

Zero-Shot Multi-View Indoor Localization via Graph Location Networks

Meng-Jiun Chiou
National University of Singapore
Singapore, Singapore
mengjiun.chiou@u.nus.edu

Zhenguang Liu*
Zhejiang Gongshang University
Hangzhou, China
liuzhenguang2008@gmail.com

Yifang Yin
National University of Singapore
Singapore, Singapore
idsyin@nus.edu.sg

An-An Liu
Tianjin University
Tianjin, China
anan0422@gmail.com

Roger Zimmermann
National University of Singapore
Singapore, Singapore
rogerz@comp.nus.edu.sg

ABSTRACT

Indoor localization is a fundamental problem in location-based applications. Current approaches to this problem typically rely on Radio Frequency technology, which requires not only supporting infrastructures but human efforts to measure and calibrate the signal. Moreover, data collection for all locations is indispensable in existing methods, which in turn hinders their large-scale deployment. In this paper, we propose a novel neural network based architecture Graph Location Networks (GLN) to perform infrastructure-free, multi-view image based indoor localization. GLN makes location predictions based on robust location representations extracted from images through message-passing networks. Furthermore, we introduce a novel zero-shot indoor localization setting and tackle it by extending the proposed GLN to a dedicated zero-shot version, which exploits a novel mechanism Map2Vec to train location-aware embeddings and make predictions on novel unseen locations. Our extensive experiments show that the proposed approach outperforms state-of-the-art methods in the standard setting, and achieves promising accuracy even in the zero-shot setting where data for half of the locations are not available. The source code and datasets are publicly available.¹

CCS CONCEPTS

•Computing methodologies → Computer vision; •Human-centered computing → Ubiquitous and mobile computing;

KEYWORDS

indoor localization; zero-shot learning; graph neural networks

ACM Reference format:

Meng-Jiun Chiou, Zhenguang Liu, Yifang Yin, An-An Liu, and Roger Zimmermann. 2020. Zero-Shot Multi-View Indoor Localization via Graph

*Corresponding author

¹<https://github.com/coldmanck/zero-shot-indoor-localization-release>

Permission to make digital or hard copies of part or all of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for third-party components of this work must be honored. For all other uses, contact the owner/author(s).

MM '20, Seattle, WA, USA

© 2020 Copyright held by the owner/author(s). 978-1-4503-7988-5/20/10...\$15.00
DOI: 10.1145/3394171.3413856

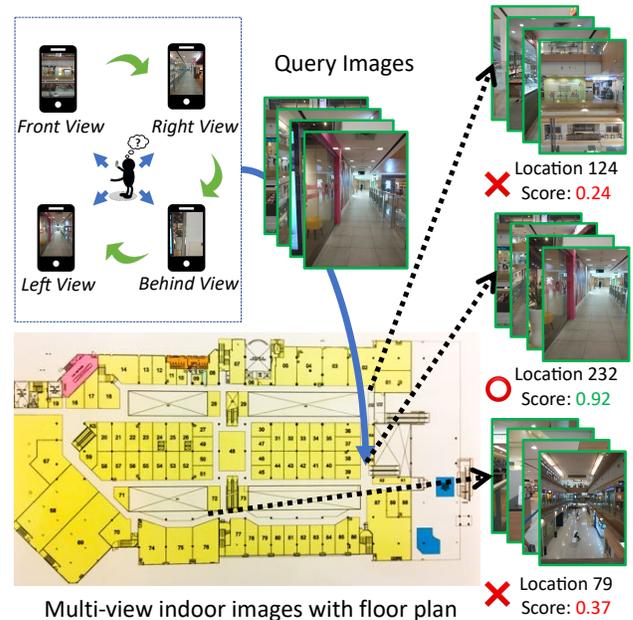


Figure 1: An illustration of a multi-view image-based indoor localization system. Current location of a user is predicted with the query images (photos of the user's surroundings).

Location Networks. In *Proceedings of Proceedings of the 28th ACM International Conference on Multimedia, Seattle, WA, USA, October 12–16, 2020 (MM '20)*, 10 pages.

DOI: 10.1145/3394171.3413856

1 INTRODUCTION

Indoor localization seeks to localize a user or a device in an indoor environment. Accurate indoor localization systems could enable various applications, e.g., guiding users in an underground parking lot to find a space, and in a large airport to get to the right boarding gate on time [41]; however, it remains to be an open challenge. While Global Positioning System (GPS) has been widely adopted to precisely localize a device in an outdoor environment with a 1- to 5-meter localization accuracy [18], it cannot be simply applied indoors since the GPS signal is significantly weakened after passing through

roofs and walls. Researchers have explored other techniques such as WiFi [39, 46], radio frequency identification (RFID) [16], optics [25], acoustics [19] and magnetism [39]. However, most of the existing approaches require additional infrastructures, such as WiFi access points, RF transmitters, or specially-designed optic/acoustic receivers. Moreover, manual periodic re-calibration is indispensable for RF-based methods since the signals are prone to fluctuation, which in turn is harder to maintain.

Purely image-based approaches [11, 24, 30] are proposed to alleviate part of the deployment costs by utilizing only indoor images. However, most of the existing methods still require special devices to collect data [24] and utilize fully supervised models to make predictions on each location [11, 30]. These limitations cause their approaches to be not only time-consuming but also labor-intensive when deployed in large-scale indoor environments. In this context, an interesting and fundamental question arises: *is it possible to infer the location of the user while data is collected for only several locations?* To answer this question, we will consider *zero-shot learning* where models recognize novel locations by transferring knowledge learned from seen to unseen classes.

Multi-view images are photos of different views at indoor locations and comprise rich location contexts. To leverage this information, a multi-view image- and geomagnetism-based localization strategy is proposed in [26] to transform the problem of indoor localization into a graph retrieval problem. However, they treat different camera views as of equal importance, which is usually not true especially when some views are similar and others are more representative. For example, for two neighboring locations (within one meter) at the same corridor, the views that parallel the corridor are extremely similar while the perpendicular views may consist of more identifiable objects (*e.g.*, different doors and windows). Treating them equally undermines the representativeness of their graph features during retrieval. Moreover, the geomagnetic signals which they rely on are unstable, hard to collect and prone to change with time.

Recently, deep learning [22] has achieved remarkable success in various areas, including but not limited to image-level understanding [15, 40], object-level detection [31, 32], human-level estimation [34, 44], text classification [4, 17] and audio understanding [45, 52]. To overcome the aforementioned problems, we exploit the strong representation power of neural networks and propose an infrastructure free, neural network-based architecture Graph Location Networks (GLN) to perform multi-view indoor localization. Given photos of different views, GLN computes robust node representations by aggregating and updating features from neighboring nodes, in which identifiable features permeate the whole graph. A location prediction is made by feeding the representations to a single fully-connected layer. Our proposed approach requires neither any infrastructure nor special devices but only a camera phone to collect the photo database. In addition to evaluating on the publicly available multi-view indoor dataset (*i.e.*, ICUBE [26]), we provide a benchmark dataset WCP that has been collected in a shopping center. We show that our approach outperforms the baseline and existing methods in terms of localization accuracy by a large margin.

Furthermore, to motivate researches in reducing data collection labor costs we introduce a novel task named *zero-shot indoor localization*, in which half of the locations are masked during training while a system is required to predict the precise user location. We propose a three-step framework to tackle this task and demonstrate the efficiency of it by extending our GLN to a dedicated zero-shot version. Specifically, to transfer the knowledge from the seen locations to the unseen ones, we propose the Map2Vec mechanism that trains location-aware embeddings for both seen and unseen classes by incorporating their geometric contexts of the floor plan. These embeddings are then leveraged to train a compatibility function that maps image-class pairs to scalar scores. Finally, a prediction is made by picking out the best class maximizing the score function of the query image. We demonstrate that, trained through the proposed framework, our model not only surpasses the baseline by a large margin but also achieves promising localization accuracy, *e.g.*, 56.3% 5-meter accuracy on the ICUBE dataset, while the query locations are never seen during training. To the best of our knowledge, our work is the first exploration of enabling zero shot recognition for indoor localization.

The key contributions of our work are summarized as follows: (a) We propose a novel, neural network based architecture Graph Location Networks which performs effective, infrastructure-free multi-view indoor localization. (b) We introduce zero-shot indoor localization and propose a training framework to tackle it. We demonstrate the efficiency by extending our proposed architecture to a dedicated zero-shot version. (c) We contribute an additional multi-view image-based indoor localization dataset. Our extensive experiments shows that the proposed approach significantly outperforms state-of-the-art methods in the fully supervised setting and achieves competitive localization accuracy in the zero-shot setting.

2 RELATED WORK

2.1 Indoor Localization

Indoor localization has been a popular topic ever since the outdoor localization was mostly tackled [18]. Most of the previous efforts rely on RF technology which requires additional transmitters/receivers to estimate the location [16, 19, 25, 39, 46] or special devices to collect data [8], causing large-scale deployment to be costly and prohibitive. More recently and related to our work, a multi-view image- and geomagnetism-based method has been proposed to formulate indoor localization into a graph retrieval problem [26]; however, it does not consider the difference between views and thus fails to capture a robust representation. While purely image-based techniques do not require additional facilities, they either need special devices to do data collection in advance [24] or require a user to take photos with specific reference objects [11]. Therefore both are not ideal methods to achieve a pervasive indoor positioning system. In addition, for all existing image-based approaches it is inevitable to collect data of all locations of interest, resulting in costly deployment for large-scale indoor environments. In our work, to implement a truly infrastructure-free indoor localization system, we adopt a purely image-based approach which does not require any special device, only a camera phone, to construct

an image database. Furthermore, we introduce zero-shot indoor localization to reduce data collection labor costs. Note that while our method is closely related to outdoor place recognition [3, 27, 43] and is possible to incorporate corresponding techniques (e.g. NetVLAD layer [3], contrastive loss [7] or 2D-3D hybrid method [36]) to improve the performance, we focus on the graph-based location network with zero-shot setting in this work and leave as possible extensions in future work.

2.2 Graph-based Methods

2.2.1 Graph analysis. Graph has garnered a lot of attention from researchers due to its nature of being suitable for representing data in various real-life applications [54], including protein-protein interaction [10], social relationship networks [14], natural science [6, 35], and knowledge graphs [13]. The typical problems that graph analysis is dealing with include node classification, link prediction and clustering. Graph Neural Networks (GNNs) [5, 20, 48] have become the de facto standard for processing graph-based data for their ability to work on large-scale graphs by borrowing the ideas of weight-sharing and local connections from Convolutional Neural Networks (CNNs) [54].

2.2.2 Graph embedding. Nodes in a graph can be represented as feature vectors by incorporating the information of the graph topology and initial node feature [12]. In our work, we leverage GNNs in two scenarios: (a) to perform message passing on a locally-connected location graph for a more robust representation, and (b) to train location embeddings to encode position information for all locations to perform zero-shot recognition.

2.3 Zero-Shot Learning

Unlike traditional supervised learning, zero-shot learning aims to recognize the instance classes that have never been seen by the model during training [50]. There has been an increasing interest in zero-shot learning and its applications [21, 33, 49, 51] since it is not unusual that data is only available for some classes. To transfer knowledge to unseen classes, a compatibility function is learned to relate semantic attributes to features [2, 42]. Specifically, in our work, we learn a compatibility function which maps image features to the location embeddings, and the predicted location is chosen as the one maximizing the compatibility score. Following the definition in [50], we perform *generalized zero-shot learning* since our search space contains both training and test classes during testing.

3 METHODOLOGY

In this section, we first formulate the indoor localization problem under fully supervised setting, followed by introducing our proposed method, Graph Location Networks (GLN), which serves as the backbone under both settings. We then demonstrate how to extend our approach to a dedicated zero-shot version to perform indoor positioning on locations of unseen classes.

3.1 Problem Formulation

We formulate the image-based indoor localization problem as follows. Given $\mathbf{x} \in \mathcal{X}$ that denotes a set of images and \mathcal{X} the space of

all sets of images, we are to predict the location $y \in \mathcal{Y}$ for \mathbf{x} , where $\mathcal{Y} = \{y_1, \dots, y_k\}$ is the set of all k locations. The goal is to learn a function $f: \mathcal{X} \rightarrow \mathcal{Y}$ that maps the input \mathbf{x} to the target class y . In our settings, \mathbf{x} comprises images of four different directions, i.e., images of the front, behind, right and left at a location.

3.2 Standard Graph Location Networks

The main idea of our proposed approach is that different views of a location possessing distinct information can be used to form a holistic representation. To take advantage of this, we formulate a locally-connected graph in which features are being refined during the message passing and finally producing a robust location representation for classification. We define Graph Location Networks (GLN) as an indoor localization approach which includes three major modules: feature extraction module, location prediction module, and especially message passing module to exploit the aforementioned graph to make accurate location prediction. We explain each module in detail in the following subsections. Figure 2 shows an overview of GLN.

3.2.1 Feature Extraction Module. Given a set of images of the front, behind, right and left views $\mathbf{x} = \{x_1, x_2, x_3, x_4\}$ at a specific location, we utilize a Convolutional Neural Networks based backbone φ that takes \mathbf{x} as input to extract high-dimensional features $\mathbf{r} = \varphi(\mathbf{x}) = \{r_1, r_2, r_3, r_4\}$, $r_i \in \mathcal{R}^d$, where d is the feature dimension. The choice of the backbone network and hyperparameters of the model are given in section 4.1.

3.2.2 Message Passing Module. We define a quadrilateral graph $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ for four views by $\mathcal{V} = \{v_1, v_2, v_3, v_4\}$ and $\mathcal{E} = \{e_{12}, e_{23}, e_{34}, e_{41}\}$, where e_{ij} denotes an undirected edge between nodes v_i and v_j . Node v_i represents a specific direction and its hidden state is initialized with r_i of that direction. To obtain a robust location representation, our system has to effectively exploit and combine neighboring features. Graph Neural Networks (GNNs) have been shown to be able to aggregate information of neighbor nodes and update the node’s hidden state accordingly [20, 37]. We employ GNNs to pass messages within the graph and refine the hidden states of the nodes. Let $h_i^l \in \mathcal{R}^F$ and $h_i^{l+1} \in \mathcal{R}^{F'}$ be the hidden state of node i at layer l and $l + 1$, the updating procedure of hidden state h_i of node v_i is defined as follows:

$$h_i^l = \begin{cases} r_i, & \text{if } l = 1 \\ \sigma \left(\sum_{j \in \mathcal{N}(i)} \frac{1}{\alpha_{ij}} W^{l-1} h_j^{l-1} \right), & \text{otherwise} \end{cases} \quad (1)$$

where $\mathcal{N}(i)$ denotes the set of neighboring nodes of node v_i , σ is a nonlinear activation function, $\alpha_{ij} = \sqrt{|\mathcal{N}(i)\mathcal{N}(j)|}$ is a normalization constant and W^l represents a shared weight matrix for node-wise feature transformation at layer l .

However, each neighboring node (i.e., $\mathcal{N}(i)$) should not have an equal affect to node i , e.g., some neighbors may share more overlapped scenes than others. Attention mechanism [47, 48] has been demonstrated to be effective to capture relational representation. We introduce a graph self-attention mechanism to assign different weights to each neighbor according to its importance to node i . Specifically, we update h_i at layer l with the weight α_{ij}^l , which is

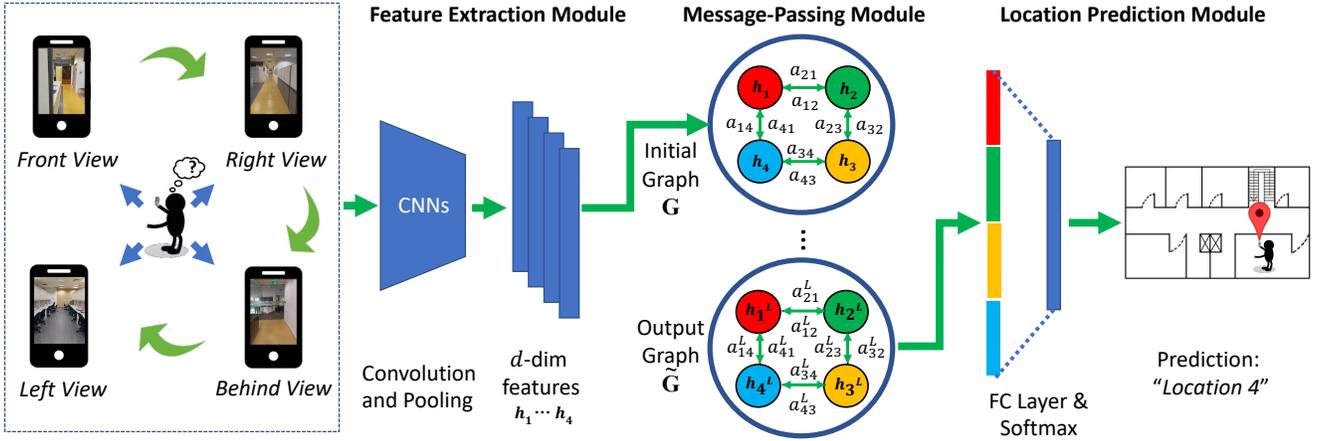


Figure 2: The architecture of our Graph Location Networks (GLN) based indoor localization system. First, features of the front, behind, right and left views extracted through CNNs are taken as input by a multi-view quadrilateral graph. An attentional message-passing algorithm is performed on the graph to extract robust location representation, which is then passed into a fully-connected layer followed by a softmax function to make prediction.

defined as follows:

$$\alpha_{ij}^l = \frac{\exp(\sigma(a[W^l h_i^l, W^l h_j^l]))}{\sum_{k \in \mathcal{N}(i)} \exp(\sigma(a[W^l h_i^l, W^l h_k^l]))}, \quad (2)$$

where $[\cdot]$ denotes the concatenation operation, and a is a shared attention mechanism that computes the importance of node j 's feature to node i .

3.2.3 Location Prediction Module. After L layers of message-passing the final hidden states \mathbf{h} are fused to form a single robust representations \mathbf{x}^r as the following:

$$\mathbf{x}^r = [h_1^L, h_2^L, h_3^L, h_4^L]. \quad (3)$$

\mathbf{x}^r is then passed into a single fully-connected layer mapping the concatenated feature vector into the location space, followed by a Softmax function to generate a probability distribution over all classes: $p' = f(\text{Softmax}(\text{FC}(\mathbf{x}^r)))$. We adopt the Softmax loss as the objective function as follows:

$$\mathcal{L} = - \sum_i p_i \log(p'_i), \quad (4)$$

where p'_i is the i -th row of p' and p_i denotes the ground truth label for location i .

3.3 Zero-Shot Graph Location Networks

In this section, we describe the proposed learning framework that enables indoor localization models to perform zero-shot prediction. We use our proposed GLN as the backbone model.

In the zero-shot setting, \mathcal{Y} is divided into two disjointed sets: $\mathcal{Y}_s \subset \mathcal{Y}$, $y_s \in \{y_1, \dots, y_n\}$ denotes a set of n seen classes, and $\mathcal{Y}_u \subset \mathcal{Y}$, $y_u \in \{y_{n+1}, \dots, y_k\}$ represents a set of $(k - n)$ unseen classes, where $\mathcal{Y} = \mathcal{Y}_s \cup \mathcal{Y}_u$, $\mathcal{Y}_s \cap \mathcal{Y}_u = \emptyset$. Note that we assign \mathcal{Y}_s and \mathcal{Y}_u alternately (one every other) on the map. Refer to Figures 3 for an illustration.

For zero-shot indoor localization, the goal is to enable the system to recognize photos of *unseen classes* \mathcal{Y}_u through training only on photos of *seen classes* \mathcal{Y}_s . It is impossible to employ traditional supervised learning methods to train a model that can recognize the unseen classes without seeing them before. Instead, we leverage the information that is available to both groups (*i.e.*, floor plans) to bridge them together. There are three key steps to perform zero-shot indoor localization: (a) training the Map2Vec location embeddings, (b) learning a compatibility function with the embeddings and GLN, and (c) performing zero-shot recognition.

3.3.1 Map2Vec Location Embedding. To overcome the aforementioned problem, we propose the *Map2Vec* mechanism to learn location-aware graph embeddings to correlate the seen and unseen classes. Figure 3(a) shows an illustration of this procedure. For a given map (floor plan) with k locations, we define a graph $\mathcal{G}' = (\mathcal{V}', \mathcal{E}')$ where each vertex v' represents a location and each edge e' is a path between locations. Similar to GLN in the standard setting, we adopt the Graph Neural Networks as in Eq. 1 to train graph structure-aware node embeddings and initialize each of the hidden states using the coordinate $o_i \in \mathbb{R}^2$ for location i . After L layers of message-passing, we extract the final hidden state $h_i^L \in \mathbb{R}^k$ as the location embedding for class y_i : $\psi(y_i)$.

3.3.2 Compatibility Function. To take advantage of the learned location embeddings, we aim at doing *knowledge transfer* so that the indoor localization knowledge can be transferred from seen to unseen classes. To carry out the knowledge transfer, we utilize a compatibility function $F : \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}$ which is a mapping from an image-class pair to a scalar score for the specific class. Figure 3(b) shows an illustration of learning the compatibility function. Since only the samples from seen classes are used for learning the compatibility function, it should be in a class-agnostic form. We follow [42] and define the compatibility function in a bilinear form as follows:

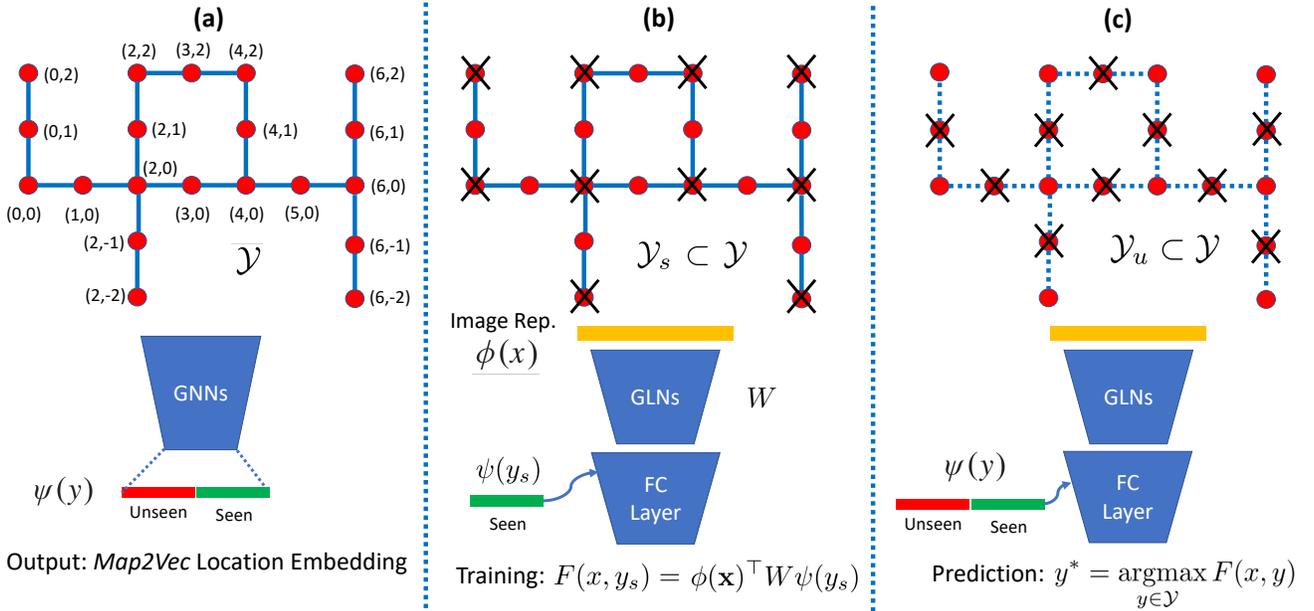


Figure 3: The key steps for training a zero-shot indoor localization model. (a) Train Map2Vec location embeddings for a given map (floor plan). (b) Learn a compatibility function with only the seen classes (circles without multiplication sign). (c) Perform zero-shot prediction by assigning an input to the location that maximizes the compatibility function. Dotted lines mean that the edge information is not available during testing. The "GLNs" block can be replaced with any indoor localization model.

$$F(x, y_s) = \phi(x)^\top W \psi(y_s), \quad (5)$$

where $\phi(x) \in \mathbb{R}^d$ is the image representation of an image x from seen classes, $\psi(y_s) \in \mathbb{R}^k$ is the location embedding of a seen class y_s and $W \in \mathbb{R}^{d \times k}$ is the weights that we are actually learning. In this context, W is in fact our GLN that takes in d -dimensional image representation and output k -dimensional logits. Similar to standard indoor localization, we adopt cross entropy loss as the objective function.

3.3.3 Zero-Shot Recognition. Once the compatibility function is learned, we can utilize it to make predictions on unseen classes. Refer to Figure 3(c) for an illustration. Zero-shot indoor localization is achieved by assigning the query image a location class y^* that maximizes $F(x, y)$:

$$y^* = \operatorname{argmax}_{y \in \mathcal{Y}} F(x, y), \quad (6)$$

Unlike [42], we predict on all possible locations $y \in \mathcal{Y}$ instead of on only unseen classes \mathcal{Y}_u to simulate real-world scenarios.

4 EXPERIMENTS

In this section, we conduct extensive experiments to evaluate the proposed method. Towards this aim, we first explain the implementation details, evaluation datasets and metrics. We then compare our GLN based indoor localization systems with existing models under both standard and zero-shot settings.

4.1 Implementation Details

We implement our model based on *PyTorch* [29] framework and train on a single NVIDIA Titan X. To extract image representation, we adopt ResNet-152 [15] and utilize off-the-shelf weights from *torchvision* package of PyTorch. We employ data augmentation technique to randomly flip, rotate by 10 degrees and resize to 256×256 pixels, followed by randomly cropping a patch of 224×224 pixels. The output of the CNNs is a 2,048-dimensional feature for each image ($d = 2048$). For both of our standard GLN and zero-shot GLN, we adopt Graph Convolutional Networks [20] as the backbone of message passing process and the attention mechanism in Graph Attention Networks (GATs) [48], and we utilize the implementation provided by *PyTorch Geometric* [9]. An undirected edge is implemented with two directed edges of opposite directions in the experiments. We observe one layer of message propagation ($L = 1$) empirically gives the best performance. The dimension of the latent representation h_i^l is 256. We utilize ReLU nonlinear activation for the original GLN and adopt LeakyReLU for the attentional GLN at each layer, while both are followed by a batch normalization layer and a dropout layer to stabilize training. Attention mechanism a is implemented with a single FC layer. We train the model in an end-to-end manner with learning rate 3×10^{-4} with Adam optimizer of the exponential decay rate 0.9 and 0.999 for the first- and the second-moment estimates respectively.

4.2 Evaluation Datasets and Metrics

We evaluate our proposed method on two datasets: ICUBE [26] that is publicly available and WCP that is collected by ourselves.

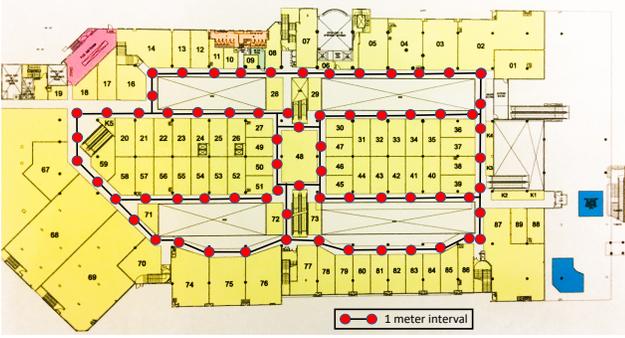


Figure 4: An illustration of WCP dataset where the red vertices represent locations, and black edges denote the adjacency of vertices. Note that the locations are not drawn to scale and are for illustrative purposes only.

4.2.1 *ICUBE dataset.* The ICUBE dataset contains 2,896 photos of 214 locations in an academic building. For standard indoor localization, to perform a fair comparison, we closely follow the original paper [26] to divide the dataset into a training set of 1,712 images and test set of 1,184 images. While in the zero-shot setting, we set aside 1,368 images of 102 locations as seen classes, where 1,092 of them are for training and the other 276 of them are for validation during training the compatibility function. The remaining 1,528 images of 112 locations are set as unseen classes to be used in zero-shot recognition.

4.2.2 *WCP dataset.* The WCP dataset consists of 3,280 photos of 394 locations in a shopping center. We assign 2,624 images for training and the other 656 images for testing in the standard indoor localization experiment. In the zero-shot setting, 1,696 images of 204 locations are assigned as seen classes, in which 1,360 and 336 of them are for training and validation compatibility function, respectively. The other 1,584 images of 190 locations are unseen classes. Overall, WCP is more difficult than ICUBE due to its higher number of classes and more complicated scenes such as shops and restaurants. Both datasets are collected in 1-meter distance interval and have a corresponding map that has vertices of locations and edges of adjacency. Figure 4 shows an illustration of the WCP dataset.

4.2.3 *Evaluation Metrics.* We report **one-meter-level accuracy** and **Cumulative Distribution Function of localization error (CDF@k)** at distance k . For zero-shot indoor localization, to perform a more detailed evaluation of the models’ strengths and weaknesses, we utilize multiple metrics including **CDF@k**, **Recall@k** that sees if the ground truth presents in top k predictions ordered by confidence scores, and **Median Error Distance (MED)** which calculates the error distance of 50-percentile predictions. Note that the distance unit is 1-meter for CDF and MED.

4.3 Quantitative Results

4.3.1 *Standard Indoor Localization.* To compare with existing indoor localization methods, we choose not only those that are purely based on images but also those based on signals. Pedes [23]

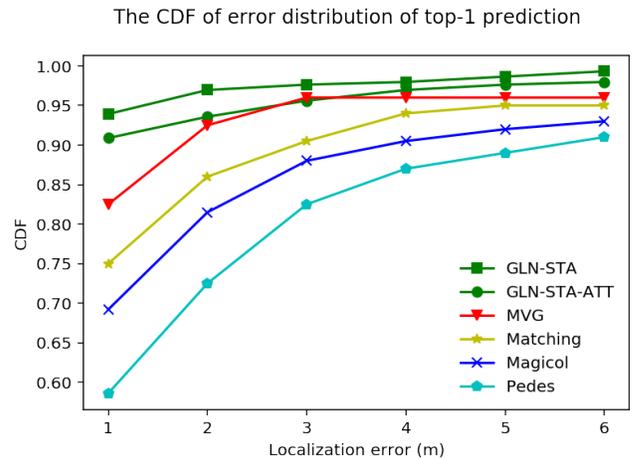


Figure 5: The cumulative distribution function (CDF) curves of the localization error of the previous and our approaches in standard indoor localization setting on ICUBE dataset.

Dataset	Method	Meter-level Accuracy
ICUBE	Pedes [23]	58.30%
	Magical [39]	69.20%
	Matching [30]	75.00%
	MVG [26]	82.50%
	GLN-STA	93.92%
	GLN-STA-ATT	90.88%
MALL-1†	Sextant [11]	47%
MALL-2‡	GeoImage [24]	53%
WCP	GLN-STA	79.88%
	GLN-STA-ATT	79.88%

Table 1: Performance comparison with state-of-the-art models on ICUBE, WCP and the respective MALL datasets. Results of previous approaches on ICUBE are taken from [26], while results on distinct MALL datasets are taken from their respective papers. †MALL-1 consists of 108 locations and 686 images. ‡Mall-2 contains 20,000 images (locations).

is a pedestrian dead reckoning localization method using inertial sensors. Magical [39] incorporates geomagnetic field and WiFi signal to perform indoor positioning. Matching [30] performs image comparison by scoring with multiple off-the-shelf algorithms. MVG [26] is a multi-view localization method via graph retrieval based on images and geomagnetism. Note that Magical and MVG are not purely image-based methods. GLN-STA is our original GLN and GLN-STA-ATT is the GLN with self attention mechanism. The upper part of Table 1 shows the meter-level accuracy compared to existing image-based methods, where our GLN variants surpass the others with significantly higher within one-meter accuracy on the ICUBE dataset and improve the state-of-the-art by 13.8%.

We observe that the usage of attention mechanism does not help on both datasets. Note that the scale of our quadrilateral graph

Dataset	Method	Recall@k					CDF@k					MED
		k=1	k=2	k=3	k=5	k=10	k=1	k=2	k=3	k=5	k=10	
ICUBE	Baseline-coord	0.00	0.01	0.02	0.03	0.03	3.53	3.73	5.96	11.65	23.95	23.00
	GLN-ZS	8.12	14.40	22.78	30.89	46.60	19.90	33.77	45.81	56.28	74.87	3.76
	GLN-ZS-ATT	8.38	14.92	23.30	32.20	45.81	18.59	34.55	43.71	55.24	73.04	4.09
WCP	Baseline-coord	0.00	0.00	0.00	0.00	0.00	1.01	1.01	2.78	3.79	8.84	27.00
	GLN-ZS	2.02	6.06	7.83	12.37	24.75	8.84	13.38	17.42	22.98	50.25	9.97
	GLN-ZS-ATT	2.02	4.55	8.33	13.64	24.50	9.09	13.38	19.70	25.00	51.52	9.93

Table 2: Results of zero-shot indoor localization in comparison of Recall@k, CDF@k and Median Error Distance (MED) on ICUBE and WCP datasets. Note that numbers of recall and CDF are in % (the higher the better), while the numbers of median error distance are in meter (the lower the better). MED results are estimated with linear interpolation.

is very different from the common graph datasets (e.g. citation networks [38], WebKB graphs [1]) that have thousands of nodes and edges. Moreover, it was observed in [28, 53] that the common instantiation of the attention mechanism on GNNs (i.e. GATs) does not necessarily bring performance boost over standard GCNs on distinct graph datasets. Therefore, more investigation into the way of instantiating and applying graph attention mechanism in our architecture is needed and we leave it as our future work.

Figure 5 shows the full localization error curve (CDF@k), where our GLN perform consistently better than previous approaches, regardless of the requirement of infrastructures.

We also list the additional results of the previous approaches that cannot be reproduced to evaluate on our datasets due to their infrastructure requirements.² Sextant [11] leverages image matching algorithms to identify and match with the pre-selected reference objects. GeoImage [24] performs image matching against a geo-referenced 3D image dataset. Their localization accuracy on the respective shopping mall datasets and our GLN variants on the WCP dataset is showed at the lower part of Table 1. Our GLN-variants achieve significantly better localization performance than the previous methods without any infrastructure requirement.

4.3.2 Zero-Shot Indoor Localization. To simulate the real-word case, while we only perform the localization on data of unseen classes \mathcal{Y}_{te} , we still make predictions on all possible locations. To demonstrate that our proposed approach helps in zero-shot localization, we implement a baseline method Baseline-coord that utilize the coordinates o but not the Map2Vec location embeddings $\psi(y)$ to train the compatibility function. Baseline-coord uses the same standard GLN as the backbone architecture.

The experimental results of zero-shot indoor localization on the ICUBE and WCP dataset are shown in table 2. GLN-ZS is the original GLN and GLN-ZS-ATT is the attentional GLN, both in the zero-shot setting. On both datasets, GLN-ZS and GLN-ZS-ATT outperform the baseline approach by a large margin. In specific on the ICUBE dataset, GLN-ZS significantly outperforms Baseline-coord by achieving 56.3% 5-meter accuracy (CDF@5) and median

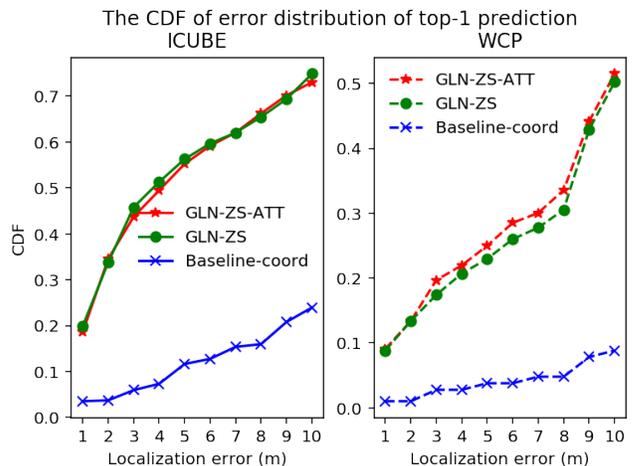


Figure 6: The cumulative distribution function (CDF) curves of the localization error of the zero-shot indoor localization experiments on ICUBE (left) and WCP (right) datasets.

error of 3.76 meters, which are considered promising since all test locations are never seen during training.

Similar to the observation in the experiments of the standard setting, GLN-ZS-ATT has similar performance to GLN-ZS on both ICUBE and WCP. In addition to the possible reasons we discussed in the previous section, we find that it may also results from the relatively monotonic scenes in the ICUBE dataset so that the attention mechanism does not help much on distinguishing views. In contrast, GLN-ZS-ATT has slight performance improvements over GLN-ZS in terms of CDF@k and MED on WCP.

Overall, compared to experiments on ICUBE, both variants of GLN perform less powerful on the metrics, presumably due to higher variance of scenes and a larger number of classes.

Figure 6 shows the full CDF@k curves of zero shot GLN variants and Baseline-coord, where ours perform consistently better. For example on the ICUBE dataset, GLN-ZS shows strong performance improvements ranging from 3.1 to 5.6 times higher CDF@k than the baseline, demonstrating the benefit of the Map2Vec embedding. In

²Note that since they were evaluated on distinct shopping center datasets, the results may not be directly comparable and serve for reference purposes.

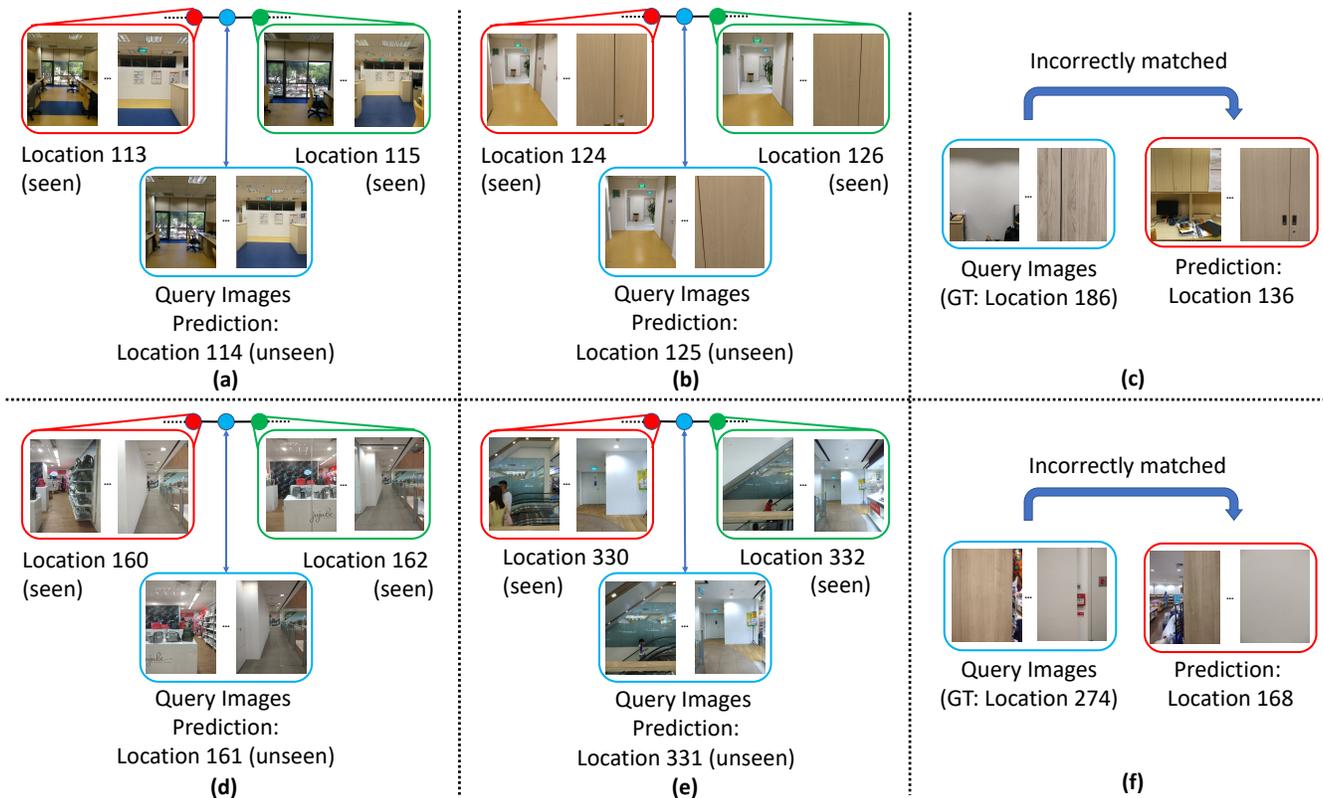


Figure 7: Qualitative results of zero-shot indoor localization on ICUBE (the top row) and WCP (the bottom row) dataset. The first two columns show examples of successful localization cases by utilizing the adjacency of seen classes to unseen classes, where the red, blue and green circles represent three adjacent locations. The last column shows examples of unsuccessful localization cases where our system is misled, especially when there are more query photos lacking distinguishable features.

addition, as mentioned above, GLN-ZS-ATT is shown to have more consistent performance improvements over GLN-ZS especially on the harder WCP dataset.

4.4 Qualitative Results for Zero Shot GLN

To better identify the strengths and weaknesses of our proposed zero-shot approach, we perform qualitative analysis for GLN-ZS in zero-shot indoor localization setting. The left two columns of Figure 7 show examples of successful localization where the correct prediction is made by inferring that the location (e.g., location 114 in Fig. 7(a)) of the query images is between two neighboring seen locations (e.g., location 113 and 115) by referring to the Map2Vec location embeddings. While some of the views that parallel corridor are extremely similar (e.g., the second view of Fig. 7(d) and (e)), our model is able to extract robust representation by passing identifiable features from neighboring views (e.g., the first view of Fig. 7(d) and (e)) to infer the correct location. However, our system could still be misled especially when there are more query photos lacking distinguishable features. The last column (Fig. 7(c) and (f)) shows unsuccessful localization cases where more images contain no distinguishable features. For instance, in Fig. (c) the first

query photo consists of mostly a white wall and the second photo contains merely the surface of a cabinet.

5 CONCLUSION

In this paper, we first propose a novel neural network based architecture, namely Graph Location Networks (GLN) to perform multi-view indoor localization. GLN takes in photos of different views and make location predictions based on robust location representations with the message-passing mechanism. To reduce prohibitive labor cost when deployed in large-scale indoor environments, we introduce a novel task named zero-shot indoor localization and propose a effective learning framework which is used to adapt GLN to a dedicated zero-shot version to make predictions on unseen locations. We evaluate our proposed approach not only on the publicly available ICUBE dataset but also on our own benchmark dataset WCP that we make publicly available to facilitate researches in multi-view indoor localization systems. Experimental results show that our proposed method achieves state-of-the-art results in the standard setting and performs well with promising accuracy in the zero-shot setting.

ACKNOWLEDGMENTS

This research is partly supported by the Natural Science Foundation of Zhejiang Province, China (No. LQ19F020001), the National Natural Science Foundation of China (No. 61902348, U1609215, 61976188, 61672460), and Singapore's Ministry of Education (MOE) Academic Research Fund Tier 1, grant number T1 251RES1713.

REFERENCES

- [1] 2001. CMU World Wide Knowledge Base (Web-zKB) project. <http://www.cs.cmu.edu/afs/cs.cmu.edu/project/theo-11/www/wkwb/>. Accessed: 2020-08-04.
- [2] Zeynep Akata, Scott E. Reed, Daniel Walter, Honglak Lee, and Bernt Schiele. 2015. Evaluation of Output Embeddings for Fine-grained Image Classification. In *IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2015, Boston, MA, USA, June 7-12, 2015*. IEEE Computer Society, 2927–2936. <https://doi.org/10.1109/CVPR.2015.7298911>
- [3] Relja Arandjelovic, Petr Gronát, Akihiko Torii, Tomás Pajdla, and Josef Sivic. 2016. NetVLAD: CNN Architecture for Weakly Supervised Place Recognition. In *2016 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2016, Las Vegas, NV, USA, June 27-30, 2016*. IEEE Computer Society, 5297–5307. <https://doi.org/10.1109/CVPR.2016.572>
- [4] Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio. 2015. Neural Machine Translation by Jointly Learning to Align and Translate. In *3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings*, Yoshua Bengio and Yann LeCun (Eds.). <http://arxiv.org/abs/1409.0473>
- [5] Peter W Battaglia, Jessica B Hamrick, Victor Bapst, Alvaro Sanchez-Gonzalez, Vinicius Zambaldi, Mateusz Malinowski, Andrea Tacchetti, David Raposo, Adam Santoro, Ryan Faulkner, et al. 2018. Relational Inductive Biases, Deep Learning, and Graph Networks. *arXiv preprint arXiv:1806.01261* (2018).
- [6] Peter W. Battaglia, Razvan Pascanu, Matthew Lai, Danilo Jimenez Rezende, and Koray Kavukcuoglu. 2016. Interaction Networks for Learning about Objects, Relations and Physics. In *Advances in Neural Information Processing Systems 29: Annual Conference on Neural Information Processing Systems 2016, December 5-10, 2016, Barcelona, Spain*, Daniel D. Lee, Masashi Sugiyama, Ulrike von Luxburg, Isabelle Guyon, and Roman Garnett (Eds.), 4502–4510. <http://papers.nips.cc/paper/6418-interaction-networks-for-learning-about-objects-relations-and-physics>
- [7] Sean Bell and Kavita Bala. 2015. Learning Visual Similarity for Product Design with Convolutional Neural Networks. *ACM Trans. Graph.* 34, 4 (2015), 98:1–98:10. <https://doi.org/10.1145/2766959>
- [8] Jaewoo Chung, Matt Donahoe, Chris Schmandt, Ig-Jae Kim, Pedram Razavai, and Micaela Wiseman. 2011. Indoor Location Sensing using Geo-magnetism. In *Proceedings of the 9th international conference on Mobile systems, applications, and services*. ACM, 141–154.
- [9] Matthias Fey and Jan E. Lenssen. 2019. Fast Graph Representation Learning with PyTorch Geometric. In *ICLR Workshop on Representation Learning on Graphs and Manifolds*.
- [10] Alex Fout, Jonathon Byrd, Basir Shariat, and Asa Ben-Hur. 2017. Protein Interface Prediction using Graph Convolutional Networks. In *Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, 4-9 December 2017, Long Beach, CA, USA*, Isabelle Guyon, Ulrike von Luxburg, Samy Bengio, Hanna M. Wallach, Rob Fergus, S. V. N. Vishwanathan, and Roman Garnett (Eds.), 6530–6539. <http://papers.nips.cc/paper/7231-protein-interface-prediction-using-graph-convolutional-networks>
- [11] Ruipeng Gao, Yang Tian, Fan Ye, Guojie Luo, Kaigui Bian, Yizhou Wang, Tao Wang, and Xiaoming Li. 2016. Sextant: Towards Ubiquitous Indoor Localization Service by Photo-Taking of the Environment. *IEEE Trans. Mob. Comput.* 15, 2 (2016), 460–474. <https://doi.org/10.1109/TMC.2015.2418205>
- [12] Palash Goyal and Emilio Ferrara. 2018. Graph Embedding Techniques, Applications, and Performance: A Survey. *Knowl. Based Syst.* 151 (2018), 78–94. <https://doi.org/10.1016/j.knsys.2018.03.022>
- [13] Takuo Hamaguchi, Hidekazu Oiwa, Masashi Shimbo, and Yuji Matsumoto. 2017. Knowledge Transfer for Out-of-Knowledge-Base Entities: A Graph Neural Network Approach. In *Proceedings of the Twenty-Sixth International Joint Conference on Artificial Intelligence, IJCAI 2017, Melbourne, Australia, August 19-25, 2017*, Carles Sierra (Ed.), ijcai.org, 1802–1808. <https://doi.org/10.24963/ijcai.2017/250>
- [14] William L. Hamilton, Zitao Ying, and Jure Leskovec. 2017. Inductive Representation Learning on Large Graphs. In *Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, 4-9 December 2017, Long Beach, CA, USA*, Isabelle Guyon, Ulrike von Luxburg, Samy Bengio, Hanna M. Wallach, Rob Fergus, S. V. N. Vishwanathan, and Roman Garnett (Eds.), 1024–1034. <http://papers.nips.cc/paper/6703-inductive-representation-learning-on-large-graphs>
- [15] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 2016. Deep Residual Learning for Image Recognition. In *2016 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2016, Las Vegas, NV, USA, June 27-30, 2016*. IEEE Computer Society, 770–778. <https://doi.org/10.1109/CVPR.2016.90>
- [16] Sverre Holm. 2009. Hybrid Ultrasound-RFID Indoor Positioning: Combining the Best of Both Worlds. In *2009 IEEE International Conference on RFID*. IEEE, 155–162.
- [17] Armand Joulin, Edouard Grave, Piotr Bojanowski, and Tomas Mikolov. 2017. Bag of Tricks for Efficient Text Classification. In *Proceedings of the 15th Conference of the European Chapter of the Association for Computational Linguistics, EACL 2017, Valencia, Spain, April 3-7, 2017, Volume 2: Short Papers*, Mirella Lapata, Phil Blunsom, and Alexander Koller (Eds.). Association for Computational Linguistics, 427–431. <https://doi.org/10.18653/v1/e17-2068>
- [18] Jii Kárník and Jakub Streit. 2016. Summary of Available Indoor Location Techniques. *IFAC-PapersOnLine* 49, 25 (2016), 311–317. <https://doi.org/10.1016/j.ifacol.2016.12.055>
- [19] Hong-Shik Kim and Jong-Suk Choi. 2008. Advanced Indoor Localization using Ultrasonic Sensor and Digital Compass. In *2008 International Conference on Control, Automation and Systems*. IEEE, 223–226.
- [20] Thomas N. Kipf and Max Welling. 2017. Semi-Supervised Classification with Graph Convolutional Networks. In *5th International Conference on Learning Representations, ICLR 2017, Toulon, France, April 24-26, 2017, Conference Track Proceedings*. OpenReview.net. <https://openreview.net/forum?id=SJU4ayYgl>
- [21] Christoph H. Lampert, Hannes Nickisch, and Stefan Harmeling. 2014. Attribute-Based Classification for Zero-Shot Visual Object Categorization. *IEEE Trans. Pattern Anal. Mach. Intell.* 36, 3 (2014), 453–465. <https://doi.org/10.1109/TPAMI.2013.140>
- [22] Yann LeCun, Yoshua Bengio, and Geoffrey Hinton. 2015. Deep Learning. *nature* 521, 7553 (2015), 436–444.
- [23] Fan Li, Chunshui Zhao, Guanzhong Ding, Jian Gong, Chenxing Liu, and Feng Zhao. 2012. A Reliable and Accurate Indoor Localization Method Using Phone Inertial Sensors. In *Proceedings of the 2012 ACM Conference on Ubiquitous Computing (UbiComp '12)*. ACM, New York, NY, USA, 421–430. <https://doi.org/10.1145/2370216.2370280>
- [24] Jason Zhi Liang, Nicholas Corso, Eric Turner, and Avideh Zakhor. 2013. Image Based Localization in Indoor Environments. In *Fourth International Conference on Computing for Geospatial Research and Application, COM.Geo '13, San Jose, CA, USA, July 22-24, 2013*. IEEE, 70–75. <https://doi.org/10.1109/COMGEO.2013.11>
- [25] Xiaohan Liu, Hideo Makino, and Kenichi Mase. 2010. Improved Indoor Location Estimation Using Fluorescent Light Communication System with a Nine-Channel Receiver. *IEICE Trans. Commun.* 93-B, 11 (2010), 2936–2944. <https://doi.org/10.1587/transcom.E93.B.2936>
- [26] Zhenguang Liu, Li Cheng, Anan Liu, Luming Zhang, Xiangnan He, and Roger Zimmermann. 2017. Multiview and Multimodal Pervasive Indoor Localization. In *Proceedings of the 2017 ACM on Multimedia Conference, MM 2017, Mountain View, CA, USA, October 23-27, 2017*, Qiong Liu, Rainer Lienhart, Haohong Wang, Sheng-Wei “Kuan-Ta” Chen, Susanne Boll, Yi-Ping Phoebe Chen, Gerald Friedland, Jia Li, and Shuicheng Yan (Eds.). ACM, 109–117. <https://doi.org/10.1145/3123266.123436>
- [27] Stephanie Lowry, Niko Sünderhauf, Paul Newman, John J Leonard, David Cox, Peter Corke, and Michael J Milford. 2015. Visual Place Recognition: A Survey. *IEEE Transactions on Robotics* 32, 1 (2015), 1–19.
- [28] Hesham Mostafa and Marcel Nassar. 2020. Permutohedral-GCN: Graph Convolutional Networks with Global Attention. *arXiv preprint arXiv:2003.00635* (2020).
- [29] Adam Paszke, Sam Gross, Soumith Chintala, Gregory Chanan, Edward Yang, Zachary DeVito, Zeming Lin, Alban Desmaison, Luca Antiga, and Adam Lerer. 2017. Automatic Differentiation in PyTorch. In *NIPS Autodiff Workshop*.
- [30] Nishkam Ravi, Pravin Shankar, Andrew Frankel, Ahmed M. Elgammal, and Liviu Iftode. 2006. Indoor Localization Using Camera Phones. In *Seventh IEEE Workshop on Mobile Computing Systems & Applications, WMCSA'06, Semiahmoo Resort, Washington, USA, April 6-7, 2006*. IEEE Computer Society, 19. <https://doi.org/10.1109/WMCSA.2006.11>
- [31] Joseph Redmon, Santosh Kumar Divvala, Ross B. Girshick, and Ali Farhadi. 2016. You Only Look Once: Unified, Real-Time Object Detection. In *2016 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2016, Las Vegas, NV, USA, June 27-30, 2016*. IEEE Computer Society, 779–788. <https://doi.org/10.1109/CVPR.2016.91>
- [32] Shaoqing Ren, Kaiming He, Ross B. Girshick, and Jian Sun. 2015. Faster R-CNN: Towards Real-Time Object Detection with Region Proposal Networks. In *Advances in Neural Information Processing Systems 28: Annual Conference on Neural Information Processing Systems 2015, December 7-12, 2015, Montreal, Quebec, Canada*, Corinna Cortes, Neil D. Lawrence, Daniel D. Lee, Masashi Sugiyama, and Roman Garnett (Eds.), 91–99. <http://papers.nips.cc/paper/5638-faster-r-cnn-towards-real-time-object-detection-with-region-proposal-networks>
- [33] Bernardino Romera-Paredes and Philip H. S. Torr. 2015. An embarrassingly simple approach to zero-shot learning. In *Proceedings of the 32nd International Conference on Machine Learning, ICML 2015, Lille, France, 6-11 July 2015 (JMLR Workshop and Conference Proceedings)*, Francis R. Bach and David M. Blei (Eds.), Vol. 37. JMLR.org, 2152–2161. <http://proceedings.mlr.press/v37/romera-paredes15.html>

- [34] Weijian Ruan, Wu Liu, Qian Bao, Jun Chen, Yuhao Cheng, and Tao Mei. 2019. POINet: Pose-Guided Ovonic Insight Network for Multi-Person Pose Tracking. In *Proceedings of the 27th ACM International Conference on Multimedia, MM 2019, Nice, France, October 21-25, 2019*, Laurent Amsaleg, Benoit Huet, Martha A. Larson, Guillaume Gravier, Hayley Hung, Chong-Wah Ngo, and Wei Tsang Ooi (Eds.). ACM, 284–292. <https://doi.org/10.1145/3343031.3350984>
- [35] Alvaro Sanchez-Gonzalez, Nicolas Heess, Jost Tobias Springenberg, Josh Merel, Martin A. Riedmiller, Raia Hadsell, and Peter W. Battaglia. 2018. Graph Networks as Learnable Physics Engines for Inference and Control. In *Proceedings of the 35th International Conference on Machine Learning, ICML 2018, Stockholm, Sweden, July 10-15, 2018 (Proceedings of Machine Learning Research)*, Jennifer G. Dy and Andreas Krause (Eds.), Vol. 80. PMLR, 4467–4476. <http://proceedings.mlr.press/v80/sanchez-gonzalez18a.html>
- [36] Paul-Edouard Sarlin, Frédéric Debraine, Marcin Dymczyk, and Roland Siegwart. 2018. Leveraging Deep Visual Descriptors for Hierarchical Efficient Localization. In *2nd Annual Conference on Robot Learning, CoRL 2018, Zürich, Switzerland, 29-31 October 2018, Proceedings (Proceedings of Machine Learning Research)*, Vol. 87. PMLR, 456–465. <http://proceedings.mlr.press/v87/sarlin18a.html>
- [37] Franco Scarselli, Marco Gori, Ah Chung Tsoi, Markus Hagenbuchner, and Gabriele Monfardini. 2009. The Graph Neural Network Model. *IEEE Trans. Neural Networks* 20, 1 (2009), 61–80. <https://doi.org/10.1109/TNN.2008.2005605>
- [38] Prithviraj Sen, Galileo Namata, Mustafa Bilgic, Lise Getoor, Brian Gallagher, and Tina Eliassi-Rad. 2008. Collective Classification in Network Data. *AI Magazine* 29, 3 (2008), 93–106. <https://doi.org/10.1609/aimag.v29i3.2157>
- [39] Yuanhao Shu, Cheng Bo, Guobin Shen, Chunhui Zhao, Liquan Li, and Feng Zhao. 2015. Magicol: Indoor Localization Using Pervasive Magnetic Field and Opportunistic WiFi Sensing. *IEEE J. Sel. Areas Commun.* 33, 7 (2015), 1443–1457. <https://doi.org/10.1109/JSAC.2015.2430274>
- [40] Karen Simonyan and Andrew Zisserman. 2015. Very Deep Convolutional Networks for Large-Scale Image Recognition. In *3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings*, Yoshua Bengio and Yann LeCun (Eds.). <http://arxiv.org/abs/1409.1556>
- [41] Deborah Snoonian. 2003. Smart buildings. *IEEE Spectrum* 40, 8 (2003), 18–23.
- [42] Gencer Sumbul, Ramazan Gokberk Cinbis, and Selim Aksoy. 2018. Fine-Grained Object Recognition and Zero-Shot Learning in Remote Sensing Imagery. *IEEE Trans. Geosci. Remote. Sens.* 56, 2 (2018), 770–779. <https://doi.org/10.1109/TGRS.2017.2754648>
- [43] Niko Sünderhauf, Sareh Shirazi, Feras Dayoub, Ben Upcroft, and Michael Milford. 2015. On the Performance of ConvNet Features for Place Recognition. In *2015 IEEE/RSJ International Conference on Intelligent Robots and Systems, IROS 2015, Hamburg, Germany, September 28 - October 2, 2015*. IEEE, 4297–4304. <https://doi.org/10.1109/IROS.2015.7353986>
- [44] Alexander Toshev and Christian Szegedy. 2014. DeepPose: Human Pose Estimation via Deep Neural Networks. In *2014 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2014, Columbus, OH, USA, June 23-28, 2014*. IEEE Computer Society, 1653–1660. <https://doi.org/10.1109/CVPR.2014.214>
- [45] Aäron van den Oord, Sander Dieleman, Heiga Zen, Karen Simonyan, Oriol Vinyals, Alex Graves, Nal Kalchbrenner, Andrew W. Senior, and Koray Kavukcuoglu. 2016. WaveNet: A Generative Model for Raw Audio. In *The 9th ISCA Speech Synthesis Workshop, Sunnysvale, CA, USA, 13-15 September 2016*. ISCA, 125. http://www.isca-speech.org/archive/SSW_2016/abstracts/ssw9_DS-4_van_den_Oord.html
- [46] Deepak Vasisht, Swarun Kumar, and Dina Katabi. 2016. Decimeter-Level Localization with a Single WiFi Access Point. In *13th USENIX Symposium on Networked Systems Design and Implementation, NSDI 2016, Santa Clara, CA, USA, March 16-18, 2016*, Katerina J. Argyraki and Rebecca Isaacs (Eds.). USENIX Association, 165–178. <https://www.usenix.org/conference/nsdi16/technical-sessions/presentation/vasisht>
- [47] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. 2017. Attention is All you Need. In *Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, 4-9 December 2017, Long Beach, CA, USA*, Isabelle Guyon, Ulrike von Luxburg, Samy Bengio, Hanna M. Wallach, Rob Fergus, S. V. N. Vishwanathan, and Roman Garnett (Eds.). 5998–6008. <http://papers.nips.cc/paper/7181-attention-is-all-you-need>
- [48] Petar Velickovic, Guillem Cucurull, Arantxa Casanova, Adriana Romero, Pietro Liò, and Yoshua Bengio. 2018. Graph Attention Networks. In *6th International Conference on Learning Representations, ICLR 2018, Vancouver, BC, Canada, April 30 - May 3, 2018, Conference Track Proceedings*. OpenReview.net. <https://openreview.net/forum?id=rjXmpikCZ>
- [49] Wei Wang, Vincent W. Zheng, Han Yu, and Chunyan Miao. 2019. A Survey of Zero-Shot Learning: Settings, Methods, and Applications. *ACM Trans. Intell. Syst. Technol.* 10, 2 (2019), 13:1–13:37. <https://doi.org/10.1145/3293318>
- [50] Yongqin Xian, Christoph H. Lampert, Bernt Schiele, and Zeynep Akata. 2019. Zero-Shot Learning - A Comprehensive Evaluation of the Good, the Bad and the Ugly. *IEEE Trans. Pattern Anal. Mach. Intell.* 41, 9 (2019), 2251–2265. <https://doi.org/10.1109/TPAMI.2018.2857768>
- [51] Yongqin Xian, Bernt Schiele, and Zeynep Akata. 2017. Zero-Shot Learning - The Good, the Bad and the Ugly. In *2017 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2017, Honolulu, HI, USA, July 21-26, 2017*. IEEE Computer Society, 3077–3086. <https://doi.org/10.1109/CVPR.2017.328>
- [52] Yifang Yin, Meng-Jiun Chiou, Zhenguang Liu, Harsh Shrivastava, Rajiv Ratn Shah, and Roger Zimmermann. 2019. Multi-Level Fusion based Class-aware Attention Model for Weakly Labeled Audio Tagging. In *Proceedings of the 27th ACM International Conference on Multimedia, MM 2019, Nice, France, October 21-25, 2019*, Laurent Amsaleg, Benoit Huet, Martha A. Larson, Guillaume Gravier, Hayley Hung, Chong-Wah Ngo, and Wei Tsang Ooi (Eds.). ACM, 1304–1312. <https://doi.org/10.1145/3343031.3351090>
- [53] Li Zhang, Heda Song, and Haiping Lu. 2018. Graph Node-Feature Convolution for Representation Learning. *CoRR abs/1812.00086* (2018). <http://arxiv.org/abs/1812.00086>
- [54] Jie Zhou, Ganqu Cui, Zhengyan Zhang, Cheng Yang, Zhiyuan Liu, and Maosong Sun. 2018. Graph Neural Networks: A Review of Methods and Applications. *CoRR abs/1812.08434* (2018). <http://arxiv.org/abs/1812.08434>