

# Supervised Contrastive CSI Representation Learning for Massive MIMO Positioning

Junquan Deng, Wei Shi, Jianzhao Zhang, Xianyu Zhang, and Chuan Zhang

**Abstract**—Similarity metric is crucial for massive MIMO positioning utilizing channel state information (CSI). In this letter, we propose a novel massive MIMO CSI similarity learning method via deep convolutional neural network (DCNN) and contrastive learning. A contrastive loss function is designed considering multiple positive and negative CSI samples drawn from a training dataset. The DCNN encoder is trained using the loss so that positive samples are mapped to points close to the anchor’s encoding, while encodings of negative samples are kept away from the anchor’s in the representation space. Evaluation results of fingerprint-based positioning on a real-world CSI dataset show that the learned similarity metric improves positioning accuracy significantly compared with other known state-of-the-art methods.

**Index Terms**—6G positioning, massive MIMO, channel state information, contrastive learning.

## I. INTRODUCTION

ACQUIRING location information of mobile devices is essential in many smart city and internet-of-things (IoT) applications, including traffic monitoring, asset tracking, autonomous driving, emergency rescue and so on. 5G and beyond wireless networks feature the widely use of massive multiple-input multiple output (MIMO) transmission technique, which not only increases the spectrum efficiencies of wireless transmissions, but also equips the networks with higher sensing and positioning capabilities. Massive MIMO channel state information (CSI) with channel responses on multiple antennas, has been utilized for positioning using geometrical methods [1]–[3], fingerprinting [4]–[11], channel charting (CC) [12]–[17] and direct mapping via deep neural network (DNN) [18]–[20]. Geometrical methods are based on Direction-of-Arrival (DOA) and Time-of-Arrival (TOA) estimates, which require rigorous array calibration and accurate synchronization among distributed network entities. Fingerprinting, CC and direct mapping, alleviate these challenges, by using a labeled dataset and machine learning methods, including k-nearest neighbors (kNN), supervised dimensionality reduction and deep neural network, to predict CSI’s associated locations. CSI similarity metric plays a vital role in fingerprinting and CC-based positioning methods. In fact, the massive MIMO positioning methods in [4]–[12], [14]–[17] build on the assumption that two CSI samples measured at nearby locations should be close to each other with a specific similarity metric.

JD, WS, JZ and XZ are with Sixty-third Research Institute, National University of Defense Technology, Nanjing, China (e-mail: jqdeng@nudt.edu.cn; w.shi@nudt.edu.cn; jianzhao63s@nudt.edu.cn; zhangxy\_sat@126.com). CZ is with LEADS of Southeast University, and Purple Mountain Laboratories, Nanjing, China (e-mail: chzhang@seu.edu.cn). This work is supported by NSFC under grant 61901497 and 62131005, in part by China Postdoctoral Science Foundation (No. 2021MD703980) and Research Project of NUDT under grant ZK 19-09.

A joint angle-delay similarity coefficient was proposed in [7], and correlation matrix distance (CMD) was applied in [8], both using the angle-delay channel power matrix (ADCPM) as CSI feature. In [15], [16], the CSI similarity metric is based on power angular profile (PAP), which is estimated by a multiple signal classification (MUSIC) algorithm. In [12], the Frobenius distance between channel covariance matrices scaled by a vary factor is adopted to measure CSI difference. The above-mentioned CSI similarity metrics assume antenna arrays with perfect linear structures and calibration, making them questionable for practical massive MIMO systems. Other hand-crafted CSI similarity metrics have been designed for general MIMO systems. For example, average Euclidean distance between CSI magnitude values over multiple antennas and sub-carriers is used in [4], [21]. Furthermore, the left singular vector corresponding to the largest singular value of MIMO channel frequency response matrix is selected as the CSI feature in [6], and inner product of such features is used to measure the CSI similarity; Such a similarity metric is closely related to Chordal distance [22] on the Grassmannian manifold.

The CSI estimated at the MIMO receiver from a transmitter position is affected by the surrounding environment and the characteristics of radio frequency (RF) chains, which are hard to be parameterized by a deterministic model. Direct mapping via DNN tries to learn a function  $\bar{\mathbf{p}} = \text{Proj}(\text{CSI})$  to predict the position  $\mathbf{p}$  from CSI, which inherently describes the uncertain environment and RF characteristics. However, to train reliable DNNs, dense CSI samples are needed [18]–[20], and their positioning accuracies are inferior to classical kNN methods as indicated in [21], [23], [24], especially when the training CSI samples are sparse.

In this letter, we propose a new CSI similarity model based on contrastive learning [25], [26], for single-site massive MIMO positioning. The goal is to learn reliable CSI features so that CSI samples collected at neighboring locations are projected by a deep convolutional neural network (DCNN) to neighboring points in a intermediate Euclidean feature space. In the leaning phase, we construct positive and negative CSI samples for each CSI based on their location information in the training dataset. The DCNN is trained using a novel contrastive loss, which considers multiple positive and multiple negative samples. The advantages of our proposed CSI similarity model are: **i)** universally applicable to different wireless systems with different types of CSI data as long as spatial consistency [27] holds for the wireless channels; **ii)** no specific knowledge of the antenna array structure and array calibration are needed; **iii)** adaptive to complex radio environment with non-line-of-sight (NLOS) and multi-path conditions. The code is released at <https://github.com/dengjunquan/SupConCSI>.

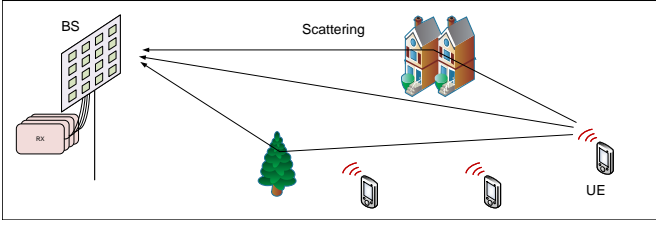


Fig. 1: Single-site massive MIMO positioning scenario.

## II. SYSTEM MODEL AND PROBLEM FORMULATION

We consider a single-site massive MIMO positioning system, where a base station (BS) with  $B$  antennas receives pilot signals from single-antenna user equipments (UEs), as depicted in Fig. 1. A typical MIMO orthogonal frequency division multiplexing (OFDM) air interface is adopted, and the raw CSI from an UE is represented by the estimated spatial-frequency channel response matrix as

$$\mathbf{H} = \mathbf{H}_o + \mathbf{H}_e = [\mathbf{h}_0, \mathbf{h}_1, \dots, \mathbf{h}_{N-1}] + \mathbf{H}_e \in \mathbb{C}^{B \times N}, \quad (1)$$

where  $N$  is the number of OFDM sub-carriers,  $\mathbf{h}_i \in \mathbb{C}^{B \times 1}$  is the channel response vector on the  $i$ th sub-carrier, and  $\mathbf{H}_e$  represents the channel estimation error. The values of  $\mathbf{H}$  are complex number and affected by a range of factors, including powers, DOAs and delays of multi-path components, and antenna array structure, antenna pattern & coupling effects, carrier frequency offset (CFO), symbol timing offset (STO) [28], which are difficult to be parameterized in practical systems. We make no assumption on details of these factors, and assume that the relationship between  $\mathbf{H}$  and the UE position  $\mathbf{p}$  can be viewed as a black box model  $\mathbf{H} = \mathcal{G}(\mathbf{p}, \mathcal{X})$ , with  $\mathcal{X}$  representing the unknown factors.

In a fingerprint-based massive MIMO positioning system, a database consist of CSI samples and their corresponding UE positions is first constructed via a dedicated site survey or crowdsourcing. Denoting the database as  $\{\mathbf{H}_i, \mathbf{p}_i\}_{i=1 \dots I}$ , our goal is to learn a similarity metric  $S(\mathbf{H}_i, \mathbf{H}_j)$  for  $i \neq j$ , such that if  $\mathbf{p}_j$  is a  $k$ -nearest neighbor of  $\mathbf{p}_i$  ( $i = 1 \dots I$ ) in the geographical space,  $\mathbf{H}_j$  should also be a  $k$ -nearest neighbor of  $\mathbf{H}_i$  with this metric. This should hold for a new CSI measurement, i.e., if  $\mathbf{H}$  is measured from an unknown location  $\mathbf{p}$  and  $\mathbf{p}_j$  ( $j = 1 \dots I$ ) is a  $k$ -nearest neighbor of  $\mathbf{p}$ ,  $\mathbf{H}_j$  should also be a  $k$ -nearest neighbor of  $\mathbf{H}$ .

## III. CONTRASTIVE CSI REPRESENTATION LEARNING

Contrastive learning (CL) methods have achieved great successes in computer vision [25], [26], its goal is to learn a representing space where positive samples stay close to the anchor sample, while negative ones are far apart. The mapping function is typically implemented by a neural network. For each anchor, its positives can be generated using data augmentation techniques in an unsupervised setting [25] or leveraging the available labels in a supervised version [26], while its negatives can be drawn randomly from the sample set.

Following a typical CL procedure as in [26], we now detail the contrastive CSI representation learning pipeline for massive MIMO positioning, as shown in Fig. 2. It has three key steps: **i**) construction of positive and negative samples based on the available position information for an anchor CSI; **ii**) map the

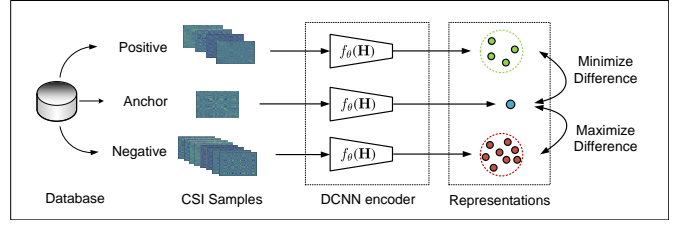


Fig. 2: Contrastive CSI representation learning pipeline.

CSI samples to the feature space via an encoder  $\mathbf{z} = f_\theta(\mathbf{H})$ ; **iii**) use a contrastive loss and an optimizer to update the weights of the encoder.

### A. Construction of Positive and Negative Samples

For a CSI sample  $\mathbf{H}_a$  in a mini-batch  $\{\mathbf{H}_a, \mathbf{p}_a\}_{a \in \mathcal{A}}$  of the training dataset  $\{\mathbf{H}_i, \mathbf{p}_i\}_{i=1 \dots I}$ , its positive samples  $\{\mathbf{H}_p\}_{p \in \mathcal{P}_a}$  are the  $|\mathcal{P}_a|$  nearest samples in the training dataset according to Euclidean distance  $\|\mathbf{p}_i - \mathbf{p}_a\|_2$ . Meanwhile, the negative ones  $\{\mathbf{H}_n\}_{n \in \mathcal{N}_a}$  are randomly selected from the training dataset with  $\|\mathbf{p}_i - \mathbf{p}_a\|_2 > d_{th}$ , with  $d_{th}$  is a predefined distance threshold. In what follows, the mini-batch, positive and negative samples are represented by the index sets  $\mathcal{A}, \mathcal{P}_a, \mathcal{N}_a \subset \{1 \dots I\}$ .

### B. DCNN-Based Encoder Network

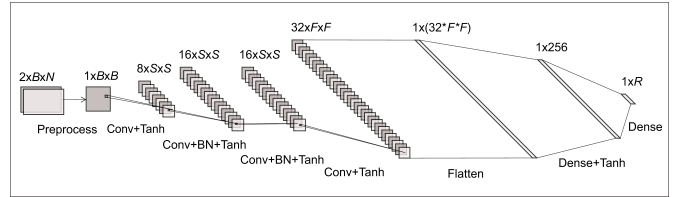


Fig. 3: Proposed DCNN-based CSI encoder  $f_\theta(\mathbf{H})$ .

In order to learn robust intermediate CSI features for positioning purpose, we have designed a DCNN-based encoder network  $f_\theta(\mathbf{H})$  as depicted in Fig. 3, which is consist of a fixed CSI preprocessing layer, a DCNN block with 2D convolution layers and a projection head with two dense layers. The dimension of output feature is  $1 \times R$  and is much smaller than the original CSI size of  $2 \times B \times N$  as the input  $\mathbf{H}$  is a complex matrix with real and imaginary parts. Noticing that autocorrelation over spatial or delay domains are more stable than the raw MIMO-OFDM channel response matrix [11], [15], in the preprocessing layer we first compute the autocorrelation matrix  $\mathbf{C} = \mathbf{H}\mathbf{H}^\dagger$ . As  $\mathbf{C}$  is Hermitian, we then use

$$\mathbf{R} = \Re\{\text{triu}(\mathbf{C})\} + \sqrt{-1} \Im\{\text{tril}(\mathbf{C})\} + \text{diag}(\mathbf{C}) \quad (2)$$

as the input for convolution layers, with  $\text{triu}(\cdot)$ ,  $\text{tril}(\cdot)$ ,  $\text{diag}(\cdot)$  denoting upper triangular, lower triangular and diagonal parts of a matrix, and  $\Re\{\cdot\}$ ,  $\Im\{\cdot\}$  the real and imaginary parts of a matrix.

The DCNN block has four 2D convolution layers, each followed by an activation layer. Different from conventional image input for a CNN, here entities of  $\mathbf{R}$  can be positive or negative, so we choose the **Tanh** function for all activation layers. To avoid the problems of over-fitting and vanishing or exploding gradient, we add two batch normalization (BN)

layers in the middle as shown in Fig. 3. More details of the network parameters are given in Section IV.

### C. Supervised Contrastive Loss

Given the CSI encoder network  $\mathbf{z} = f_\theta(\mathbf{H})$  and the positive and negative samples  $\mathcal{P}_a$  and  $\mathcal{N}_a$  related to an anchor  $\mathbf{H}_a$  in a mini-batch  $\mathcal{A}$  of size  $A$ , we propose a contrastive loss function which is given by

$$\mathcal{L} = \frac{1}{A} \sum_{a \in \mathcal{A}} \left[ - \sum_{p \in \mathcal{P}_a} \log \frac{e^{\mathbf{z}_a \cdot \mathbf{z}_p / \tau}}{\sum_{n \in \mathcal{N}_a} e^{\mathbf{z}_a \cdot \mathbf{z}_n / \tau} + \sum_{r \in \mathcal{P}_a} e^{\mathbf{z}_a \cdot \mathbf{z}_r / \tau}} \right], \quad (3)$$

where  $\mathbf{z}_a$ ,  $\mathbf{z}_p$  and  $\mathbf{z}_n$  are the representations of an anchor, positive and negative CSI samples, and the dot symbol denotes inner product. The fraction before taking logarithm represents the probability that  $\mathbf{z}_a$  selects  $\mathbf{z}_p$  over all positives and negatives as one of its neighbours in feature space. In this regard, minimizing  $\mathcal{L}$  would make  $\mathbf{z}_a$  and  $\mathbf{z}_p$  close. The temperature  $\tau$  is used for tuning how concentrated the features are in the feature space. With a low temperature,  $\mathcal{L}$  is dominated by the small distances and widely separated features contribute less.

The proposed loss is inspired by both the InfoNCE loss [29] and the soft nearest neighbor (SNN) loss [30], commonly used in self-supervised representation learning. The differences are that: 1) InfoNCE is suitable for instance discrimination tasks and involves only a single positive sample; 2) In SNN loss, the summation in the numerator over positives is located inside the logarithm operation; Moreover, we use inner product  $\mathbf{z}_i \cdot \mathbf{z}_j$  instead of  $-\|\mathbf{z}_i - \mathbf{z}_j\|_2^2$  in SNN [30] to measure the similarity of the embedded features, which is more stable and efficient for gradient-descent computation.

To give more details of  $\mathcal{L}$ , we derive its gradient with respect to  $\mathbf{z}_a$ . We focus on the loss  $\mathcal{L}_a$  related to  $\mathbf{z}_a$ , which is expressed as inside the brackets in (3). The gradient is given by

$$\frac{\partial \mathcal{L}_a}{\partial \mathbf{z}_a} = \frac{|\mathcal{P}_a|}{\tau} \left[ \sum_{p \in \mathcal{P}_a} - \left[ \frac{1}{|\mathcal{P}_a|} - x_{ap} \right] \mathbf{z}_p + \sum_{n \in \mathcal{N}_a} x_{an} \mathbf{z}_n \right], \quad (4)$$

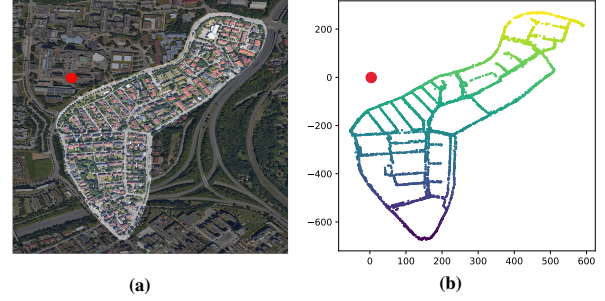
where  $x_{ap} = e^{\mathbf{z}_a \cdot \mathbf{z}_p} / \sum_{j \in \mathcal{N}_a \cup \mathcal{P}_a} e^{\mathbf{z}_a \cdot \mathbf{z}_j / \tau}$  and the same form for  $x_{an}$ . In the initialization phase, the encoder  $f_\theta(\mathbf{H})$  would generate random features, then  $x_{ap}$  and  $x_{an}$  would have values around  $1/(|\mathcal{P}_a| + |\mathcal{N}_a|)$ . When there are enough negatives,  $x_{ap}$  is smaller than  $1/|\mathcal{P}_a|$  for most of the positives during training, so  $\mathbf{z}_a$  will be dragged towards the mean positive representation vector, while be pulled away from the negatives via gradient-descent. With a well-trained encoder  $f_\theta(\mathbf{H})$ , the supervised contrastive learning-based CSI similarity metric (**SupCon**) is

$$S(\mathbf{H}_i, \mathbf{H}_j) \stackrel{\text{def}}{=} \|f_\theta(\mathbf{H}_i) - f_\theta(\mathbf{H}_j)\|_2^{-1} = \|\mathbf{z}_i - \mathbf{z}_j\|_2^{-1}, \quad (5)$$

and will be used for downstream kNN-based positioning tasks.

## IV. EXPERIMENTS

We perform single-site positioning task on a real outdoor massive MIMO dataset provided by IEEE Communications Theory Workshop data competition [31]. The CSI measurements were taken in a residential area of about 0.3 km<sup>2</sup>, as depicted in Fig. 4 a). The CSI data was collected by a channel sounder at a carrier frequency of 1.27 GHz, the measurement details



**Fig. 4:** a) Map of the measurement campaign in a residential area, where BS is marked with a red dot; b) Locations associated with the CSI measurements, which are coloured w.r.t. their Y coordinates [m].

**TABLE I:** Settings for the DCNN Encoder and Its Training.

Conv. layer 1	kernel size = 5, stride = 2, padding = 1
Conv. layer 2	kernel size = 3, stride = 1, padding = 1
Conv. layer 3	kernel size = 3, stride = 1, padding = 1
Conv. layer 4	kernel size = 3, stride = 1, padding = 0
Feature space	$R = 32$
Optimizer	Adam
Learning rate	0.0005 (divided by 2 every 5 epochs)
Weight decay	0.0001
No. of epochs	20
Batch size	$A = 32$
Distance th.	$d_{th} = 25$ m
Temperature	$\tau = 1.5$

are given in [28]. The data are made up of channel frequency responses between a moving single-antenna transmitter and a fixed receiver which has a  $8 \times 8$  patch antenna array. The uplink OFDM channel has a bandwidth of 20 MHz with 1024 sub-carriers, of which 100 are used for guard bands, and  $N = 924$  are effective. There are 8 antennas which were malfunctioning, so the effective antenna number is  $B = 56$ . Notice that the position of the effective antennas in the array is not provided.

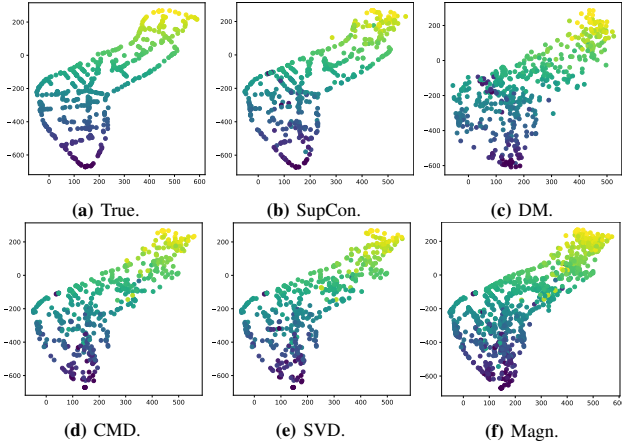
The data contains a total of 4979 labeled CSI measurements with locations given by a GPS device, as shown in Fig. 4 b). To train the encoder  $f_\theta(\mathbf{H})$  and evaluate its positioning capability, we randomly split the data into a training and a test set. The ratio  $\rho$  of test samples is set to be 10% (with 498 samples) or 20% (with 996 samples), similar to the settings in [21] and [6].

More details of the DCNN-based CSI encoder depicted in Fig. 3 are given in Table I. With those settings, the output size of the convolution part is  $32 \times 25 \times 25 = 20000$ , which is then reduced by the projection head to a dimension of  $R = 32$ . We implement the encoder and train it using the loss  $\mathcal{L}$  via Pytorch on a desktop with a GeForce GTX 1660 Ti GPU. The default numbers of positives and negatives we used are  $|\mathcal{P}_a| = 16$  and  $|\mathcal{N}_a| = 64$ . More details for training are also given in Table I. It takes less than 10 minutes to learn a successful CSI encoder for the downstream kNN-based positioning task. Finally, We feed all CSI samples in  $\{\mathbf{H}_i, \mathbf{p}_i\}_{i=1 \dots I}$  to the trained encoder, and obtain a CSI fingerprint database  $\{\mathbf{z}_i, \mathbf{p}_i\}_{i=1 \dots I}$ .

### A. KNN-Based Positioning with the Learned CSI Similarity

For a test CSI  $\mathbf{H}$  measured from an unknown location  $\mathbf{p}$ , we first compute its representation  $\mathbf{z} = f_\theta(\mathbf{H})$ , and find its k-nearest neighbors (denoted by  $\mathcal{K} \subset \{1 \dots I\}$ ) in fingerprint database  $\{\mathbf{z}_i, \mathbf{p}_i\}_{i=1 \dots I}$ , then the predicted position is given by

$$\hat{\mathbf{p}} = \frac{\sum_{i \in \mathcal{K}} S(\mathbf{H}_i, \mathbf{H}) \mathbf{p}_i}{\sum_{i \in \mathcal{K}} S(\mathbf{H}_i, \mathbf{H})} = \frac{\sum_{i \in \mathcal{K}} \|\mathbf{z}_i - \mathbf{z}\|_2^{-1} \mathbf{p}_i}{\sum_{i \in \mathcal{K}} \|\mathbf{z}_i - \mathbf{z}\|_2^{-1}}, \quad (6)$$



**Fig. 5:** a) True positions of 498 test samples in the 2D XY-coordinates [m]; b) Predicted positions via the supervised contrastive learning (SupCon) based similarity metric and kNN with  $k = 4$  neighbors; c) Direct Mapping via a neural network in Fig. 3 with output size  $R$  equals 2; d) kNN with correlation matrix distance (CMD); e) kNN with SVD-based similarity; f) kNN with channel magnitude-based similarity.

and the positioning error is evaluated as  $\|\mathbf{p} - \bar{\mathbf{p}}\|_2$ . Using 498 CSI measurements for test (with  $I = 4481$  for training), their predicted locations are plotted in Fig. 5 b). As compared with the true positions shown in Fig. 5 a), we achieve a mean positioning error of 35.2m, which is better than the lowest errors reported in [6], [24] and [21] as 42m, 40m and 37m respectively. In [6], its best result was achieved by converting the Nadaraya-Watson estimator to a three-layer neural network and using a SVD-based similarity metric. In [24], its best performance was given by kNN with a sophisticated similarity metric based on signal subspace, whose computation complexity is rather high. While in [21], the 37m mean error is achieved via ensemble learning with multi-layer perceptron neural networks (MLP NN), and hand-crafted CSI features extracted from channel magnitude information via polynomial regression and Fourier fitting. Compared with these methods, our proposed positioning pipeline not only has better positioning performance, but also is easier to be implemented.

### B. Comparison with other CSI Similarity Metrics

To gain more insights about the learned similarity metric, in our experiments, we compare it with three other CSI similarity metrics, i.e. CMD, SVD-based, and channel magnitude-based, by replacing  $S(\mathbf{H}_i, \mathbf{H})$  in (6). The CMD similarity [8] is

$$S_{\text{CMD}}(\mathbf{H}_i, \mathbf{H}) = \frac{\text{Tr}\{(\mathbf{H}_i \mathbf{H}_i^\dagger)(\mathbf{H} \mathbf{H}^\dagger)\}}{\|(\mathbf{H}_i \mathbf{H}_i^\dagger)\|_{\text{F}} \|(\mathbf{H} \mathbf{H}^\dagger)\|_{\text{F}}}, \quad (7)$$

with  $\text{Tr}\{\cdot\}$  denotes matrix trace and  $\|\cdot\|_{\text{F}}$  the Frobenius norm. The SVD-based similarity [6] is defined as

$$S_{\text{SVD}}(\mathbf{H}_i, \mathbf{H}) = |\mathbf{v}_i^\dagger \mathbf{v}|, \quad (8)$$

where  $\mathbf{v}$  is the left singular vector w.r.t. the largest singular value of  $\mathbf{H}$ . The magnitude-based similarity is defined as

$$S_{\text{Magn}}(\mathbf{H}_i, \mathbf{H}) = \left[ \sum_b \sum_n (|\mathbf{H}_i|_{b,n} - |\mathbf{H}|_{b,n})^2 \right]^{-1}, \quad (9)$$

where  $|\mathbf{H}|_{b,n}$  is the  $(b, n)$ th element of the channel magnitude matrix  $|\mathbf{H}|$  of  $\mathbf{H}$ .

**TABLE II:** Mean Positioning Errors [m].

	DM	SupCon	CMD	SVD	Magn.
$\rho = 10\%$	54.3	<b>35.2</b>	48.3	50.8	56.6
$\rho = 20\%$	55.5	<b>36.4</b>	48.7	51.4	57.3

For the 498 test samples, their positioning results are shown as in Fig. 5 d), e), f). Compared with our CSI similarity, CMD has a mean positioning error of 48.3m, SVD-based similarity has one of 50.8m, while it is 56.6m for channel magnitude-based, as summarized in Table II. Here, we have set  $k = 4$  neighbors for kNN as it produced the best results for all similarity metrics in our experiments. The cumulative distribution functions (CDFs) of their positioning errors, together with that for our method are given in Fig. 6. The error distribution has a long tail because the CSI samples from the northeast and southernmost regions exhibit very low SNRs. With our method, 66% of the test samples have a positioning error smaller than 25m, while other methods have less than 47% of the test samples achieve this goal.

### C. Comparison with Direct Mapping

The object of direct mapping (DM) [18]–[20] is to train a CSI encoder such that CSI is mapped to a point in the 2D (or 3D) geographic space directly. For DM, we use the same network architecture as in Fig. 3, with an output dimension of  $R = 2$ , while other network parameters are the same as for contrastive CSI feature learning. To train the DM encoder  $\mathbf{z} = f_{\theta, \text{DM}}(\mathbf{H})$ , a loss function  $\mathcal{L}_{\text{DM}} = \frac{1}{A} \sum_{a \in \mathcal{A}} \|\mathbf{z}_a - \mathbf{p}_a\|_2$  is used. Different from Table I, we use an initial learning rate of 0.01 and train the DM encoder for 100 epochs. For a new CSI  $\mathbf{H}$ , the predicted position is then given by  $\bar{\mathbf{p}} = f_{\theta, \text{DM}}(\mathbf{H})$ .

The positioning results using DM is shown in Fig. 5 c). It exhibits a mean error of 54.3m, which is close to kNN with CMD, SVD-based and magnitude-based similarities, but much worse than our contrastive learning-based approach. It indicates that the proposed DCNN-based CSI encoder network do have the capability to learn an inverse mapping from high-dimensional CSI space to the geographic space for the black box model  $\mathbf{H} = \mathcal{G}(\mathbf{p}, \mathcal{X})$ . However, due to the sparsity of the training data and over fitting problem, the learned DM encoder has high generalization errors on the test set. Instead, the proposed contrastive learning method tries to learn an intermediate CSI feature space where CSI's representations exhibit similar neighbour relations as in the geographic space, supervised by the relative distance information of positives and negatives. As shown in Fig. 7, with only one positive and one negative, a mean positioning error of 45.1m can be achieved and the performance improves as both the number of positives  $|\mathcal{P}_a|$  and number of negatives  $|\mathcal{N}_a|$  increase.

## V. CONCLUSIONS

We have proposed a novel contrastive CSI representation learning method for massive MIMO positioning. A contrastive loss function has been designed, involving multiple positive and negative samples. A versatile DCNN-based CSI encoder we designed has been trained using this loss to learn robust CSI representations for fingerprinting-based positioning task. We have conducted extensive experiments on a real-world massive



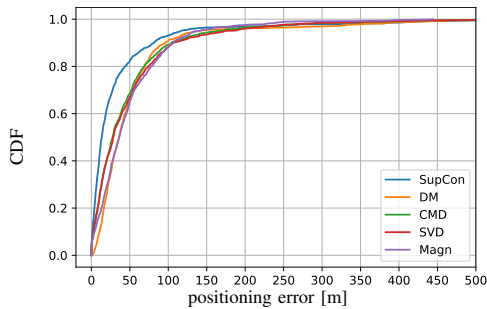


Fig. 6: CDF of positioning errors for kNN with different CSI similarity metrics (i.e. SupCon, CMD, SVD-based and magnitude-based), and direct mapping (DM). Number of test samples is 498.

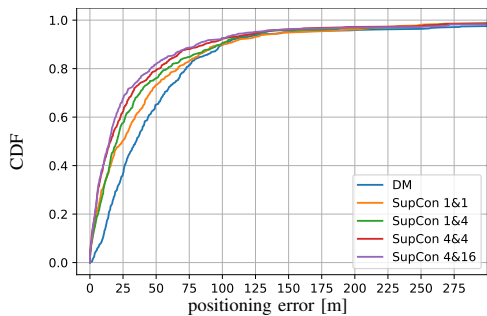


Fig. 7: CDF of positioning errors on the test set for direct mapping (DM), and kNN based on SupCon CSI similarities with different numbers ( $|\mathcal{P}_a|$  &  $|\mathcal{N}_a|$ ) of positives and negatives. Number of test samples is 498.

MIMO dataset measured in a complex outdoor environment. The results show that the learned CSI similarity metric improves the positioning accuracy significantly compared with other known methods. No specific knowledge of the array structure and array calibration are needed for the proposed positioning pipeline. Instead of using a static kNN estimator in the final step, one can also consider using a trainable regressor as in [6] to further improve the positioning accuracy by end-to-end fine-tuning. Another promising research direction is to combine the learned similarity metric with the semi-supervised channel charting methods [15], [16], considering neighborhood relationships in the feature space is generally non-linear.

## REFERENCES

- [1] J. Xiong and K. Jamieson, "ArrayTrack: A Fine-Grained indoor location system," in *10th USENIX Symposium on Networked Systems Design and Implementation (NSDI 13)*, Apr. 2013, pp. 71–84.
- [2] N. Garcia, H. Wymeersch, E. G. Larsson, A. M. Haimovich, and M. Coulon, "Direct localization for massive MIMO," *IEEE Transactions on Signal Processing*, vol. 65, no. 10, pp. 2475–2487, 2017.
- [3] R. Mendrzik, H. Wymeersch, G. Bauch, and Z. Abu-Shaban, "Harnessing NLOS components for position and orientation estimation in 5G millimeter wave MIMO," *IEEE Transactions on Wireless Communications*, vol. 18, no. 1, pp. 93–107, 2019.
- [4] A. Sobehy, E. Renault, and P. Muhlethaler, "CSI-MIMO: K-nearest neighbor applied to indoor localization," in *IEEE International Conference on Communications (ICC)*, 2020, pp. 1–6.
- [5] J. Fan, J. Zhang, and X. Dou, "Single-site indoor fingerprint localization based on MIMO-CSI," *China Communications*, vol. 18, no. 8, pp. 199–208, 2021.
- [6] L. Le Magoarou, "Similarity-based prediction for channel mapping and user positioning," *IEEE Communications Letters*, vol. 25, no. 5, pp. 1578–1582, 2021.
- [7] X. Sun, C. Wu, X. Gao, and G. Y. Li, "Fingerprint-based localization for massive MIMO-OFDM system with deep convolutional neural networks," *IEEE Transactions on Vehicular Technology*, vol. 68, no. 11, pp. 10 846–10 857, 2019.

- [8] C. Wu, X. Yi, W. Wang, L. You, Q. Huang, X. Gao, and Q. Liu, "Learning to localize: A 3D CNN approach to user positioning in massive MIMO-OFDM systems," *IEEE Transactions on Wireless Communications*, vol. 20, no. 7, pp. 4556–4570, 2021.
- [9] J. Vieira, E. Leitingner, M. Sarajlic, X. Li, and F. Tufvesson, "Deep convolutional neural networks for massive MIMO fingerprint-based positioning," in *IEEE 28th Annual International Symposium on Personal, Indoor, and Mobile Radio Communications (PIMRC)*, 2017, pp. 1–6.
- [10] Q. Li, X. Liao, M. Liu, and S. Valaee, "Indoor localization based on CSI fingerprint by siamese convolution neural network," *IEEE Transactions on Vehicular Technology*, vol. 70, no. 11, pp. 12 168–12 173, 2021.
- [11] E. Gönültaş, E. Lei, J. Langerman, H. Huang, and C. Studer, "CSI-based multi-antenna and multi-point indoor positioning using probability fusion," *IEEE Transactions on Wireless Communications*, pp. 1–1, 2021.
- [12] C. Studer, S. Medjkouh, E. Gonultaş, T. Goldstein, and O. Tirkkonen, "Channel charting: Locating users within the radio environment using channel state information," *IEEE Access*, vol. 6, pp. 47 682–47 698, 2018.
- [13] J. Deng, S. Medjkouh, N. Malm, O. Tirkkonen, and C. Studer, "Multipoint channel charting for wireless networks," in *52nd Asilomar Conference on Signals, Systems, and Computers*, 2018, pp. 286–290.
- [14] E. Lei, O. Castañeda, O. Tirkkonen, T. Goldstein, and C. Studer, "Siamese neural networks for wireless positioning and channel charting," in *57th Annual Allerton Conference on Communication, Control, and Computing (Allerton)*, 2019, pp. 200–207.
- [15] J. Deng, O. Tirkkonen, J. Zhang, X. Jiao, and C. Studer, "Network-side localization via semi-supervised multi-point channel charting," in *International Wireless Communications and Mobile Computing (IWCMC)*, 2021, pp. 1654–1660.
- [16] J. Deng, W. Shi, J. Hu, and X. Jiao, "Semi-supervised t-SNE for millimeter-wave wireless localization," in *7th International Conference on Computer and Communications (ICCC)*, 2021, pp. 1015–1019.
- [17] P. Ferrand, A. Decurninge, L. G. Ordoñez, and M. Guillaud, "Triplet-based wireless channel charting: Architecture and experiments," *IEEE Journal on Selected Areas in Communications*, vol. 39, no. 8, pp. 2361–2373, 2021.
- [18] P. Ferrand, A. Decurninge, and M. Guillaud, "DNN-based localization from channel estimates: Feature design and experimental results," in *IEEE Global Communications Conference*, 2020, pp. 1–6.
- [19] S. D. Bast, A. P. Guevara, and S. Pollin, "CSI-based positioning in massive MIMO systems using convolutional neural networks," in *IEEE 91st Vehicular Technology Conference (VTC2020-Spring)*, 2020, pp. 1–5.
- [20] G. Cerar, A. Švigelj, M. Mohorčič, C. Fortuna, and T. Javornik, "Improving CSI-based massive MIMO indoor positioning using convolutional neural network," in *Joint European Conference on Networks and Communications 6G Summit (EuCNC/6G Summit)*, 2021, pp. 276–281.
- [21] A. Sobehy, É. Renault, and P. Muhlethaler, "Generalization aspect of accurate machine learning models for CSI-based localization," *Annals of Telecommunications*, Jun 2021.
- [22] K. Ko and J. Lee, "Multiuser MIMO user selection based on chordal distance," *IEEE Transactions on Communications*, vol. 60, no. 3, pp. 649–654, 2012.
- [23] W. Y. Al-Rashdan and A. Tahat, "A comparative performance evaluation of machine learning algorithms for fingerprinting based localization in DM-MIMO wireless systems relying on big data techniques," *IEEE Access*, vol. 8, pp. 109 522–109 534, 2020.
- [24] P. Garau Burguera, "Logical radio maps for user localization in a real outdoor radio environment," Master's thesis, Aalto University, 2020.
- [25] T. Chen, S. Kornblith, M. Norouzi, and G. Hinton, "A simple framework for contrastive learning of visual representations," in *International conference on machine learning*, 2020, pp. 1597–1607.
- [26] P. Khosla, P. Teterwak, C. Wang, A. Sarna, Y. Tian, P. Isola, A. Maschinot, C. Liu, and D. Krishnan, "Supervised contrastive learning," in *Advances in Neural Information Processing Systems*, vol. 33, 2020, pp. 18 661–18 673.
- [27] M. Kurras, S. Dai, S. Jaeckel, and L. Thiele, "Evaluation of the spatial consistency feature in the 3GPP geometry-based stochastic channel model," in *IEEE Wireless Communications and Networking Conference (WCNC)*, 2019, pp. 1–6.
- [28] M. Gauger, M. Arnold, and S. ten Brink, "Massive MIMO channel measurements and achievable rates in a residential area," in *24th International ITG Workshop on Smart Antennas*, 2020, pp. 1–6.
- [29] A. Van den Oord, Y. Li, O. Vinyals et al., "Representation learning with contrastive predictive coding," *arXiv preprint arXiv:1807.03748*, vol. 2, no. 3, p. 4, 2018.
- [30] N. Frosst, N. Papernot, and G. Hinton, "Analyzing and improving representations with the soft nearest neighbor loss," in *36th International Conference on Machine Learning*, vol. 97, Jun 2019, pp. 2012–2020.
- [31] "IEEE CTW 2020 Challenge," <https://data.ieeeemc.org/Ds4Detail>, 2020.