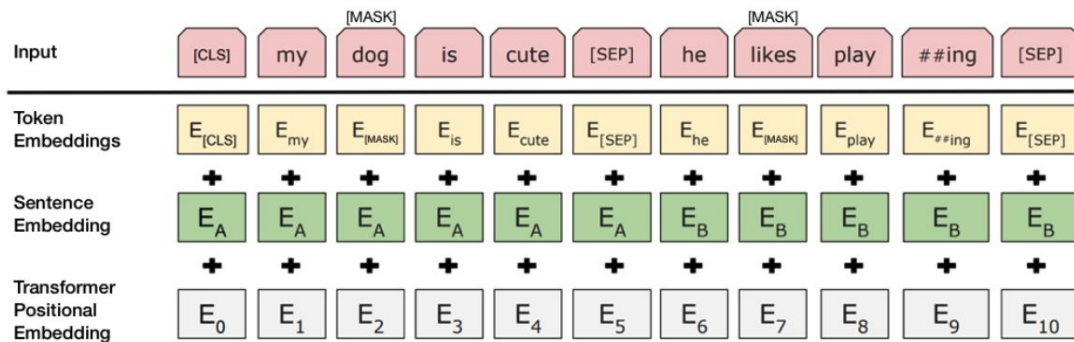
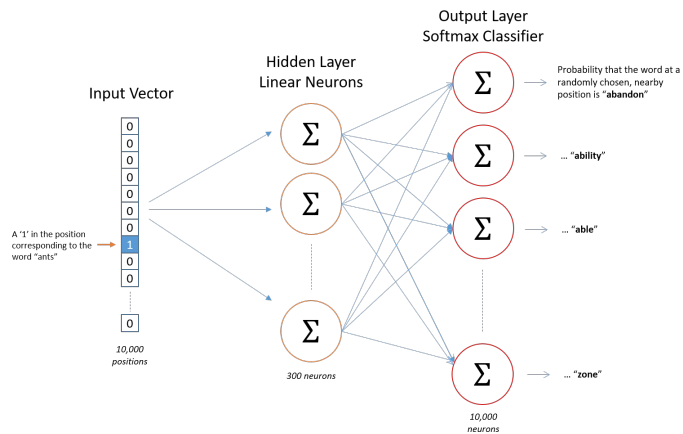
# Unsupervised Text Encoding

Position Encoding: encode word positions in sentences with multiple sinusoid functions with different frequencies



BERT (Devlin et al. 2019)

# Unsupervised Text Encoding



Word2Vec (Mikolov et al., 2013)<sup>1</sup>

BERT (Devlin et al. 2019)

<sup>1</sup><http://mccormickml.com/2016/04/19/word2vec-tutorial-the-skip-gram-model/>

# Unsupervised Location Encoding

## 1. Radial Basis Function (RBF)

$$K(\mathbf{x}, \mathbf{x}') = \exp\left(-\frac{\|\mathbf{x} - \mathbf{x}'\|^2}{2\sigma^2}\right)$$

- choosing the correct scale is challenging
- Need to memorize the training samples

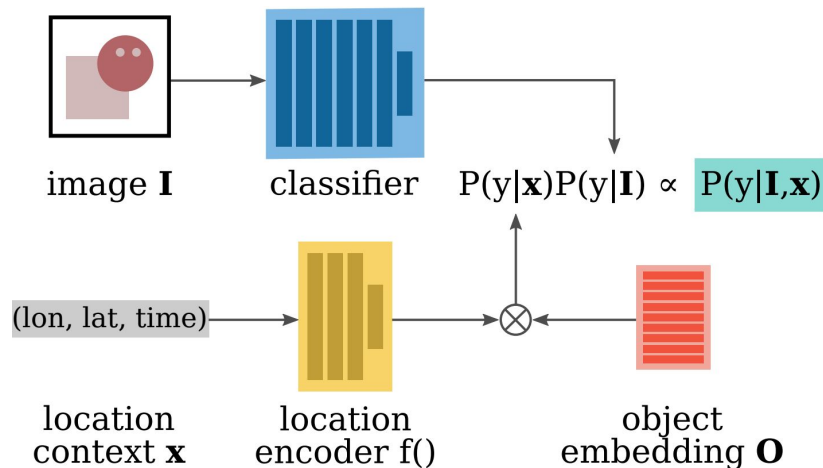
## 2. Tile-based approaches (Berg et al. 2014):

discretize the study area into regular grids

- choosing the correct scale is challenging
- does not scale well in terms of memory

## 3. Directly feed the coordinates into a FFN (inductive single-scale location encoder)

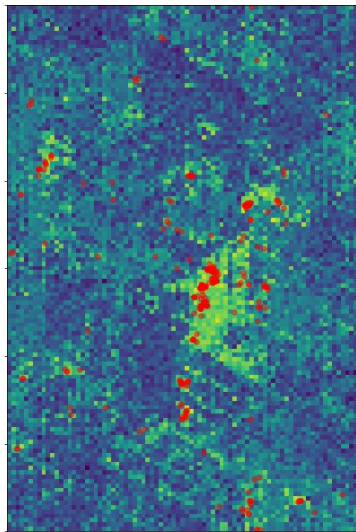
- hard to capture fine grained distributions



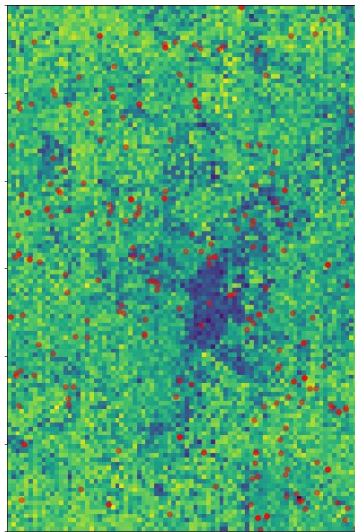
Geo-aware Image Classification (Mac Aodha et al., 2019)

# Key challenge for location encoding

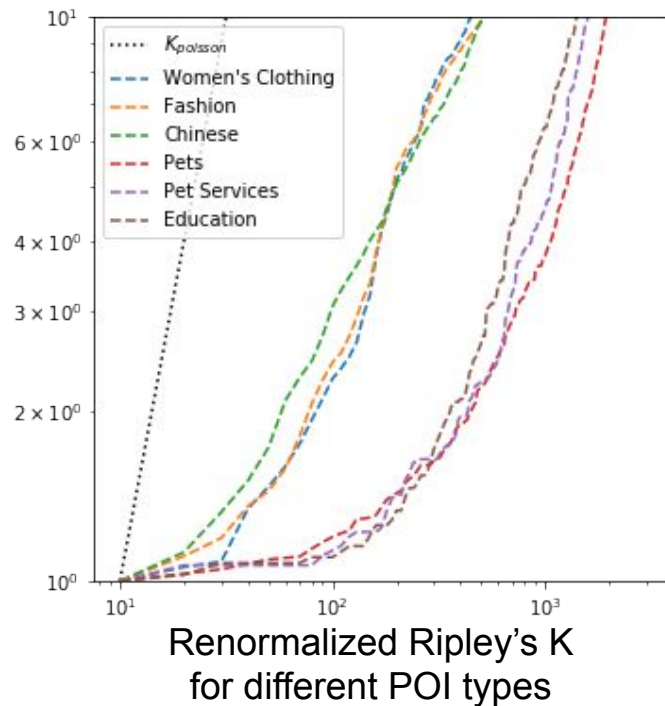
- Joint modeling distributions with very different characteristics
- => **multi-scale location representations**



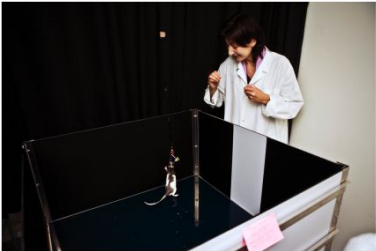
Women's Clothing  
(Clustered Distribution)



Education  
(Even Distribution)



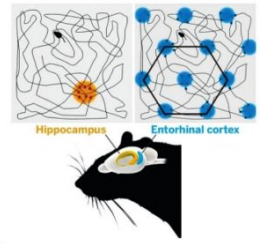
# Grid Cell Based Multi-Scale Location Encoding



(a)



(b)

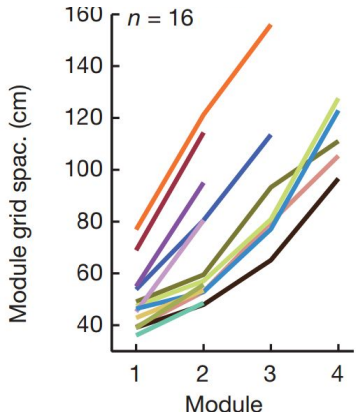


(c)



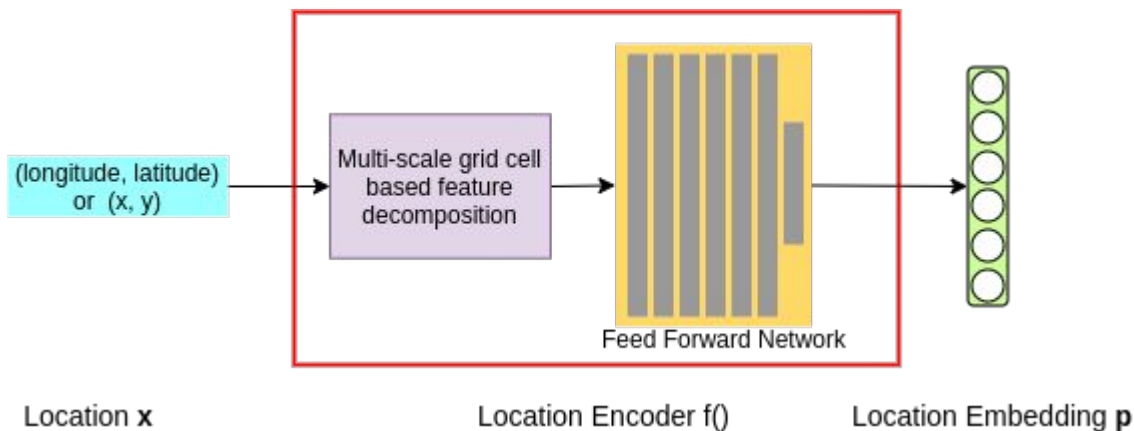
(d)

- **Grid cells** in mammals provide a **multi-scale periodic representation** that functions as a metric for location encoding.
- It can be simulated by summing **three cosine grating functions** oriented 60 degree apart (a **simple Fourier model of the hexagonal lattice**).



Mean grid spacing for all modules (M1–M4) in all animals (colour-coded)

# Space2Vec



**Given a location  $\mathbf{x}$ :**

$$Enc_{theory}^{(x)}(\mathbf{x}) = \text{NN}(PE^{(t)}(\mathbf{x}))$$

$$PE^{(t)}(\mathbf{x}) = [PE_0^{(t)}(\mathbf{x}); \dots; PE_s^{(t)}(\mathbf{x}); \dots; PE_{S-1}^{(t)}(\mathbf{x})]$$

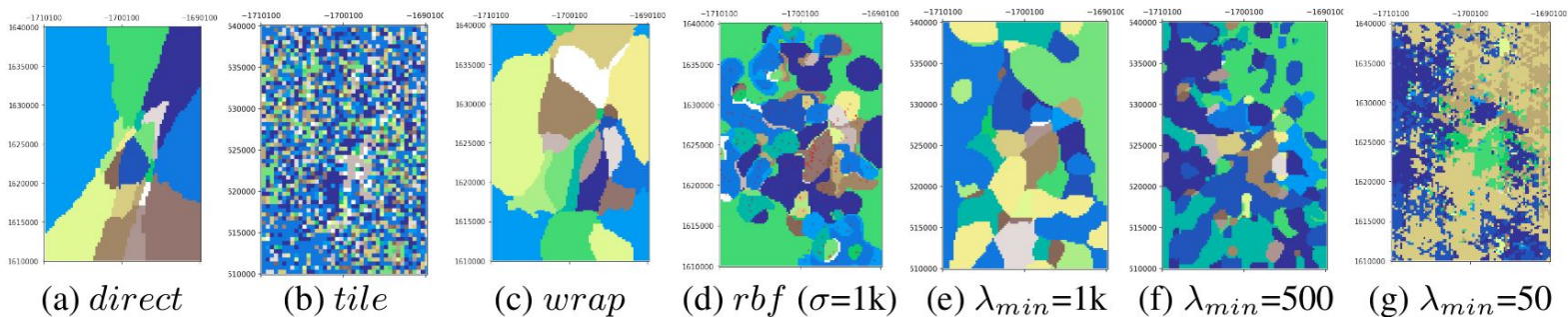
$$PE_s^{(t)}(\mathbf{x}) = [PE_{s,1}^{(t)}(\mathbf{x}); PE_{s,2}^{(t)}(\mathbf{x}); PE_{s,3}^{(t)}(\mathbf{x})]$$

$$PE_{s,j}^{(t)}(\mathbf{x}) = \left[ \cos\left(\frac{\langle \mathbf{x}, \mathbf{a}_j \rangle}{\lambda_{min} \cdot g^{s/(S-1)}}\right); \sin\left(\frac{\langle \mathbf{x}, \mathbf{a}_j \rangle}{\lambda_{min} \cdot g^{s/(S-1)}}\right) \right] \forall j = 1, 2, 3;$$

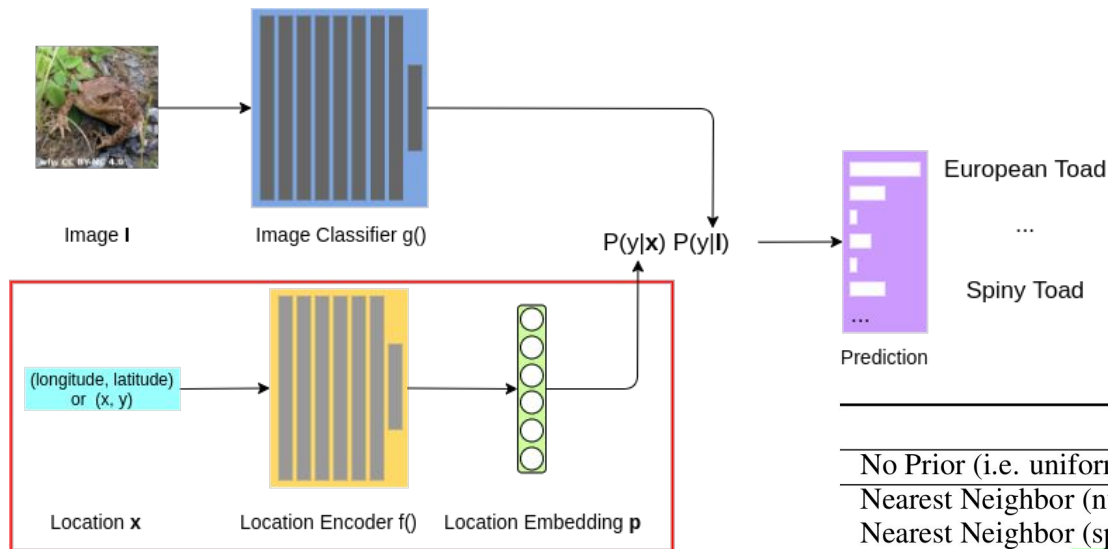


# Point of Interest Type Classification

POI Groups	Clustered ( $r \leq 100m$ )	Middle ( $100m < r < 200m$ )	Even ( $r \geq 200m$ )
<i>direct</i>	0.080 (-0.047)	0.108 (-0.030)	0.084 (-0.047)
<i>wrap</i>	0.106 (-0.021)	0.126 (-0.012)	0.122 (-0.009)
<i>tile</i>	0.108 (-0.019)	0.135 (-0.003)	0.111 (-0.020)
<i>rbf</i>	0.112 (-0.015)	0.136 (-0.002)	0.119 (-0.012)
<i>theory</i>	0.127 (-)	0.138 (-)	0.131 (-)
# POI	16,016	7,443	3,915
Root Types	Restaurants; Shopping; Food; Nightlife; Automotive; Active Life; Arts & Entertainment; Financial Services	Beauty & Spas; Health & Medical; Local Services; Hotels & Travel; Professional Services; Public Services & Government	Home Services; Event Planning & Services; Pets; Education



# Geo-Aware Image Classification



(Mac Aodha et al., 2019)

	BirdSnap <sup>†</sup>	NABirds <sup>†</sup>
No Prior (i.e. uniform)	70.07	76.08
Nearest Neighbor (num)	77.76	79.99
Nearest Neighbor (spatial)	77.98	80.79
Adaptive Kernel (Berg et al., 2014)	78.65	81.11
<i>tile</i> (Tang et al., 2015) (location only)	77.19	79.58
<i>wrap</i> (Mac Aodha et al., 2019) (location only)	78.65	81.15
<i>rbf</i> ( $\sigma=1k$ )	78.56	81.13
<i>grid</i> ( $\lambda_{min}=0.0001$ , $\lambda_{max}=360$ , $S = 64$ )	<b>79.44</b>	81.28
<i>theory</i> ( $\lambda_{min}=0.0001$ , $\lambda_{max}=360$ , $S = 64$ )	79.35	<b>81.59</b>

# Geo-Aware Image Classification

Our **multi-scale location encoding** (*grid* and *theory*) can outperform 1) RBF (*rbf*); 2) tile-based approaches (*tile*); 3) single-scale location encoding (*wrap*).

	BirdSnap <sup>†</sup>	NABirds <sup>†</sup>
No Prior (i.e. uniform)	70.07	76.08
Nearest Neighbor (num)	77.76	79.99
Nearest Neighbor (spatial)	77.98	80.79
Adaptive Kernel (Berg et al. 2014)	78.65	81.11
<i>tile</i> (Tang et al. 2015) (location only)	77.19	79.58
<i>wrap</i> (Mac Aodha et al. 2019) (location only)	78.65	81.15
<i>rbf</i> ( $\sigma=1k$ )	78.56	81.13
<i>grid</i> ( $\lambda_{min}=0.0001, \lambda_{max}=360, S=64$ )	<b>79.44</b>	81.28
<i>theory</i> ( $\lambda_{min}=0.0001, \lambda_{max}=360, S=64$ )	79.35	<b>81.59</b>

