

Improving CSI-based Massive MIMO Indoor Positioning using Convolutional Neural Network

Gregor Cerar^{*†}, Aleš Švigelj^{*†}, Mihael Mohorčič^{*†}, Carolina Fortuna^{*}, and Tomaž Javornik^{*†}

^{*}Department of Communication Systems, Jožef Stefan Institute, SI-1000, Slovenia.

[†]Jožef Stefan International Postgraduate School, SI-1000, Slovenia.

{gregor.cerar | ales.svigelj | miha.mohorcic | carolina.fortuna | tomaz.javornik}@ijs.si

Abstract—Multiple-input multiple-output (MIMO) is an enabling technology to meet the growing demand for faster and more reliable communications in wireless networks with a large number of terminals, but it can also be applied for position estimation of a terminal exploiting multipath propagation from multiple antennas. In this paper, we investigate new convolutional neural network (CNN) structures for exploiting MIMO-based channel state information (CSI) to improve indoor positioning. We evaluate and compare the performance of three variants of the proposed CNN structure to five NN structures proposed in the scientific literature using the same sets of training-evaluation data. The results demonstrate that the proposed residual convolutional NN structure improves the accuracy of position estimation and keeps the total number of weights lower than the published NN structures. The proposed CNN structure yields from 2 cm to 10 cm better position accuracy than known NN structures used as a reference.

Index Terms—MIMO, wireless, localization, deep learning, neural network (NN), position estimation, residual networks, convolutional networks, fingerprinting

I. INTRODUCTION

Multiple-input multiple-output (MIMO) is an enabling technology to meet the growing demand for faster and more reliable wireless communications in wireless networks with a vast number of wireless terminals. The idea is to introduce space diversity through an increased number of antennas at the transmitter and/or receiver side, thus improving the wireless link reliability or increasing radio link capacity by exploiting multipath propagation.

The MIMO approach is already part of WiFi devices based on IEEE 802.11n standard. It is also extensively used in the third-generation (3G) and the fourth generation (4G) mobile broadband networks. Furthermore, the fifth-generation (5G) and future sixth-generation (6G) broadband networks are putting even more emphasis on multi-antenna technologies extending the concept to massive MIMO.

In addition to improving the reliability and capacity of wireless links, the multiple antennas can be exploited for indoor or outdoor positioning. Accurate indoor positioning (and position estimation in general) is a highly desirable feature of future wireless networks [1], [2]. It is a key enabler for a wide range of applications, including navigation, smart factories and cities, surveillance, security, IoT, sensor networks, and future reconfigurable intelligent surfaces (RIS). Additionally, indoor positioning can be leveraged for improved beamforming and channel estimation in wireless communications.

The positioning techniques can be classified into five classes; namely, i) proximity-based, ii) angle-based, iii) range-based, iv) fingerprinting-based and v) device-free localization [3]–[5]. The *proximity-based* approaches rely on the information about which objects are detected in the observer vicinity. The *angle-based* approaches rely on the angle-of-arrival (AoA) of the received signal as obtained by the multi-antenna system. The *range-based* approaches rely on either time-of-arrival (ToA) or time-difference-of-arrival (TDoA) of the received signal or on received signal strength (RSS). Since the ToA approach requires large bandwidth (*e.g.* 20 MHz correspond to 15 m accuracy), it is not always feasible. Next, the *fingerprinting-based* approaches rely on the access to accurate channel state information (CSI), where the radio frequency (RF) fingerprints consist of signal measurements obtained at known positions within the deployment area. Fingerprint positioning can work with a single base station (BS), if CSI includes spatial channel information obtained from several antennas. Since the RF fingerprints typically form high-dimensional datasets, the fingerprinting-based localization is a good candidate for utilizing machine learning (ML) methods. Finally, in the *device-free* localization, the system detects and tracks any entity based on an entity’s impact on CSI.

The fingerprinting-based and the device-free localization classes rely on the CSI datasets that can be obtained by extensive measurement campaigns, which prove to be time consuming and expensive. Alternatively, the analytical and simulation-based approaches (*e.g.* radio ray-tracing [6]) can supplement measurements to some extent. The simulated CSI accuracy, however, depends on the span of considered phenomena (see [2, p. 22]) and knowledge about a particular radio environment.

In this paper, motivated by CTW 2019 challenge [7], we investigate new machine learning approaches for improving the indoor positioning using CSI obtained by a single massive MIMO antenna. In particular, we propose a new convolutional neural network (CNN) structure with three variants in its internal layers, we compare its performance to some of the existing neural network (NN) structures on the same CSI dataset, and discuss possible further improvements.

The main contributions of this paper are the following:

- the design of a new deep CNN structure able to accurately estimate the position of a transmitter in a room,

- the re-implementation of the NN structures published in scientific literature and comparison of their performance to the newl proposed CNN structure and
- the evaluation procedure with four different distributions of training/evaluation samples based on publicly available CSI dataset from the CTW 2019 challenge [7].

The rest of the paper is structured as follows. In Section II, we briefly analyze some representative studies on deep learning application for position estimation and provide the problem definition. Section III describes the publicly available CTW 2019 challenge dataset, and Section IV outlines the new proposed CNN structure with its variants. Performance evaluation and comparison to five representative NNs is done in Section V. Discussion on lessons learned and conclusions are provided in Section VI.

II. RELATED WORK AND PROBLEM DEFINITION

The use of ML methods for indoor positioning gained much attention from the research community in recent years, so we selected just a few studies that are relevant and representative to our approach. Savic and Larsson in [3] focus on existing methods for position estimation using classical ML approaches, where they briefly present k-Nearest Neighbour, Support Vector Machine and Gaussian Process Regression as suitable candidates. Recent fingerprint-based positioning studies [8]–[10] consider larger antenna array and high number of considered subcarriers. A high number of antennas and subcarriers inevitably produce rich fingerprint samples and thus large final dataset size. Huge datasets make traditional ML approaches difficult to utilize, opening an opportunity to more recent ML approaches such as deep NNs that utilize either batch processing or stream processing on large datasets.

In the literature, we found several NN architectures proposed for more accurate position estimation. In [4], [8], authors experiment with fully (*i.e.* densely) connected NNs (FCNN). However, feature extraction process using convolutional NNs (CNNs) proved far more effective from the perspective of the performance and the number of weights compared to FCNN [4], [9]–[13]. The top-performing NNs in related work use similar sequence of layers: a convolutional layer, an activation layer and an average pooling layer.

The main application of convolution layers is in image recognition domain. Images typically have sharp edges between surfaces, and from the gradient perspective they are much more “dynamic” compared to CSI, where changes are much slower as shown in Figure 2 for the CTW2019 dataset. In the NN design, image-related tasks typically utilise convolutional layers with kernel shape (n, n) , while a significant part of related work for position estimation utilises kernel shape $(1, n)$.

Concerning position estimation accuracy, [8]–[10] show that an increasing number of antennas and therefore an amount of CSI improves the overall accuracy of position estimation. Furthermore, [10] suggests generating additional features derived from raw CSI, in particular “time series” from inverse FFT and representation with polar coordinates, but in our study we

did not obtain any performance gains when using additional features. Also, in [10] authors show in their experiment that linear and rectangular antenna arrays perform slightly better in terms of accuracy than distributed antenna arrays around the testing area. However, their experiment was limited to a single plane, where height was unchanged, and event recording was done solely in front of the antenna arrays.

Several metrics that serve either as loss metrics at the training phase or as evaluation metrics are applied for performance evaluation. Most often used metrics in related work are mean distance error (MDE) (1), normalized mean distance error (NMDE) (2) and root-mean-squared-error (RMSE) (3). MDE gives the Euclidean distance error between the ground truth position p and the estimated position \hat{p} . In (1), $\|\cdot\|_2$ thus stands for the Euclidean norm. NMDE is used because samples farther away from the antenna array are hard to estimate accurately. Thus, by normalizing, farther away samples receive less penalty for the error.

$$\text{MDE} = \mathbb{E} [\|p - \hat{p}\|_2], \quad (1)$$

$$\text{NMDE} = \mathbb{E} \left[\frac{\|p - \hat{p}\|_2}{\|p\|_2} \right]. \quad (2)$$

$$\text{RMSE} = \sqrt{\mathbb{E} [\|p - \hat{p}\|_2^2]}. \quad (3)$$

In our study, the task of the applied ML approach was to predict the transmitter’s location ($p = p(x, y, z)$) from the publicly available offline dataset [14] containing CSI measurements at different positions. The proposed CNN structures were trained, to gain experience, and evaluated, to estimate the performance, on different non-overlapping subsets of the CSI dataset.

III. DATASET DESCRIPTION

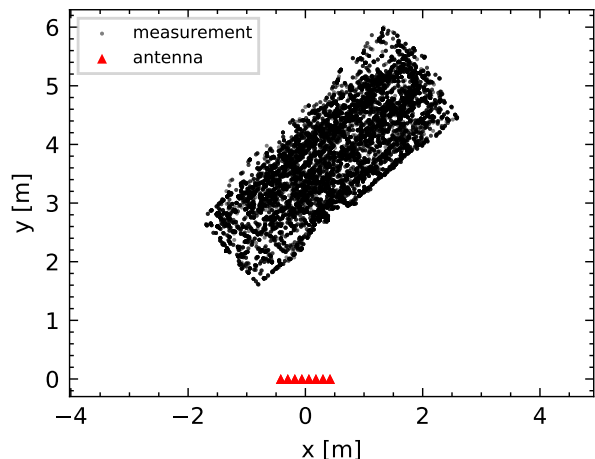


Figure 1: Top-down view on 17,486 samples and antenna orientation.

As explained in Section II, for training of the newly proposed CNN and for the performance evaluation presented in Section V we used the openly available dataset from

CTW 2019 challenge [14]. The dataset was acquired by a massive MIMO channel sounder [9] in a setup visually depicted in [7]. The CSI was measured between a moving transmitter and 8×2 antenna array. The transmitter implemented on SDR was placed on a vacuum-cleaner robot. The robot drove in a random path on approximately $4 \text{ m} \times 2 \text{ m}$ size table. The transmitted signal consisted of OFDM pilots with a bandwidth of 20 MHz and 1024 subcarriers at the central frequency 1.25 GHz. 100 sub-carriers were used as guard bands, 50 on each side of the frequency band. The measurement setup is summarized in Table I, while Figure 1 depicts the top-down view on the the antenna orientation and positions of 17.486 CSI samples available in the dataset.

Table I: Dataset Summary

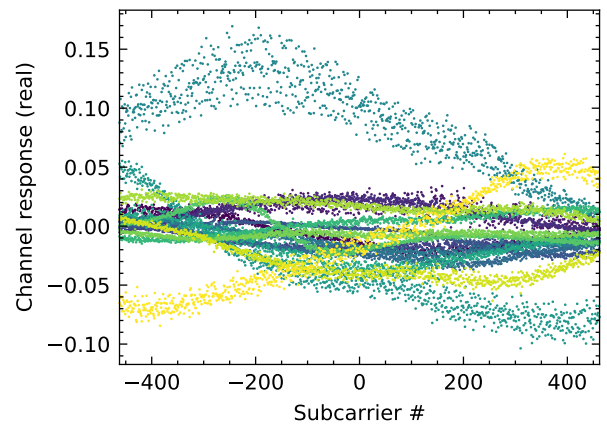
Property	Value
No. antennas	16
Central frequency (f_c)	1.25 GHz
Bandwidth	20 MHz
No. subcarriers	924 useful; ≈ 20 kHz in between
CSI data shape	16 (ant.) \times 924 (subc.) \times 2 (Re/Im)
SNR information	16 (ant.) \times 1 SNR value given in dB
No. all samples	17486 samples

As an example, the CSI sample #130 is presented in Figure 2, namely, real part in Figure 2a and imaginary part in Figure 2b. Both figures present sub-carriers ordered linearly (by their frequency) from left to right, where the value on x-axis present padding (approx. 20 kHz) from the central frequency (1.25 GHz). Each of the seemingly continuous curves (despite slight value deviations) distinguished by colour presents individual channel responses obtained by one of the 16 antennas.

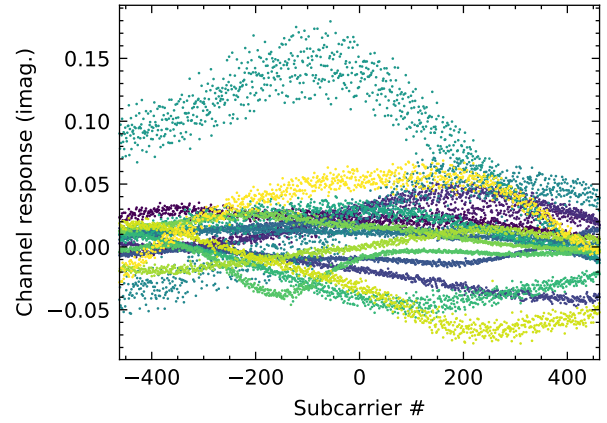
Since we are dealing with a limited amount of data, we divided the dataset to non-overlapping training and evaluation subsets. Moreover, to investigate the influence of training data selection, we generated four different training/evaluation sets from the same original dataset using (i) a random dataset split with overlapping training and evaluation areas, (ii) a long-edge dataset split with narrow evaluation area, (iii) a short-edge dataset split with wide evaluation area, and (iv) a cut-out dataset split with evaluation area within the training area. These four approaches are graphically depicted in Figures 3a, b, c and d, respectively. All approaches split the dataset to training and evaluation subsets approximately at a 9:1 ratio.

IV. PROPOSED NEURAL NETWORK STRUCTURES

In this paper, we propose three CNN structures with multiple layers designed for position estimation based on CSI samples. Unless otherwise stated, after each layer except the last one, we utilize a rectified linear unit (ReLU) activation. The first structure denoted as CNN4 uses four sequential convolutional layers with the kernel shape (1, 7) and stride (1, 3), which are followed by a dense layer with 1000 units. With each convolutional layer, the number of filters (depth) is increased by 50%.



(a) Real part of channel response



(b) Imaginary part channel response

Figure 2: Channel State Information of sample #130

The second structure denoted as CNN4R is similar to CNN4, but instead of four convolutional layers, it utilizes four residual network (ResNet) blocks inspired by ResNet-18 [15]. CNN4R's ResNet blocks use kernel size (1, 7), where each block reduces width with stride (1, 3) and internal, pattern with identity connection repeats three times.

The third structure denoted as CNN4S is based on CNN4R. However, instead of the first ResNet block it utilizes a single convolutional layer (*i.e.* stem) with kernel (1, 7) and stride (1, 2), followed by an average pooling layer with a pool size of (1, 4) and stride (1, 2), which functions as a rolling average.

V. PERFORMANCE EVALUATION

In order to benchmark the proposed CNN structures, we also implemented some representative NNs for position estimation from the related work [4], [8]–[10]. The evaluation results are presented in Table II, where RMSE and MDE are expressed in meters, and NMDE is given in percents. In general, the CNN structures exhibit better performance compared to FCNN. An exception is [16], where authors had put significant effort into data pre-processing, but unfortunately we were unable to replicate the authors' results.

Table II: Performance evaluation on CTW-2019 dataset

Approach	Weights [10^6]	Random			Narrow			Wide			Within		
		RMSE	MDE	NMDE	RMSE	MDE	NMDE	RMSE	MDE	NMDE	RMSE	MDE	NMDE
Dummy (linear), FCNN	<0.1	0.724	1.122	25.1	1.055	1.809	51.4	0.878	1.428	28.2	0.441	0.721	15.1
Arnold <i>et al.</i> [8], FCNN	32.3	0.570	0.853	19.4	1.001	1.594	45.0	0.733	1.145	23.3	0.381	0.584	12.3
Arnold <i>et al.</i> [9], CNN	7.6	0.315	0.445	10.0	0.857	1.330	37.7	0.605	0.923	18.6	0.454	0.702	14.8
Bast <i>et al.</i> [10], CNN	0.4	0.722	1.120	25.1	1.110	1.907	54.2	0.828	1.331	26.5	0.377	0.611	12.7
Chin <i>et al.</i> [4] FCNN	123.6	0.563	0.838	19.0	1.007	1.611	45.4	0.726	1.133	23.0	0.365	0.574	12.0
Chin <i>et al.</i> [4] CNN	13.7	0.100	0.093	2.1	0.854	1.326	37.8	0.530	0.808	16.3	0.381	0.620	13.0
CNN4	5.3	0.122	0.149	3.4	0.819	1.286	36.6	0.514	0.787	15.9	0.365	0.552	11.6
CNN4R	10.8	0.113	0.127	2.8	0.776	1.227	34.7	0.539	0.835	16.8	0.351	0.521	11.0
CNN4S	16.3	0.108	0.120	2.7	0.821	1.285	36.5	0.528	0.804	16.2	0.351	0.524	11.1

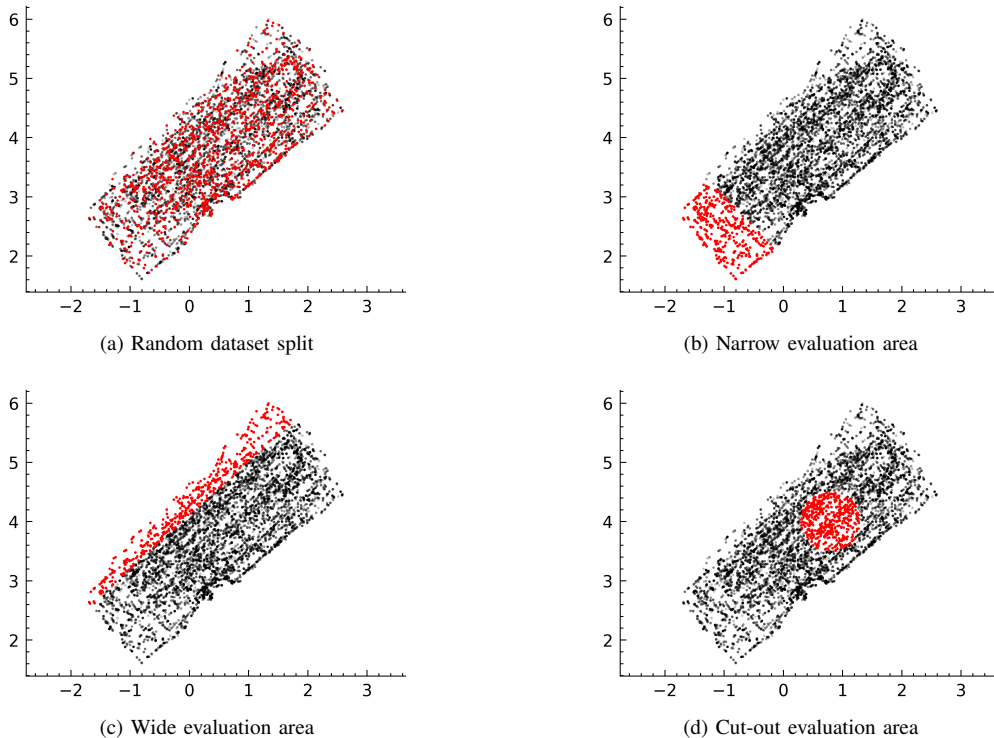


Figure 3: Four types of training/evaluation data subsets, where training data is marked black and testing data is marked red.

Each experiment ran for 250 epochs at most, where we considered early stopping (after 21 epochs without improvement) and exponential decrease of learning rate (*i.e.* $\text{rate} \leftarrow \text{rate}/10$) after 10 epochs without improvement. The batch size was 32, and MDE was used as a training loss metric. Weights were updated using stochastic gradient descent with a learning rate of 10^{-3} and 0.9 momentum.

The results in Table II show that the fingerprinting approach works best when the training set and evaluation set overlap (*i.e.* in random dataset split shown in Fig.3a). Because training and evaluation samples are close to each other, the results do not include error due to unbalanced training. The performance difference between NN structures is also the most significant

for this dataset split. We see that difference between the best and the worst performing model is almost one meter or 23 per cent for NMDE. For this scenario, the best performing model is CNN [4], but our proposed CNN structures show comparable performance. However, in the case of dataset with narrow evaluation area (Fig. 3b), our model performs slightly better than CNN [4]. The difference in MDE is up to 6cm. Even when the models are evaluated by dataset with wide validation area (Fig. 3c), our model performs slightly better than CNN [4], but the difference is only up to 2cm. Additionally, we see that our model performs slightly better when evaluated by dataset with cutout area (Fig. 3d). The estimated distance error difference with respect to CNN [4]

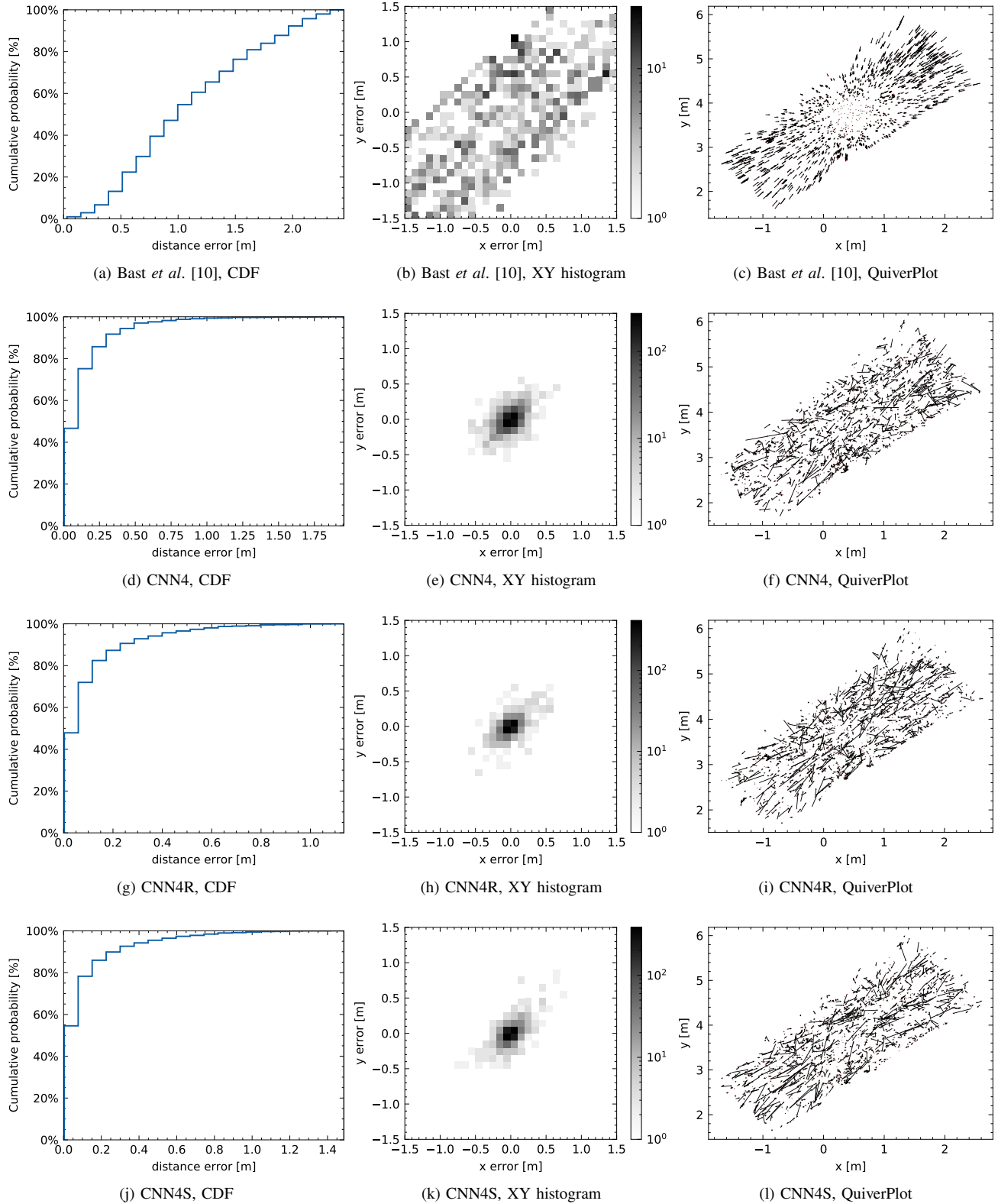


Figure 4: Error distribution for proposed CNN structures and a reference Bast *et al.* [10] structure presented in terms of CDFs of estimation error, histogram of error along X and Y axis, and quiver plot with arrows pointing from ground truth values towards error offset. Note that quivers were automatically scaled.

is up to 3 cm.

In Figure 4, we graphically present the error distribution of Bast *et al.* [10] approach and three newly proposed structures for random train/evaluate scenario. We present position estimation accuracy using three types of figures. The figures in first column present cumulative distribution function (CDF), the second column contains figures with discrete error density over X and Y axis, and in the third column, we depicted quiver plot, where quivers/arrows points from ground truth toward the estimated position in relative scale. Figures 4a and 4b show that NN structure proposed by Bast *et al.* [10] has error distribution far more spread compared to our proposed structures. Figure 4c reveals that Bast *et al.* [10] structure has decent accuracy at the centre of the dataset other points have a bias toward the centre. CNN4R structure shows the highest accuracy, which is also aligned with Table II. The MDE values are below 1.2 m (Fig. 4g). The same can be observed by looking at Figure 4h, where the histogram is narrower compared to the other two proposed models. In addition, analysis of the quiver plots in the third column in Figure 4 shows that there are no inconsistent areas present.

VI. CONCLUSIONS AND FUTURE WORK

In this paper we proposed three new CNN structures referred to as CNN4, CNN4R and CNN4S, designed for improving indoor positioning based on CSI obtained from a single massive MIMO antenna. The performance evaluation of the CNN structures shows that they are within the top ranked or the best performers in all given scenarios for different training/evaluation datasets derived from the publicly available CSI dataset from the CTW 2019 challenge [7]. Even though the performance difference is minimal, or in some cases even in the noise range, we achieve similar performance to NNs [4], [9] with a significantly lower number of trainable weights. Since we utilize strides instead of pooling operation, the computational cost is also lower than that of the comparable models.

The proposed CNN structures as well as those from the related work are not optimal. For our CNNs, we focused solely on pursuing the highest position estimation accuracy. Thus, in the process, we ignored several vital aspects of NNs that may be worth of further investigation, such as dead neurons, extreme weight values and fading/exploding gradient. While dead neurons could be prevented by using different activation functions, they can pose an opportunity for pruning NN, which would reduce the number of weights, making it sparse and consequently decreasing the neural network's size. The extreme weight values can be tackled using regularisation, but it would significantly increase the training time and the number of tunable parameters. To adequately address the fading/exploding gradient, besides using the residual connection, we see great potential in the recently introduced self-normalised NNs. However, their full potential has yet to be explored.

ACKNOWLEDGMENTS

The Slovenian Research Agency supported this work under grants P2-0016 and J2-2507. The authors would like to thank Maximilian Arnold from the University of Stuttgart, one of the CTW dataset creators, for kindly responding to our questions during the initial exploration of the dataset.

REFERENCES

- [1] F. Lemic, J. Martin, C. Yarp, D. Chan, V. Handziski, R. Brodersen, G. Fettweis, A. Wolisz, and J. Wawrzynek, "Localization as a feature of mmwave communication," in *2016 International Wireless Communications and Mobile Computing Conference (IWCMC)*, 2016, pp. 1033–1038.
- [2] F. Wen, H. Wymeersch, B. Peng, W. P. Tay, H. C. So, and D. Yang, "A survey on 5g massive mimo localization," *Digital Signal Processing*, vol. 94, pp. 21 – 28, 2019, special Issue on Source Localization in Massive MIMO. [Online]. Available: <http://www.sciencedirect.com/science/article/pii/S1051200419300569>
- [3] V. Savic and E. G. Larsson, "Fingerprinting-based positioning in distributed massive mimo systems," in *2015 IEEE 82nd Vehicular Technology Conference (VTC2015-Fall)*, 2015, pp. 1–5.
- [4] W. L. Chin, C. C. Hsieh, D. Shiung, and T. Jiang, "Intelligent indoor positioning based on artificial neural networks," *IEEE Network*, vol. 34, no. 6, pp. 164–170, 2020.
- [5] K. Bregar, A. Hrovat, M. Mohorčič, and T. Javornik, "Self-calibrated uwb based device-free indoor localization and activity detection approach," in *2020 European Conference on Networks and Communications (EuCNC)*, 2020, pp. 176–181.
- [6] R. Novak, "Discrete Method of Images for 3D Radio Propagation Modeling," *3D Research*, vol. 7, no. 3, p. 26, Sep. 2016. [Online]. Available: <http://link.springer.com/10.1007/s13319-016-0102-y>
- [7] M. Arnold, S. Dörner, S. Cammerer, and S. Ten Brink. IEEE CTW 2019 – Positioning Algorithm Competition. [Online]. Available: http://attend.ieee.org/ctw-2019/wp-content/uploads/sites/105/2019/02/CTW2019_UserPos_Comp.pdf
- [8] —, "On deep learning-based massive mimo indoor user localization," in *2018 IEEE 19th International Workshop on Signal Processing Advances in Wireless Communications (SPAWC)*, 2018, pp. 1–5.
- [9] M. Arnold, J. Hoydis, and S. t. Brink, "Novel massive mimo channel sounding data applied to deep learning-based indoor positioning," in *SCC 2019; 12th International ITG Conference on Systems, Communications and Coding*, 2019, pp. 1–6.
- [10] S. D. Bast, A. P. Guevara, and S. Pollin, "Csi-based positioning in massive mimo systems using convolutional neural networks," in *2020 IEEE 91st Vehicular Technology Conference (VTC2020-Spring)*, 2020, pp. 1–5.
- [11] J. Vieira, E. Leitinger, M. Sarajlic, X. Li, and F. Tufvesson, "Deep convolutional neural networks for massive mimo fingerprint-based positioning," in *2017 IEEE 28th Annual International Symposium on Personal, Indoor, and Mobile Radio Communications (PIMRC)*, 2017, pp. 1–6.
- [12] M. Arnold, S. Dörner, S. Cammerer, J. Hoydis, and S. ten Brink, "Towards practical fdd massive mimo: Csi extrapolation driven by deep learning and actual channel measurements," in *2019 53rd Asilomar Conference on Signals, Systems, and Computers*, 2019, pp. 1972–1976.
- [13] M. Widmaier, M. Arnold, S. Dörner, S. Cammerer, and S. ten Brink, "Towards practical indoor positioning based on massive mimo systems," in *2019 IEEE 90th Vehicular Technology Conference (VTC2019-Fall)*, 2019, pp. 1–6.
- [14] Institut für Nachrichtenübertragung, University of Stuttgart. IEEE CTW 2019 Challenge. [Online]. Available: <https://data.ieeeimc.org/Ds1Detail>
- [15] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016, pp. 770–778.
- [16] A. Sobehy, E. Renault, and P. Muhlethaler, "Ndr: Noise and dimensionality reduction of csi for indoor positioning using deep learning," in *2019 IEEE Global Communications Conference (GLOBECOM)*, 2019, pp. 1–6.