

WSNNet: Stacked Bidirectional LSTM with Residual Attention for Indoor Localization of Wireless Sensor Network

HYUNGTAELIM¹, (Student Member, IEEE), CHANGGYU PARK¹, (Student Member, IEEE), AND HYUN MYUNG.¹, (Senior Member, IEEE)

¹Urban Robotics Laboratory, Korea Advanced Institute of Science and Technology, Daejeon 34141, South Korea.

Corresponding author: Hyun Myung (hmyung@kaist.ac.kr).

This material is based upon work supported by the Ministry of Trade, Industry & Energy(MOTIE, Korea) under Industrial Technology Innovation Program. No.10067202, 'Development of Disaster Response Robot System for Lifesaving and Supporting Fire Fighters at Complex Disaster Environment'.

ABSTRACT

As verified experimentally, this new proposal represents a significant improvement in accuracy, computation time, and robustness against outliers.

INDEX TERMS Enter key words or phrases in alphabetical order, separated by commas. For a list of suggested keywords, send a blank e-mail to keywords@ieee.org or visit http://www.ieee.org/organizations/pubs/ani_prod/keywrd98.txt

I. INTRODUCTION

WIRELESS sensor networks(WSN) signifies a number of the sensors which send and receive the signal each other. In recent years, advancement in micro-electro-mechanical systems(MEMS) contributes improvement of sensors' performance and miniaturization at the same time in such a way as to have enabled the wireless sensors to communicate with other wireless sensors efficiently and conveniently. By virtue of this low-cost, small-size and acceptively accurate performance, WSNs are widely utilized in various areas, such as rescue works, surveillance, pollution monitoring, inspection of structures, and so on [1]–[4].

Especially, WSNs are also employed in various location-aware applications, e.g., tracking and localization of objects. Because WSNs can be easily installed, they have been suggested as a solution for localization on the indoor environment [5], [6] where the signals of the Global Positioning Systems(GPS) cannot be received. WSNs consist of two types of nodes: anchor nodes(ANs) that interchange signals to tag nodes(TNs) for estimating positions of TNs and TNs whose position are estimated by the gathered data. Several methods for localizing TNs are have been proposed and these are broadly categorized into two categories based on whether the range between the ANs and TNs are measured or not. One method, called range-based, is to directly measure the

anchor-to-tag distance or angle, e.g., time of arrival(TOA) [7], time difference of arrival(TDOA) [8], time of flight(TOF) [6], arrival of angle(AOA) [9] and the other method is to calculate position of the TNs using connectivities of the WSNs, e.g., hop counts [10] or weighted centroid [11]. Usually, range-based method mostly have better performance than range-free method [12] and particularly, TOF-based technique is commonly used in practice because it does not need to synchronize, and is easy to implement, and consume less resources than other approaches [13]

However, these range-based approaches suffer from *rank deficiency* problem, [14], which means range measurements only consist of single value to represent distance between each TN and AN respectively in such a way as to be deficient to describe exact position or orientation of the MNs. Besides, only single value could represent the range measurement, these measurements have huge uncertainties caused by non-line-of-sight(NLOS) problem and multipath fading channel(MPF) problem [13]. To alleviate these issues, many studies are conducted how to localize more precisely. For example, in [13], confidence-based intersection method is proposed for robust trilateration method. In [15], particle filter algorithm is introduced to estimate the position of a MN attached to a aerial robot while suppressing the noise. Also, many machine learning techniques have been intro-

duced: one authors utilized support vector machine(SVM) for localizaiton [16]–[19], other author developed method support vector regression(SVR) for localization [20], [21].

In the meantime, as deep learning age has come [22], WSNs fields also have introduced various kinds of neural network architectures [8], [10], [23]–[25]. An main advantage of neural networks-based approaches is that the neural networks enables themselves to recognize the position of the MNs quite accurately through the training even if there is no mathematical description. The neural networks also cover all noises that occurs when ANs and MNs interchange their signals. The authors show that their nerual networks-based approach have better performance than traditional approaches or machine learning-based approaches, yet most networks are based on The Multilayer Perceptron(MLP). In addition, the some authors just show the feasibility of their approaches, i.e., they trained and tested on the simulations situation. Plus, some authors let the neural networks infer the estimates on the two-dimensional space, but gap of the complexity of estimating position between on two-dimensional space and on three-dimensional using *rank-deficient* measurements cannot be negligible.

In this paper, we propose a stacked bidirectional Long Short-Term Memory(stacked Bi-LSTM) with residual attention for more accurate localization of a MN. Our contributions can be summarized as follows. First of all, Using deep learning, our structure directly learns the end-to-end mapping between range data measured by ANs-to-MN distances and the position of MN on three dimensional space. Unlike other approaches whose authors implement MLP-based neural networks, we exploit LSTM architecture in such a way as not only to take present range data but also take previous temporal data as input for estimating the MN's position. Second, we verify that effect of the residual attention in the way that the residual attention layer makes the network better understand context of the feature when dealing with low-dimensional input. Finally, we implement the networks presented at related works and compare which neural architecture perperforms better when real-world data are given as input. As a result, and our networks perform the best performance among other proposed approaches. Our system overview is shown in the figure 1.

II. RELATED WORKS

In the past few years, researchers have conducted the studies using machine learning-based approaches to reduce computational complexity and localize a mobile node more precisely. One authors utilized support vector machine(SVM) for localizaiton, [16]–[19], other author developed method support vector regression(SVR) for localization [20], [21]. In [16], authors suggested two SVMs for localization, called LSVMs, one LSVM infers x-dimension and the other LSVM infers y-dimension. To employeeing LSVMs, they divide the field into $M-l$ x-classes and $M-l$ y-classes, like grid, and this deployment has had an impact on succeeding studies [18], [19], [26]. Samadian *et al.* [27] introduced probabilistic

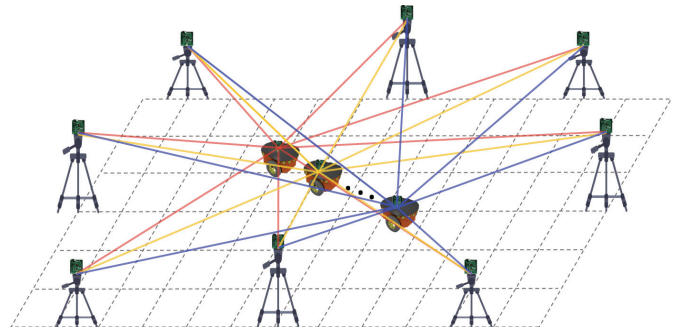


FIGURE 1: System overview.

support vector machine for localization and they showed that probabilistic vector machine has better performance than LSVM. In terms of SVR, Lee *et al.* suggested various types of SVR for localization [20], [21]

Especially, several works have been studied using neural networks to localize nodes of the range measurement sensors on the indoor space while covering range measurements' uncertainties. Regarding previous proposals, Chenna *et al.* first shows the suitability that Kalman filter could be replaced with the RNN when estimating states and tracking nodes [28]. However, they did not provide numerical analysis, so Shareef *et al.* did and conducted their experiment in the real-world. They concluded Radial Basis Function(RBF) may be the best option among the suggested Kalman filter models and RNN [29].

Similarly, many researchers also have achived considerable improvement to localize position of mobile node by exploiting MLP in WSNs fields [8], [10], [23]–[25], [29], [30]. In case of range-based method, Rahman *et al.* [23] have considered the neural networks for mapping between RSS and corresponding position of sensor nodes and let neural networks be trained by the train data gathered by the sensor nodes that are eqaully spaced over x-axis and y-axis. In [8], Singh *et al.* compared that performance of Multilayer Back propagation Network Model(MLBPN) and Radial Basis Function Network Model(RBFN) and the authors show that RBFN performs better than MLBPN when the number of the sensor nodes is larger than 220. Abdelhadi *et al.* [24] presented two artificial intelligence techniques: Sugeno-type fuzzy system and neural networks system. In addition, the authors conducted experiment on three-dimensional space in such a way as verified the feasibility of localization by utilizing nueral networks in 3D space. Kumar *et al.* [25] also introduced the neural networks and evaluated five different training techniques,e.g., Levenberg-Marquardt (LM), Bayesian Regularization (BR), Resilient Back-propagation

TABLE 1: Comparison of previous studies with our approach

Localization method	Dimension	Type of input	Train data	Test data mobility	Implementation environment
RBF [29]	2D	Temporary	Grid	Dynamic nodes ✓	Real-world ✓
MLP [23]	2D	Temporary	Grid	Static nodes	Simulation
MLBPN [8]	2D	Temporary	Grid	Static nodes	Simulation
MLP [24]	3D ✓	Temporary	Spread ✓	Static nodes	Simulation
c-FCNNs [30]	2D	Temporary	Spread ✓	Static nodes	Real-world ✓
MLP [25]	2D	Temporary	Grid	Static nodes	Real-world ✓
LPSONN [10]	2D	Temporary	Spread ✓	Static nodes	Simulation
Ours	3D ✓	Sequential ✓	Spread ✓	Dynamic nodes ✓	Real-world ✓

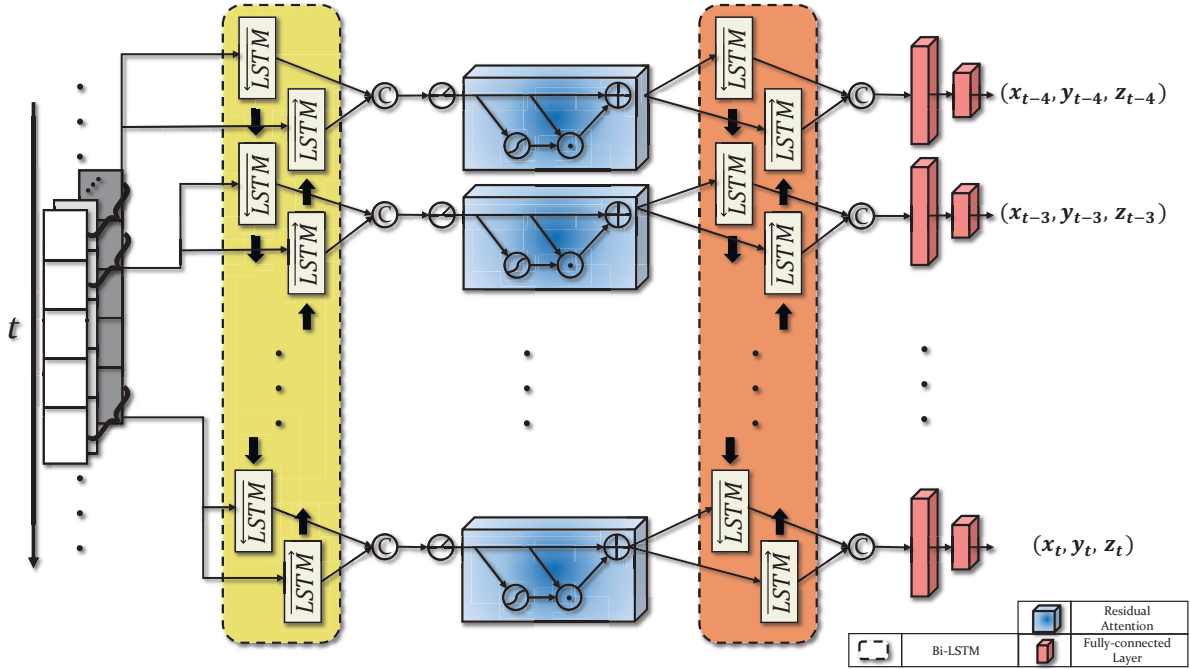


FIGURE 2: Our proposed network architecture.

(RP), Scaled Conjugate Gradient (SCG) and Gradient Descent (GD), to find optimal way to train neural networks with the best accuracy. Recently, [10] have proposed the neural networks with novel training technique, called Particle Swarm Optimization(PSO) and prove thire network, called LPSONN, has better localization accuracy than previous machine learning method, soft computing method, and previously proposed network.

However, there are some points that could have been better. First of all, in some cases, their networks were trained by range measurement data corresponding position of mobile node in simulation environment [8], [10], [23], [26], [29]. The simulation situation is almost ideal in the point that the NLOS or MPF problem does not occur. Even though in [23], they artificially generate NLOS data, but the gap of complexity between real-world and artificially generated data exists, since the generated data has less noise than those of real-world necessarily. For these reasons, it is hard to say that their networks also works well on real world situation. Therefore, to test on whether it is possible for neural networks to estimate position with covering all disturbance,

the experiment should be conducted on real-world. Second, some studies was conducted on the two-dimensional space to simplify their problem definition [8], [10], [23], [25], [29], [30]. However, gap of the complexity of estimating position between on two-dimensional space and on three-dimensional using *rank-deficient* measurements cannot be negligible. Finally, in previous studies [8], [23], [25], [29], it has a possibility of overfitting because the authors generate the grid-map as train data in such a way as to restrict their ground truth region. That is to say, their finite ground truth indicates where the sensors are placed at the equal distance interval so that neural networks may recognize the only locations included in the grid are correct even when the position of MN to be tested is quite far from the grid. Therefore their grid map train impedes the optimization of neural networks to cover all over the region.

III. WSN NET

In this chapter, we explain how our proposed residual attention-based stacked Bi-LSTM is implemented, as illustrated in Fig. 2. In detail, we introduce the neural networks

concepts that we choose for localizaing the tag node and the describe the reason why we let the neural network infer in three-dimensional space even though experiment is conducted on the mobile robot. Finally, we explain how to set the loss function of our neural network and then compare to those of other previous works.

A. LONG SHORT-TERM MEMORY

Recurrent Neural Networks(RNN) is a special artificial neural networks in the way that it has a loop, so that RNN can deal with temporal information for sequential modeling. It originally used in the natural language processing, speech recognition, and image captioning area. By virtue of a loop, RNN can remember past data and past situation and respond appropriately to the present situation based on these past information.

But unfortunately, as the time-sequential gap grows, RNNs become unable to learn the relationship of these sequential information. This issue is called the problem of *Long-Term Dependency*, which fail to propagate the previous matter into present tasks so that long-term dependency lead to a failure of learning. In other words, RNNs are not able to learn to store appropriate internal states and operate on long-term trends. That is the reason why the Long Short-Term Memory (LSTM) architecture is introduced to solve this long-term dependency problem and make the networks possible to learn longer-term contextual understandings [31]. That's why LSTM have been actively studied for many tasks in a wide area of science and engineering. In most of the deep learning research areas and numerous variations of LSTM architectures have been studied.

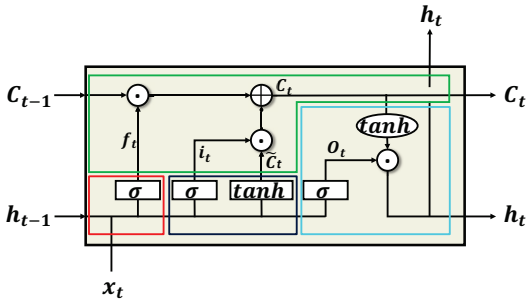


FIGURE 3: Architecture of the LSTM. It consists of 3 gates, forget gate(the part inside the red box), input gate(the part inside the blue box), and output gate. And output gate is divided into cell state layer(the part inside the green box) and output gate layer(the part inside the cyan box)

Unlike RNN that consist only of hidden state, in LSTM, cell state is added on the network. The cell state consists of the 3 gates to preserve the previous information and control the cell state: forget gate, input gate, and output gate and equations of those are as follows:

$$f_t = \sigma_s(W_{xf} \cdot x_t + W_{hf} \cdot h_{t-1} + b_f) \quad (1)$$

$$i_t = \sigma_s(W_{xi} \cdot x_t + W_{hi} \cdot h_{t-1} + b_i) \quad (2)$$

$$\tilde{c}_t = \tanh(W_{xc} \cdot x_t + W_{hc} \cdot h_{t-1} + b_c) \quad (3)$$

$$c_t = f_t \odot c_{t-1} + i_t \odot \tilde{c}_t \quad (4)$$

$$o_t = \sigma_s(W_{xo} \cdot x_t + W_{ho} \cdot h_{t-1} + b_o) \quad (5)$$

$$h_t = o_t \odot \tanh(c_t) \quad (6)$$

where σ_s is a kind of activation function, called *sigmoid*, f_t , i_t , and o_t respectively indicates the forget gate, input gate, and output gates, and c_t denotes cell states. And \odot denotes element-wise multiplication, called *Hadamard product*. Entire gates are activated by sigmoid function and cell states are activated by tanh function.

The Forget gate layer, f_t , decides how much information to forget. The sigmoid layer, which is the activation function of f_t , takes previous hidden state, h_t , and present input, x_t and outputs a number between 0 and 1. Note that 1 indicates "totally keep the previous cell state, C_{t-1} " and 0 indicate "totally forget C_{t-1} " (1). Next, the input gate, i_t , decides how much information to embrace when updating the cell state. i_t are also from the sigmoid function layer (2) and \tanh generates the new candidate cell state, \tilde{c}_t , which ranges from -1 to 1 (3). After that, c_t is updated by the cell state layer based on f_t , i_t , and \tilde{c}_t (4). In addition, output gate layer, o_t , serves as a filter, which means o_t determine what values are going to output (5) in such a way as that present hidden state, h_t , is updated based on o_t updated cell state, c_t (6).

B. BIDIRECTIONAL LSTM

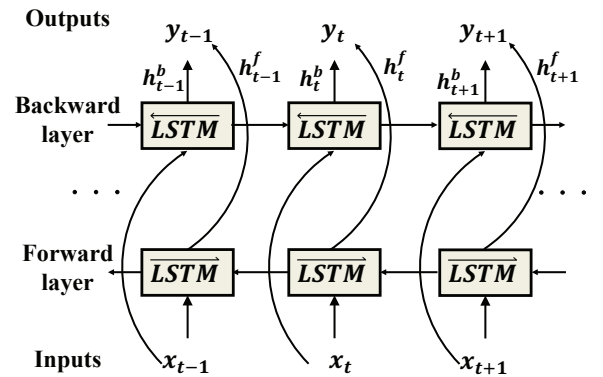


FIGURE 4: Architecture of the Bidirectional LSTM(Bi-LSTM)

One shortcoming of conventional RNNs is that they only exploit previous context to update h_t and c_t . However, in many cases dealing with sequential data, it could be efficient to extract well-described context by utilizing future context as well. Bidirectional RNNs are introduced [32] for that reason and bidirectional RNNs process the data in both directions

with two separate hidden layers. Especially, bidirectional LSTM, which we employee, has one forward LSTM and one backward LSTM running in reverse time and their features are combined at the output layer, y_t . As a result, bidirectional LSTM can produce more appropriate context considering both past and future at the same time. By virtue of this characteristics, bidirectional LSTM is popularly utilized for many tasks to model their sequentional systems [33]–[35].

As FIGURE. 4 shown, bidirectional LSTM consist of 2 LSTMs: one forward LSTM layer, \overrightarrow{LSTM} , and one backward LSTM layer, \overleftarrow{LSTM} . Let assume the hidden state of \overrightarrow{LSTM} at the time step t be h_t^f and \overleftarrow{LSTM} at the time step t be h_t^b , the hidden states and output sequence, y_t are calcaulated as follows:

$$h_t^f = \overrightarrow{LSTM}(x_t, h_{t-1}^f) \quad (7)$$

$$h_t^b = \overleftarrow{LSTM}(x_t, h_{t+1}^b) \quad (8)$$

$$y_t = \sigma_R(W_{h^f y} \cdot h_t^f + W_{h^b y} \cdot h_t^b + b_y) \quad (9)$$

where σ_R denotes activation function called ReLu. Note that in our case, we concatenate h_t^f and h_t^b to preverse their contexts seperately as our range measurement, which is gathered by the tag node and each anchor node, suffer from the *rank-deficiency*, which means range-based measurement consist of one-dimensional data [36]. Hence, we judge that it would be more helpful to increase the number of features naturally by concatenating two hidden states rather than adding them and when the network infer the position .

C. STACKED ARCHITECTURE

Recently, researchers show that the deeper the architecture of neural networks, the better their performance [37], [38] and their demonstrations has opened a deep learning area. Likewise, many authors have analyzed variations of LSTM architecture and find out that stacking multiple layers of the LSTM improve the performance for many tasks [35], [39], [40]. In other words, as the number of stacked layers is getting large, the more activation functions which rise the non-linearity within the networks are stacked in such a way as that complexity of networks increases. As a results, networks could model more complex system by virtue of these increased non-linerity.

Therefore, we also construct our networks by stacking two LSTM to increase the non-linearity. Note that stacking more than three LSTM doesn't show the improvement of performance. We suppose that activation funtions within the LSTM cause the *vanishing gradient problem* [41], which the networks fail to training due to the fact that the gradient is getting closer to zero during the backpropagation. We deem that this problem could comes from the sigmoid function and *tanh* function that compose the part of LSTM. Consequently, we put the ReLu function between LSTMs to avoid the vanishing gradient problem [42], instead of stacking LSTM.

D. RESIDUAL ATTENTION LAYER

A Attention layer is powerful module nowadays and mostly improves performance of neural network. Originally, neural networks treats information equally. But, using attention layer, neural networks can be ATTENDED what it should be examined closely. At the first time, attention is utilized at natural language processing area for improving translation performance [43]. But nowadays, attention layer is employed in many areas to improve the performance of the networks. For example, Jaderbeg *et al.* [44] introduced the attention layer to let the neural networks attend to spatial information. In addition, attention is even utilized to pose estimation and optimization [45], detection [46], and video captioning [47]

To precisely estimate the position of the tag node, it is important for the network to distinguish which is more meaningful context on time step T to help contextual understanding of our networks. So, we add the attention layer between the LSTM and the attention layer take on a role of a feature selector [48]. The equation of the attention machanism is as follows:

$$H(x) = M(x) \odot x \quad (10)$$

where x denotes the output of previous neural network layer, $H(x)$ denotes the output of the attention layer to be passed to the next neural network layer and $M(x)$ denotes the attention mask. The attention layer takes s as input and outputs the $H(x)$. By element-wise multiplying x by $M(x)$, attention layer makes the network weight more crutial information.

Despite of the improvement of the performance, the attention layer has potential risks that it may dilutes the features because attention mask ranges over 0 to 1. To alleviate this problem, residual attention layer is introduced in our network as follows [48]:

$$H(x) = (1 + M(x)) \odot x \quad (11)$$

As blue cuboid shape in the FIGURE 2 shown, this idea is originated from the Residual Net(ResNet) [38] that has skip connection in such a way as to mitigate aforementioned dilution problem and help the network to be trained well. Likewise the ResNet, residual attention also has other branch to calculate how much to attend and the branch is joined with original feature vector x . Each hidden state has each residual attention layer so that these attention modules can determine which time stamp has more fruitful meaning and deliver the output to second bidirectional LSTM.

E. TRAINING LOSS

In this section, we describe the method for training our network. Generally, let n be the number of mobile nodes and m be the number of the anchor nodes, data set are represented as follows:

$$\{(l_{11}, l_{12}, \dots, l_{1m}, P_1), \dots, (l_{n1}, l_{n2}, \dots, l_{nm}, P_n)\} \quad (12)$$

where l_{ij} denotes the the distance between i^{th} mobile node and j^{th} anchor node, P_i denotes the position of mobile node, which consist of 2D (x and y), or 3D (x, y , and z). In other words, data consist of set of distance data corresponding to the position of mobile nodes. Consequently, neural network could be optimized to be able to localize the mobile node when take distance set as input.

our neural network does not only take a set of distance data but takes sets of distance data on the time step T where T indicates sequential length of input to our network. And in our case, one mobile node is only placed on the robot. Therefore, data are formulated as follows:

$$\mathbb{L} = (L_t, Y_t) \quad (13)$$

where $L_t = \{(l_1, l_2, \dots, l_m)_t\}$ denotes input range measurement between a tag node and each anchor node at the time t . We omit the part of subscript that indicates i^{th} mobile node because we have only one mobile node. Y_t denotes the ground truth of the robot's 3D position, which is denoted as $Y_t = \{x_t, y_t, z_t\}$.

Let Θ be the parameters of our network model and assume that the trained network model could be expressed as conditional probability as follows:

$$P(Y_t | L_{t-T+1}, L_{t-T+2}, \dots, L_t) = p((x_t, y_t, z_t) | (l_1, \dots, l_m)_{t-T+1}, (l_1, \dots, l_m)_{t-T+2}, \dots, (l_1, \dots, l_m)_t) \quad (14)$$

Note that other studies only consider the input on time t , yet our approach consider temporal information about the range measurement data. Then, our final goal is to find optimal parameters Θ^* for localization by minimizing L_2 loss term. The L_2 loss term indicates mean square error (MSE) of Euclidean distance between ground truth position Y_k and estimated position \hat{Y}_k as follows:

$$\Theta^* = \arg\min_{\Theta} \frac{1}{N} \sum_{k=1}^N \|Y_k - \hat{Y}_k\|^2 \quad (15)$$

F. WHY ON THREE-DIMENSIONAL?

One may ask a question that why we infer the robot's position on the three dimensional space even test are being conducted on the mobile robot. It is true that position of the z varies very little. However, we found that localizing the mobile node on the three dimensional space using range measurement data is very weak to noise. In more detail, let assume that 4 anchor node are placed on the ground and form square with similar height. Let x_i, y_i, z_i , and d_i be the position of i^{th} anchor node and range measurement. Then equations on the 3-D space are as follows:

$$(x - x_1)^2 + (y - y_1)^2 + (z - z_1)^2 = d_1^2 \quad (16)$$

$$(x - x_2)^2 + (y - y_2)^2 + (z - z_2)^2 = d_2^2 \quad (17)$$

$$(x - x_3)^2 + (y - y_3)^2 + (z - z_3)^2 = d_3^2 \quad (18)$$

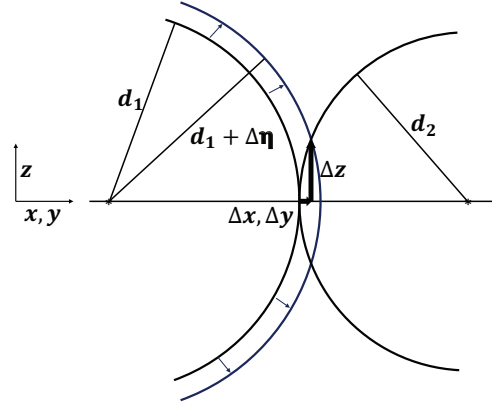


FIGURE 5: Figures from experiment (a)The anchor and tag nodes (b)Four examples of the trajectory (c) the process that makes dataset

$$(x - x_4)^2 + (y - y_4)^2 + (z - z_4)^2 = d_4^2 \quad (19)$$

where x, y , and z is the unknown position of the mobile node. And we can rewrite these equation by subtracting (19) from (16), (17), and (18)

$$A_{3D} X_{3D} = b_{3D} \quad (20)$$

where X_{3D} indicates $[x, y, z]^T$ and A_{3D} and b_{3D} are as follows:

$$A_{3D} = \begin{bmatrix} 2(x_2 - x_1) & 2(y_2 - y_1) & 2(z_2 - z_1) \\ 2(x_3 - x_1) & 2(y_3 - y_1) & 2(z_3 - z_1) \\ 2(x_4 - x_1) & 2(y_4 - y_1) & 2(z_4 - z_1) \end{bmatrix} \quad (21)$$

$$b_{3D} = \begin{bmatrix} (d_1^2 - d_2^2) - (x_1^2 - x_2^2) - (y_1^2 - y_2^2) - (z_1^2 - z_2^2) \\ (d_1^2 - d_3^2) - (x_1^2 - x_3^2) - (y_1^2 - y_3^2) - (z_1^2 - z_3^2) \\ (d_1^2 - d_4^2) - (x_1^2 - x_4^2) - (y_1^2 - y_4^2) - (z_1^2 - z_4^2) \end{bmatrix} \quad (22)$$

Unlike the case of 2D, A_{3D} consists of z components on the third column. Anchor nodes are placed in a less scattered on the z direction than x and y axis. In other words, it is hard to put the anchor nodes with exactly same height in real-world, which means $z_1 \approx z_2 \approx z_3 \approx z_4$. Consequently, values on the third column of A_{3D} converge to zero. Let assume that A_{3D} be full rank in such a way as to $\exists A_{3D}^{-1}$, then values of third row of A_{3D}^{-1} relatively have considerable numbers. As a result, this make z value very unstable in such a way as to cause huge error with respect to z adirection.

IV. EXPERIMENTS

A. EXPERIMENTAL ENVIRONMENT

Our experimental system consists of a UWB (ultra wideband) sensor tag node attached on the mobile robot platform and eight anchor nodes that have a UWB transceiver, the motion capture system with 12 cameras, and a SFF (small-form-factor) computer.

UWB sensor anchors are attached to landmarks. These become the end points of the range measurements. The anchor nodes transmit the UWB signal. A UWB sensor

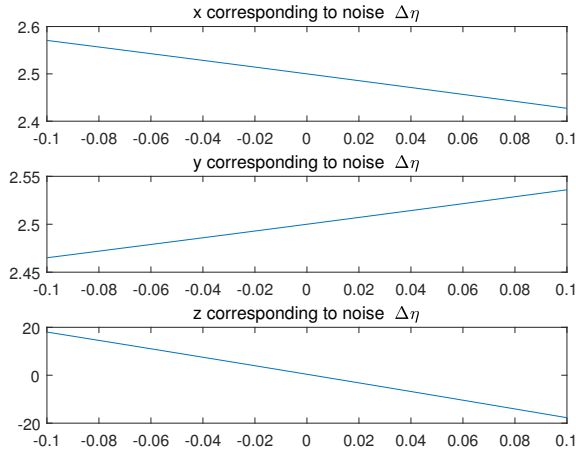


FIGURE 6: $(-0.01, 0.001, 0.24)$, $(5.02, 0.01, 0.2)$, $(5, 5.01, 0.21)$, $(0.01, 4.999, 0.23)$ true: $(2.5, 2.5, 0.4)$

tag is attached to a robot. It becomes the opposite side end point of the measurements. The tag node receives the signal and measures the range between two devices. Each UWB transceiver, DW1000 UWB-chip made by Decawave, supports 6 RF bands from 3.5 GHz to 6.5 GHz. It measures in centimeter-level accuracy. Fig. 7 shows anchor and tag nodes.

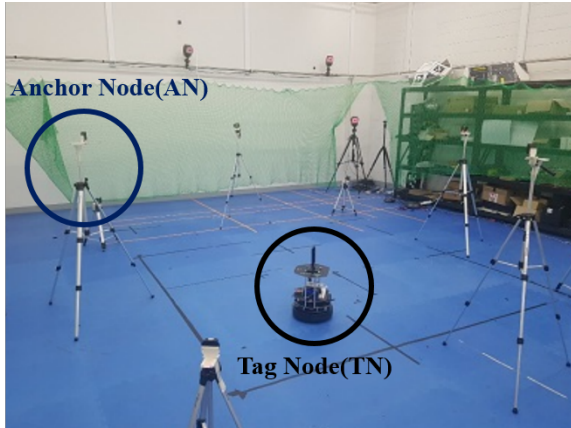


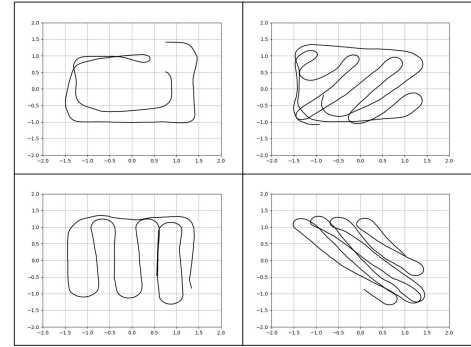
FIGURE 7: The anchor and tag nodes

We infer the position of a robot with our network. To train the network and test the results, the ground truths are needed. We get the ground truth by using the motion capture system. The system is Eagle Digital Realtime system of motion analysis corporation that operates with the principle of stereo pattern recognition that is a kind of photogrammetry based on the epipolar geometry and the triangulation methodology. We attach four markers to a robot. The system gives us the location of these markers and has $< 1\text{mm}$ accuracy.

A mobile robot used in experiment is iCleo Kobuki from Yujinrobot that has 70 cm/s maximum velocity. The SFF computer is a gigabyte Ultra compact PC. Deep learning

framework used for our network is pytorch 0.4.0 on python 3.6. The network inferences on the same setting.

The UWB tag is attached to mobile robot that has a SFF computer. The UWB anchors are attached to stands that have two different heights. The anchors are positioned randomly in the square space where the motion capture system can measure the position. The robot moves in this space also. As you can see in Fig. 8, a mobile robot manually goes on various random trajectories by experimenters.



• • •

FIGURE 8: Four example of the trajectories

B. DATA SYNCHRONIZATION FOR TRAIN/TEST DATA

During the robot is going on, the data is saved in the computer. The distance data used for input data is measured by the UWB sensors. The global position data used for ground truth is measured by the motion capture system. These two kinds of data are paired in a dataset. Because these two data are transmitted by different frequency, we need to synchronize these. This process is conducted in a SFF computer. The computer receives these two kinds of data respectively and synchronizes these by time. To synchronize, we make an independent thread that concatenates and saves these data at the same time. The data is saved at 20Hz frequency. Each trajectory becomes one dataset. All the trajectories are different. Fig. 9 shows this process. After collecting whole datasets, we separate the entire dataset to two types, some are the training datasets and others are test datasets.

Ground truth data is robot's position measured by eagle eye motion capturer, whose error is in mm units.

C. TRAINING THE NETWORKS

Pass

V. RESULTS

To improve the performance of our proposed networks, we check which sequence length is optimal for localization of MN. And then, we implemented and trained other previous networks to compare their performance when our real-world data are given as input. We set three test trajectory cases:

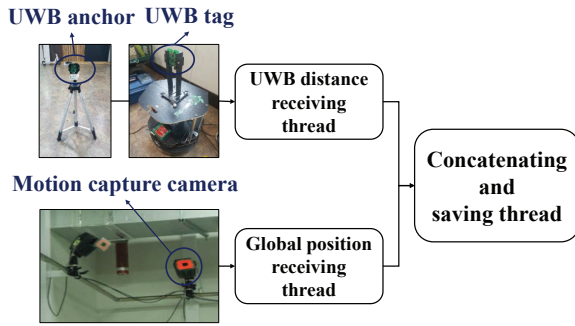


FIGURE 9: the process that makes dataset

a square path, a vertical and horizontal winding path, and a diagonal winding path.

A. PERFORMANCE ACCORDING TO THE SEQUENCE LENGTH

Even though LSTM solve the long term dependency problem, we thought that it would not always make neural network perform better as sequence length become more longer. So, as part of turning hyperparameters, we modified sequence length of our network in a variety of numbers to find optimal sequence length.

As illustrated in FIGURE. 8, The train data are gathered by the robot that arbitrary moves on the region where the motion capture camera cover. The neural network takes distances measured by each AN and a MN as input and outputs robot's position. The Root-Mean-Squared Error (RMSE) results of trajectory prediction with respect to sequence length are shown in FIGURE. 10 and specific figures are shown in Table 2.

As a result, we found that there is a trade-off between robustness and improvement of general performance according to the sequence length. Figure. 10 show that as longer the sequence length, as more inaccurate the mean performance. This is because it is hard for the neural architecture to train the characteristics of the sequential data. In other words, as sequence length become longer, the tendency of the values may be different due to accumulation of different patterns of system noises even if the data is acquired by moving to a similar path in the train data. Note that network with longer sequence length tend to have less error variance and more a ability to generalize the situation since they could utilize more extended temporal information. By doing so, the neural network have the ability to suppress the disturbance caused by noises.

Similarly, as sequence length become shorter, it is less difficult for the neural architecture to understand the characteristics of the sequential data by virtue of the less accumulation of different patterns of system noises. Note that in most cases, the network constructed by the shorter sequential length

tend to estimate the position more precisely than those with longer sequential lengths. However, it becomes a double-edged sword because temporal information is also reduced accordingly in such a way as to week to the noises, having larger error variance.

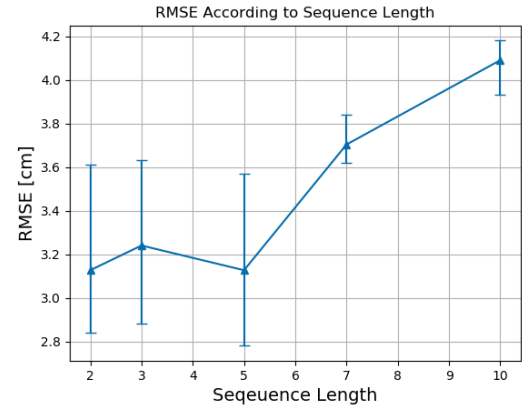


FIGURE 10: Error bar graph of RMSE with respect to sequence length

TABLE 2: Root mean squared error of each case

The results of RMSE[cm]					
	Sequence length				
	2	3	5	7	10
Test1	2.84	2.88	2.78	3.65	4.16
Test2	3.61	3.63	3.57	3.84	4.18
Test3	2.93	3.21	3.03	3.62	3.93

B. PERFORMANCE COMPARISON RESULT

The results of trajectory prediction are shown in Fig. 3(a) and Fig. 3(d) and Root-Mean-Squared Error (RMSE) are shown in Table 1. Performance is better in order of stacked Bi-LSTM, Bi-LSTM, LSTM and GRU. In case of GRU, it has only two gates which is less complex structure than LSTM [27]. However, due to GRU's less complexity, GRU has less number of neurons than LSTM so their non-linear mapping achieves less performance. Likewise, Bi-LSTM consists of two LSTMs to process sequence in two directions so that infer output using the correlation of the backward information and the forward information of the sequences of each time step with its two separate hidden layers. Thus, Bi-LSTM has better nonlinear mapping capability than LSTM. For similar reasons, stacked Bi-LSTM is the architecture that stacks two Bi-LSTMs, so inference performance is better than Bi-LSTM. As a result, the stacked Bi-LSTM showed the best performance among unit RNN architectures. Therefore, we can conclude that the performance improves as the non-linearity of the architecture increases.

!!!!!!!!!!!!!! We also verified effectiveness of attention layer. It was confirmed that the performance of the networks with the attention layer is improved compared to

the networks without the attention layer. We also provide statistical analysis from simulations demonstrating that our new approach can cope with highly noisy sensors and reduces in one order of magnitude the average errors of the method proposed !!!!!!!!!!!!!!!

VI. CONCLUSION

In this paper, we proposed a novel approach to range-only SLAM using multimodal-based RNN models and tested our architectures in two test data.

Using deep learning, our structure directly learns the end-to-end mapping between distance data and robot position. The multimodal bidirectional stacked LSTM structure exhibits the precise estimates of robot positions. We set two test trajectory cases: an square path and zigzag path. The results shows that it has better performance than established probabilistic-based approach. In both cases, performance of our networks is better that of particle filter. RMSE of our networks in test1 is 3.928cm and 4.119cm in test2. Therefore, we could check the possibility that our multimodal LSTM-based structure can substitute traditional algorithms

As a future work, because we conducted on just localization, this approach may not be operated when locations of sensors are changed. Therefore, the proposed method needs to be revised for precise estimates even though locations of anchors are changed.

Appendixes, if needed, appear before the acknowledgment.

ACKNOWLEDGMENT

The preferred spelling of the word "acknowledgment" in American English is without an "e" after the "g." Use the singular heading even if you have many acknowledgments. Avoid expressions such as "One of us (S.B.A.) would like to thank" Instead, write "F. A. Author thanks" In most cases, sponsor and financial support acknowledgments are placed in the unnumbered footnote on the first page, not here.

REFERENCES

- [1] S. Li, A. Zhan, X. Wu, and G. Chen, "Ern: Emergence rescue navigation with wireless sensor networks," in 2009 15th International Conference on Parallel and Distributed Systems. IEEE, 2009, pp. 361–368.
- [2] K. K. Khedo, R. Perseedoss, A. Mungur, et al., "A wireless sensor network air pollution monitoring system," arXiv preprint arXiv:1005.1737, 2010.
- [3] J. Zhang, G. Song, G. Qiao, T. Meng, and H. Sun, "An indoor security system with a jumping robot as the surveillance terminal," IEEE Transactions on Consumer Electronics, vol. 57, no. 4, 2011.
- [4] A. Kulaib, R. Shubair, M. Al-Qutayri, and J. W. Ng, "An overview of localization techniques for wireless sensor networks," in Innovations in Information Technology (IIT), 2011 International Conference on. IEEE, 2011, pp. 167–172.
- [5] L. Peneda, A. Azenha, and A. Carvalho, "Trilateration for indoors positioning within the framework of wireless communications," in Industrial Electronics, 2009. IECON'09. 35th Annual Conference of IEEE. IEEE, 2009, pp. 2732–2737.
- [6] J. Jung and H. Myung, "Indoor localization using particle filter and map-based nlos ranging model," in Robotics and Automation (ICRA), 2011 IEEE International Conference on. IEEE, 2011, pp. 5185–5190.
- [7] R. Kaune, "Accuracy studies for tdoa and toa localization," in Information Fusion (FUSION), 2012 15th International Conference on. IEEE, 2012, pp. 408–415.
- [8] P. Singh and S. Agrawal, "Tdoa based node localization in wsn using neural networks," in Communication Systems and Network Technologies (CSNT), 2013 International Conference on. IEEE, 2013, pp. 400–404.
- [9] K. Doğançay and H. Hmam, "Optimal angular sensor separation for aoa localization," Signal Processing, vol. 88, no. 5, pp. 1248–1260, 2008.
- [10] S. S. Banihashemian, F. Adibnia, and M. A. Sarram, "A new range-free and storage-efficient localization algorithm using neural networks in wireless sensor networks," Wireless Personal Communications, vol. 98, no. 1, pp. 1547–1568, 2018.
- [11] J. Blumenthal, R. Grossmann, F. Golatowski, and D. Timmermann, "Weighted centroid localization in zigbee-based sensor networks," in Intelligent Signal Processing, 2007. WISP 2007. IEEE International Symposium on. IEEE, 2007, pp. 1–6.
- [12] T. Han, X. Lu, and Q. Lan, "Pattern recognition based kalman filter for indoor localization using tdoa algorithm," Applied Mathematical Modelling, vol. 34, no. 10, pp. 2893–2900, 2010.
- [13] J. Li, X. Yue, J. Chen, and F. Deng, "A novel robust trilateration method applied to ultra-wide bandwidth location systems," Sensors, vol. 17, no. 4, p. 795, 2017.
- [14] F. R. Fabresse, F. Caballero, I. Maza, and A. Ollero, "An efficient approach for undelayed range-only slam based on gaussian mixtures expectation," Robotics and Autonomous Systems, vol. 104, pp. 40–55, 2018.
- [15] F. Caballero, L. Merino, I. Maza, and A. Ollero, "A particle filtering method for wireless sensor network localization with an aerial robot beacon," in Robotics and Automation, 2008. ICRA 2008. IEEE International Conference on. IEEE, 2008, pp. 596–601.
- [16] D. A. Tran and T. Nguyen, "Localization in wireless sensor networks based on support vector machines," IEEE Transactions on Parallel and Distributed Systems, vol. 19, no. 7, pp. 981–994, 2008.
- [17] R. Huan, Q. Chen, K. Mao, and Y. Pan, "A three-dimension localization algorithm for wireless sensor network nodes based on svm," in Green Circuits and Systems (ICGCS), 2010 International Conference on. IEEE, 2010, pp. 651–654.
- [18] V.-s. Feng and S. Y. Chang, "Determination of wireless networks parameters through parallel hierarchical support vector machines," IEEE Transactions on Parallel and Distributed Systems, vol. 23, no. 3, pp. 505–512, 2012.
- [19] S. Afzal and H. Beigy, "A localization algorithm for large scale mobile wireless sensor networks: a learning approach," The Journal of Supercomputing, vol. 69, no. 1, pp. 98–120, 2014.
- [20] J. Lee, W. Chung, and E. Kim, "A new kernelized approach to wireless sensor network localization," Information Sciences, vol. 243, pp. 20–38, 2013.
- [21] J. Lee, B. Choi, and E. Kim, "Novel range-free localization based on multidimensional support vector regression trained in the primal space," IEEE transactions on neural networks and learning systems, vol. 24, no. 7, pp. 1099–1113, 2013.
- [22] Y. LeCun, Y. Bengio, and G. Hinton, "Deep learning," nature, vol. 521, no. 7553, p. 436, 2015.
- [23] M. S. Rahman, Y. Park, and K.-D. Kim, "Localization of wireless sensor network using artificial neural network," in Communications and Information Technology, 2009. ISCIT 2009. 9th International Symposium on. IEEE, 2009, pp. 639–642.
- [24] M. Abdelhadi, M. Anan, and M. Ayyash, "Efficient artificial intelligent-based localization algorithm for wireless sensor networks," Journal of Selected Areas in Telecommunications, vol. 3, no. 5, pp. 10–18, 2013.
- [25] S. Kumar, R. Sharma, and E. Vans, "Localization for wireless sensor networks: A neural network approach," arXiv preprint arXiv:1610.04494, 2016.
- [26] A. Chatterjee, "A fletcher-reeves conjugate gradient neural-network-based localization algorithm for wireless sensor networks," IEEE transactions on vehicular technology, vol. 59, no. 2, pp. 823–830, 2010.
- [27] R. Samadian and S. M. Noorhosseini, "Probabilistic support vector machine localization in wireless sensor networks," ETRI Journal, vol. 33, no. 6, pp. 924–934, 2011.
- [28] S. K. Chenna, Y. K. Jain, H. Kapoor, R. S. Bapi, N. Yadaiah, A. Negi, V. S. Rao, and B. L. Deekshatulu, "State estimation and tracking problems: A comparison between kalman filter and recurrent neural networks," in International Conference on Neural Information Processing. Springer, 2004, pp. 275–281.

Method	Abs. Error				Std. Error				RMSE	Min. Error	Max. Error
	x	y	z	total	x	y	z	total			
Test1: Square path											
RBF [29]											
MLP [23]											
MLBPN [8]											
MLP [24]											
MLP [25]											
LPSONN [10]											
Ours w/o attn.											
Ours											
Test2: Vertical winding path											
RBF [29]											
MLP [23]											
MLBPN [8]											
MLP [24]											
MLP [25]											
LPSONN [10]											
Ours w/o attn.											
Ours											
Test3: Diagonal winding path											
RBF [29]											
MLP [23]											
MLBPN [8]											
MLP [24]											
MLP [25]											
LPSONN [10]											
Ours w/o attn.											
Ours											

- [29] A. Shareef, Y. Zhu, and M. Musavi, "Localization using neural networks in wireless sensor networks," in Proceedings of the 1st international conference on MOBILE Wireless MiddleWARE, Operating Systems, and Applications, ICST (Institute for Computer Sciences, Social-Informatics and ...), 2008, p. 4.
- [30] M. Bernas and B. Placzek, "Fully connected neural networks ensemble with signal strength clustering for indoor localization in wireless sensor networks," International Journal of Distributed Sensor Networks, vol. 11, no. 12, p. 403242, 2015.
- [31] S. Hochreiter and J. Schmidhuber, "Long short-term memory," Neural computation, vol. 9, no. 8, pp. 1735–1780, 1997.
- [32] M. Schuster and K. K. Paliwal, "Bidirectional recurrent neural networks," IEEE Transactions on Signal Processing, vol. 45, no. 11, pp. 2673–2681, 1997.
- [33] F. Zhang, C. Hu, Q. Yin, W. Li, H.-C. Li, and W. Hong, "Multi-aspect-aware bidirectional lstm networks for synthetic aperture radar target recognition," IEEE Access, vol. 5, pp. 26 880–26 891, 2017.
- [34] W. Li, W. Nie, and Y. Su, "Human action recognition based on selected spatio-temporal features via bidirectional lstm," IEEE Access, vol. 6, pp. 44 211–44 220, 2018.
- [35] A. Ullah, J. Ahmad, K. Muhammad, M. Sajjad, and S. W. Baik, "Action recognition in video sequences using deep bi-directional lstm with cnn features," IEEE Access, vol. 6, pp. 1155–1166, 2018.
- [36] F. R. Fabresse, F. Caballero, I. Maza, and A. Ollero, "Undelayed 3d roslam based on gaussian-mixture and reduced spherical parametrization," in Intelligent Robots and Systems (IROS), 2013 IEEE/RSJ International Conference on. Citeseer, 2013, pp. 1555–1561.
- [37] K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image recognition," arXiv preprint arXiv:1409.1556, 2014.
- [38] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in Proceedings of the IEEE conference on computer vision and pattern recognition, 2016, pp. 770–778.
- [39] A. Graves, N. Jaitly, and A.-r. Mohamed, "Hybrid speech recognition with deep bidirectional lstm," in Automatic Speech Recognition and Understanding (ASRU), 2013 IEEE Workshop on. IEEE, 2013, pp. 273–278.
- [40] A. Graves, A.-r. Mohamed, and G. Hinton, "Speech recognition with deep recurrent neural networks," in Acoustics, speech and signal processing (icassp), 2013 IEEE international conference on. IEEE, 2013, pp. 6645–6649.
- [41] R. Pascanu, T. Mikolov, and Y. Bengio, "On the difficulty of training recurrent neural networks," in International Conference on Machine Learning, 2013, pp. 1310–1318.
- [42] V. Nair and G. E. Hinton, "Rectified linear units improve restricted boltzmann machines," in Proceedings of the 27th international conference on machine learning (ICML-10), 2010, pp. 807–814.
- [43] M.-T. Luong, H. Pham, and C. D. Manning, "Effective approaches to attention-based neural machine translation," arXiv preprint arXiv:1508.04025, 2015.
- [44] M. Jaderberg, K. Simonyan, A. Zisserman, et al., "Spatial transformer networks," in Advances in neural information processing systems, 2015, pp. 2017–2025.
- [45] E. Parisotto, D. S. Chaplot, J. Zhang, and R. Salakhutdinov, "Global pose estimation with an attention-based recurrent network," arXiv preprint arXiv:1802.06857, 2018.
- [46] X. Zhu, J. Dai, L. Yuan, and Y. Wei, "Towards high performance video object detection," in Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 2018, pp. 7210–7218.
- [47] J. Xu, T. Yao, Y. Zhang, and T. Mei, "Learning multimodal attention lstm networks for video captioning," in Proceedings of the 2017 ACM on Multimedia Conference. ACM, 2017, pp. 537–545.
- [48] F. Wang, M. Jiang, C. Qian, S. Yang, C. Li, H. Zhang, X. Wang, and X. Tang, "Residual attention network for image classification," arXiv preprint arXiv:1704.06904, 2017.



HYUNGTAE LIM received the B.S. degree in the Mechanical Engineering from the Korea Advanced Institute of Science and Technology, where he is currently pursuing the M.S. degree with the robotics program and civil engineering. His research includes computer vision, localization using raw-level sensors, SLAM, and deep learning.



CHANGGYU PARK received the B.S. degree in the Civil Engineering from the Yonsei University and Technology and now he is on the M.S. degree with the the Civil Engineeering from the Korea Advanced Institute of Science and Technology. His research includes deep learning, optimization, and computer vision.



HYUN MYUNG received the B.S., M.S., and Ph.D. degrees from the Korea Advanced Institute of Science and Technology (KAIST), Daejeon, South Korea, in 1992, 1994, and 1998, respectively, all in electrical engineering. He was a Senior Researcher with the Electronics and Telecommunications Research Institute, Daejeon, from 1998 to 2002, a CTO and the Director with the Digital Contents Research Laboratory, Emersys Corporation, Daejeon, from 2002 to 2003, and a

Principle Researcher with the Samsung Advanced Institute of Technology, Yongin, South Korea, from 2003 to 2008. Since 2008, he has been an Associate Professor with the Department of Civil and Environmental Engineering, KAIST, where he is currently an Adjunct Professor in robotics program. His current research interests include structural health monitoring using robotics, soft computing, simultaneous localization and mapping, robot navigation, machine learning, deep learning, and swarm robot.

...