
Exploiting Temporal Priors for Efficient Real-time Compression and Feedback of Wireless Channels

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Machine learning based compression frameworks are rapidly gaining popularity
2 as the demand for efficient storage, processing, and transmission of large-scale
3 data continues to grow across diverse applications such as video streaming and IoT.
4 Recently, such frameworks have also sparked significant interest in wireless com-
5 munications and the task of ML based wireless channel compression is currently
6 one of the use cases being explored by the international wireless standards body,
7 3GPP, for standardization. In wireless communication systems, each user device or
8 user equipment (UE) typically estimates the wireless channel between the transmit-
9 ting base station (BS) and itself and feeds back information related to the estimated
10 channel state information (CSI) to the serving BS, which may then be utilized for
11 downstream processing. While the current 5G communication stack employs a
12 combination of matrix factorization and quantization approaches to compress the
13 CSI, autoencoders (AE) have emerged as a viable option to compress the estimated
14 spatial-frequency (SF) channel sample and send it back to the base station for
15 reconstruction. Although the AE-based approaches have shown acceptable CSI re-
16 construction performance, there is still a large room for further improvement, both
17 from an overhead reduction as well as reconstruction performance perspectives.
18 This paper proposes a new AE framework that leverages the temporal correlation
19 properties of the channel to enhance the compression process. In particular, we
20 propose an AE framework that performs temporal-spatial-frequency (TSF) com-
21 pression by utilizing priors based on historical CSI samples to efficiently compress
22 the current estimated CSI sample. End-to-end simulation results on a realistic
23 test bench demonstrate the superiority of the proposed TSF compression approach
24 relative to the state-of-the-art methods.

25 1 Introduction

26 Compression tasks are critical across numerous domains, including but not limited to, telecommuni-
27 cations, healthcare, video streaming and IoT, where efficiently storing and transmitting large amounts
28 of data is essential for maintaining performance and reducing costs. Compression techniques like
29 compressive sensing, matrix and tensor decompositions have long been used to reduce data dimen-
30 sionality and extract meaningful features; however, autoencoders (AE) offer a distinct advantage
31 by learning non-linear representations directly from raw data. Autoencoders can preserve essential
32 features while reducing redundancy, ultimately leading to more efficient, scalable and cost-effective
33 solutions in domains such as video processing, wireless communications, e.g., sensor data analytics
34 and wireless channels.

35 For time-series data compression, AEs can potentially be more efficient by leveraging Recurrent
36 Neural Networks and Transformers to exploit the inherent temporal correlations within the data to

37 achieve superior compression. Depending on the application, the availability of data and the set of
38 underlying constraints may differ. While some applications may have the flexibility to utilize multiple
39 or all temporal data samples before compressing and transmitting the data, others may require
40 real-time compression and transmission as soon as a new data sample arrives. The latter is indeed
41 more challenging from a compression perspective. For example, video data inherently benefits from
42 temporal correlations between consecutive samples (GOPs or group of pictures), where the content
43 remains similar over short periods, allowing compression techniques to exploit these redundancies
44 for efficient compression. In contrast, wireless data, e.g., sensor data or cellular channels, presents a
45 different challenge, as it requires real-time processing and transmission of data arriving sequentially
46 over time. In other words, while video has all or at least a few temporally-correlated samples available
47 a priori for analysis and compression, wireless systems must compress and transmit any observed or
48 estimated data without delay, making real-time compression more complex and time-sensitive.

49 In this work, we introduce a novel autoencoder framework that can efficiently handle real-time
50 compression of time-series data. While the proposed framework can be widely applicable to several
51 time-series data compression use-cases, e.g., wireless sensor networks, IoT, etc., we here adopt the
52 wireless channel compression as an example use-case for real-time time-series data compression. The
53 wireless channel data compression has recently gained popularity in the machine learning domain and
54 it is one of the few ML use-cases that is currently considered as a study item in the 3GPP wireless
55 standards. In the context of wireless channel compression, the processing and transmission of each
56 sample has to be completed within 1-2 milliseconds (ms) while the time difference between two
57 consecutive samples can range from 5 ms to 20 ms. Thus, each sample has to be compressed and
58 transmitted before the next sample arrives. The wireless channel compression is a crucial task in
59 massive multiple-input-multiple-output (MIMO) systems.

60 Massive MIMO is a leading technology that has the potential to meet the data rate requirements in
61 the next-generation wireless communication systems [4]. The key advantage of employing large
62 antenna arrays in massive MIMO systems is the capability of achieving striking performance gains in
63 multiuser MIMO systems. However, achieving such gains is contingent on the availability of accurate
64 channel state information (CSI) at the transmitter. This requires the receiver, e.g., user equipment, to
65 send back the estimated CSI to the transmitter, e.g., base station. The CSI feedback process incurs
66 additional system overhead that scales up with the number of transmit antennas, the number of receive
67 antennas, and the number of allocated frequency resources, thereby resulting in a considerable uplink
68 overhead that can impact the system performance.

69 Autoencoder (AE)-based CSI compression and feedback has gained significant popularity in the
70 wireless domain [2, 3, 6]. The current 5G compliant communication systems under deployment across
71 the world employ a combination of matrix factorization and quantization approaches to compress and
72 feedback the CSI [1]. AE-based solutions have the potential to offer a favorable trade-off between
73 the CSI feedback overhead and reconstruction performance. Existing AE-based approaches operate
74 on solely on the spatial-frequency (SF) channel samples, where the spatial components represent the
75 number of transmit side and receive side antennas while the frequency components refer to the number
76 of orthogonal frequencies over which the data transmission happens. In such setups, the estimated
77 CSI sample at each user is compressed using an encoder network, and then the compressed version is
78 sent to the base station for reconstruction using a decoder network. This process is repeated for every
79 CSI reporting instance, where each collected sample over time is compressed independently. While
80 SF compression was shown to provide acceptable reconstruction quality, there is still large room for
81 further improvement in the reconstruction performance or reduction in the signaling overhead.

82 One way to further improve the performance of the SF compression approach is to leverage the
83 temporal correlation properties of the channel in the compression process. Exploiting the CSI
84 temporal correlation on the top of SF compression has the potential to provide further i) improvement
85 in the reconstruction performance for a given overhead, ii) reduction in the overhead for a given
86 performance and/or ii) improvement in both performance and overhead relative to the SF compression.

87 To address the CSI feedback overhead reduction problem, we propose a new AE framework that
88 seeks to efficiently compress the current spatial-frequency CSI sample by utilizing the temporal
89 correlation of the channel in the compression process – this technique is referred to as Temporal-
90 Spatial-Frequency (TSF) compression. Preliminary simulation results on a realistic 5G compliant
91 test bench show that exploiting the past collected CSI samples in the compression task can result in
92 considerable throughput gains relative to SF compression and state-of-the-art methods.

93 **2 Problem Statement**

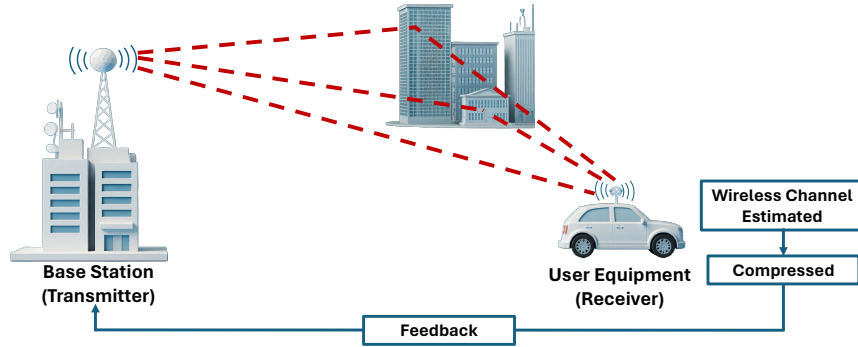


Figure 1: A representation of the downlink communication with base station (BS) as the transmitter and a self driving car as the receiver (user equipment(UE)). The transmitted signal from the BS, travels through multiple paths (red dotted lines) before being received at the UE. The wireless channel \mathbf{H}_n represents the overall impulse response associated with the signal propagation and is estimated at the UE. The channel is then compressed and sent back to the BS for further downstream processing.

94 Consider a downlink data transmission setup where a single BS equipped with N_t transmit antennas
 95 is serving (or transmitting) to a single UE with N_r receive antennas. The UE receives its data over
 96 multiple frequency components, e.g., sub-carriers (orthogonal frequencies). The UE needs to estimate
 97 and transmit the channel tensor, $\mathbf{H}_n \in \mathbb{C}^{N_r \times N_t \times N_c}$, where N_c denotes the number of frequency
 98 components. The goal is then to compress the estimated channel \mathbf{H}_n assuming that the user and base
 99 station may have access to upto N historical samples, i.e., $\mathbf{H}_{n-1}, \dots, \mathbf{H}_{n-N}$.

100 **3 TSF Framework and Model Architecture**

101 In this section, we introduce the framework designed to enable real-time compression and feedback of
 102 the current channel instance between a transmission and reception node, whilst leveraging historical
 103 channel information to improve reconstruction performance.

104 **3.1 System Framework**

105 Given a maximum look-back size of N historical samples, we consider a set of $N + 1$ encoder-decoder
 106 pairs, one associated with each possible value of the available past samples $\{0, \dots, N\}$. Thus for the
 107 channel sample at time n , i.e. \mathbf{H}_n , the k -th model is utilized wherein, $k = (n \bmod N + 1)$. The k -th
 108 model utilizes k past channel samples for both, encoding and decoding. Thus, for a value of $k = 0$,
 109 no past information is utilized for the compression, and the channel \mathbf{H} is compressed standalone.

110 A typical compression pipeline consisting of an encoder and a quantizer is used by the UE to obtain
 111 a compressed representation of the channel, \mathbf{z}_n with dimensionality D_k . The UE feeds back the
 112 compressed representation, \mathbf{z}_n , to the base station which decompresses it using its decoder. The
 113 UE is further equipped with the same decoder model being utilized by the BS, this allows both UE
 114 and BS to have the same reference for the past samples. Post decoding, both BS and UE store the
 115 reconstructed channel in a buffer. The reconstructed channels stored in the UE-side buffer are utilized
 116 as priors for compressing the channel samples at the next time instance. Since the UE and BS utilize
 117 the same priors or past samples, the encoder and decoder remain synchronized (in the absence of
 118 packet loss and noise) allowing for better compression and reconstruction of the channel data. A
 119 diagrammatic representation of the compression framework can be seen in Fig. 2.

120 As mentioned earlier, the selection of the encoder-decoder pair for a specific n is governed by
 121 $k = (n \bmod N + 1)$. This setup ensures any noise, error, or packet loss that may have been
 122 introduced during transmission (or feedback) is not accumulated for more than N samples. Further,
 123 having $N + 1$ models, each dedicated to a specific look-back period, enables us to identify and
 124 analyze the maximum performance or improvement such a setup could achieve. However, this is

125 achieved at the expense of having multiple models with potential redundancies across their learned
 126 layers/weights.

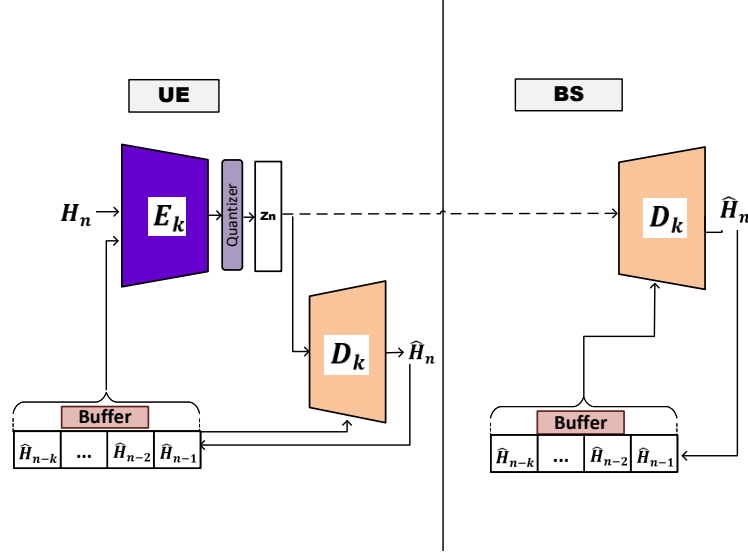


Figure 2: Proposed TSF Autoencoder framework. The UE utilizes an encoder model along with past reconstructed channel samples to compress the current channel sample. The NW and UE both utilize the same decoder network to reconstruct the channel sample with the synchronized information about the past channel data.

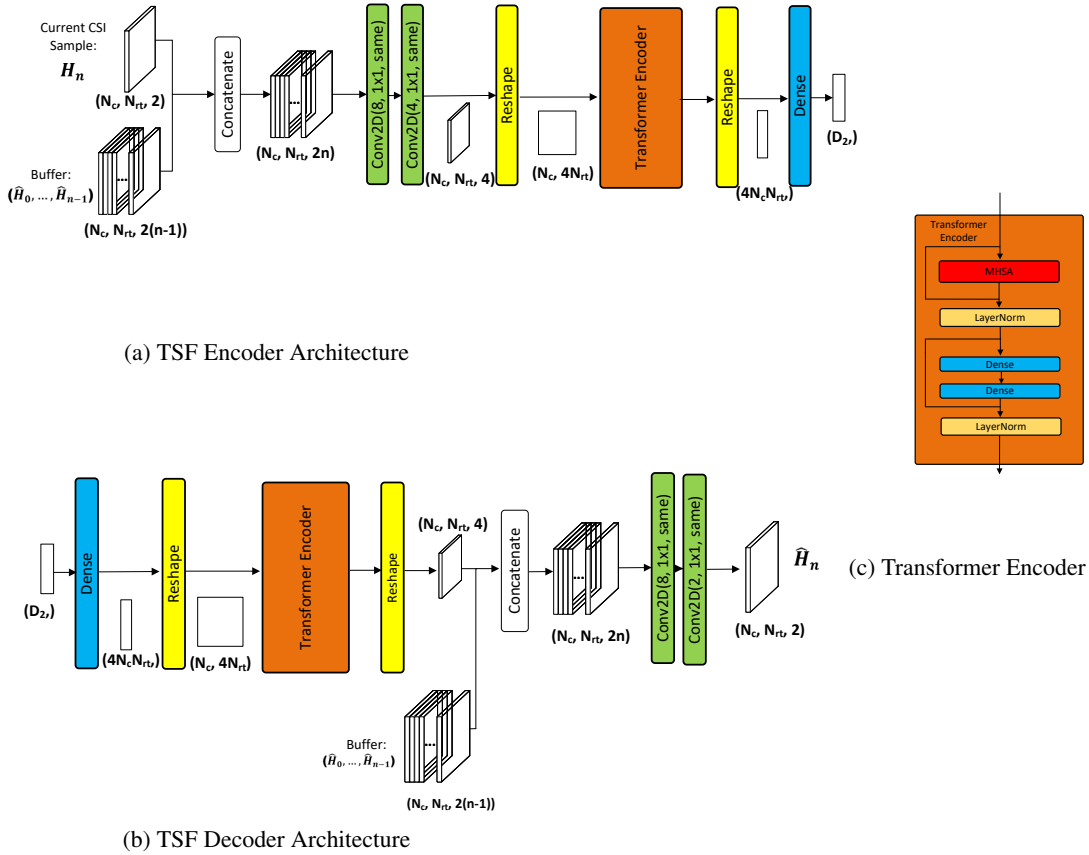
127 3.2 TSF Model Architecture

128 The full architecture of the encoder and decoder blocks are detailed in Fig. 3. The previously decoded
 129 samples $(\hat{\mathbf{H}}_{n-1}, \hat{\mathbf{H}}_{n-2}, \dots, \hat{\mathbf{H}}_{n-k})$ are combined with the current channel sample \mathbf{H}_n at the input
 130 of the encoder. The input tensor has shape $\mathbf{H}_{in} \in \mathbb{R}^{N_c \times N_{rt} \times 2(k+1)}$, where $(k+1)$ represents the
 131 current and the past k channel samples, 2 represents the real and imaginary parts of the complex
 132 channels and $N_{rt} = N_r * N_t$. We compute 2D convolutions with a 1×1 kernel, so that each element
 133 in the output is derived from a combination of elements at the same position in the input, across
 134 all the current and past samples. This representation is reshaped appropriately for input to the
 135 transformer block shown in Fig. 3c. The multi-head self-attention block extracts pairwise similarities
 136 between frequency sub-bands, N_c resulting in an $(N_c \times N_c)$ attention matrix that is used to weight
 137 the full input tensor. Channel samples typically exhibit high correlation across sub-bands therefore
 138 the attention mechanism can be viewed as removing redundancy by focusing on the most relevant
 139 sub-bands. This representation is passed to a position-wise feed-forward layer that transforms the
 140 features of each sub-band independently. The output of the transformer block is reshaped and passed
 141 through a final dense layer and discretized using 2-bit scalar quantization, giving z_n .

142 At the decoder (Fig. 3b) the prior samples are introduced into the model post the transformer. The
 143 objective being, that the earlier layers in the encoder learn to filter out the redundant or correlated
 144 information across samples and the transformer part of the encoder and decoder models only focus
 145 on compressing and reconstructing the non-redundant information. The redundant information from
 146 past samples can then be reintroduced into the data at the final stages of the decoding utilizing the
 147 past samples and 2D convolution layers.

148 3.3 Training

149 As mentioned earlier, we train a total of $N+1$ encoder-decoder pairs $\{(E_0, D_0), \dots, (E_N, D_N)\}$
 150 models. These models are trained serially. Given a set of sequential channel samples, $\mathbf{S} =$
 151 $\{\mathbf{H}_1, \mathbf{H}_2, \dots, \mathbf{H}_i, \dots, \mathbf{H}_{|\mathbf{S}|}\}$, where $|\mathbf{S}|$ represents the cardinality of the set \mathbf{S} . The first model is trained
 152 such that E_1 and D_1 minimize the reconstruction loss $\|\mathbf{H}_i - D_1(E_1(\mathbf{H}_i))\|_F^2 \forall i \in \mathbf{S}$. For the second
 153 encoder-decoder pair, we seek to minimise the reconstruction loss: $\|\mathbf{H}_i - D_2(E_2(\mathbf{H}_i, \hat{\mathbf{H}}_{i-1}))\|_F^2 \forall i \in$



(a) TSF Encoder Architecture

(b) TSF Decoder Architecture

(c) Transformer Encoder

Figure 3: TSF Model Architecture.

154 $\{2, \dots, |\mathbf{S}|\}$, where $\hat{\mathbf{H}}_{i-1} = D_1(E_1(\mathbf{H}_{i-1}))$. Generalising this to the k^{th} encoder-decoder pair,
 155 we minimise $\|\mathbf{H}_i - D_k(E_k(\mathbf{H}_i, \hat{\mathbf{H}}_{i-1}, \dots, \hat{\mathbf{H}}_{i-k}))\|_F^2 \forall i \in \{k+1, \dots, |\mathbf{S}|\}$, where, $\hat{\mathbf{H}}_{i-j} =$
 156 $D_j(E_j(\mathbf{H}_{i-j})) \forall j \in \{1, \dots, k\}$. We train the encoder-decoder pairs serially so that at the end
 157 of each training cycle, we can run inference using the trained model to generate the priors to train the
 158 next model.

159 We train our models using Adam and use a learning rate scheduler that reduces the learning rate by
 160 10% every 5 epochs, with a starting rate of 0.01. We use a batch size of 128 and train each model for
 161 100 epochs.

162 4 Results

163 In this section, we provide some preliminary results on the performance of the proposed TSF approach
 164 relative to two baselines; SF compression using an auto-encoder with the same architecture as the
 165 proposed TSF model but without using any past information and a 3GPP code-book based baseline,
 166 referred to as Rel-16 Type II codebook, which is part of the existing wireless standards. For the
 167 simulation setup, we consider the urban macro (UMa) channel scenario [5] at 4 GHz carrier frequency.
 168 The channels are collected from multiple BSs and multiple users moving at 10 km/hr speed, where
 169 each BS has 16 transmit antennas while each user has two receive antennas. We consider a bandwidth
 170 of 26 frequency components, i.e. $N_c = 26$. This makes each channel sample a complex tensor of
 171 dimensions $2 \times 16 \times 26$. For the TSF approach, we assume that the value of N is set to 3, so each
 172 user is utilizing up to 3 past channel samples in the compression process of the current sample.

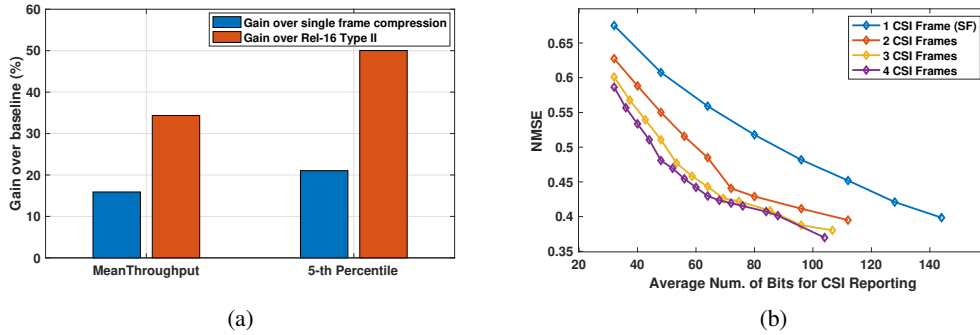


Figure 4: (a)Throughput results of the proposed TSF approach relative to Rel-16 Type II baseline and SF compression baseline. (b)NMSE performance comparison between SF and TSF compression methods with different number of feedback bits.

173 Fig. 4a shows the overall throughput performance gain of the proposed TSF approach over the SF
 174 baseline (blue bar) and the Rel-16 Type II 3GPP baseline (orange bar). To evaluate the throughput,
 175 we plug the 3 different models in a 3GPP-compliant wireless communication simulation pipeline
 176 and observe the impact of these methods on the overall data throughput or data rate achieved. For
 177 TSF compression, the first SF model used to compress the first sample has an overhead of 128 bits
 178 while the TSF model for the second, third and fourth sample has an overhead of 64 bits. This brings
 179 the average overhead of the TSF approach to 80 bits per reporting instance. The standalone SF
 180 compression model has an overhead of 80 bits for every reporting instance and likewise Rel-16 Type
 181 II. It can be seen that in terms of mean throughput, the proposed approach achieves 34% and 16%
 182 gain over the Rel-16 Type II and the SF approach, respectively. In addition, in terms of 5-th percentile,
 183 i.e., cell-edge (users with the worst channel conditions), throughput, the proposed approach achieves
 184 quite promising gains of 50% and 21% over the Rel-16 Type II and the SF approach, respectively.

185 We further showcase the performance of the TSF method by comparing the number of bits utilized
 186 to compress a CSI sample vs the achieved normalized mean squared error (NMSE) associated with
 187 CSI reconstruction. The rate-distortion-styled curve has been evaluated for the SF baseline (which
 188 assumes access to the current CSI sample only) and the TSF scheme with access to the current as
 189 well as up to 3 past CSI samples. Fig. 4b shows the NMSE performance against the average overhead
 190 associated with the CSI reporting. It can be seen that the proposed TSF approach with access to
 191 just 2 samples (current and 1 past sample) considerably outperforms SF compression. Further, as
 192 the number of past samples available for compression is increased, the performance improves more,
 193 while almost saturating when 3 past samples (a total of 4 CSI samples) are used for compression.

194 To achieve a reconstruction error benchmark of 0.4 NMSE, the SF scheme uses 144 bits on average,
 195 while the 4-sample TSF approach is able to achieve similar performance with just 88 bits. That's an
 196 overhead reduction of almost 39%.

197 5 Conclusion and Future Work

198 In this work, we introduce the interesting data compression paradigm associated with the problem
 199 of real-time channel state information (CSI) compression in wireless communication systems. To
 200 address the challenge, we propose to use the knowledge of past samples to better compress and
 201 reconstruct the channel data. We further propose a transformer-based compression model that
 202 effectively outperforms the single-sample methods and the existing methods currently utilized as
 203 part of the 5G standard. As future work, we plan to explore improved model architectures to better
 204 leverage information contained in past samples and remove the dependency associated with having
 205 the decoder as part of the encoding process.

206 References

- 207 [1] 3GPP TS 38.214. NR; Physical layer procedures for data. 3rd Generation Partnership Project;
 208 Technical Specification Group Radio Access Network.

- 209 [2] Teng-Hui Huang, Akshay Malhotra, and Shahab Hamidi-Rad. A deep learning method for joint
210 compression and unsupervised denoising of csi feedback. In *ICC 2023 - IEEE International Con-*
211 *ference on Communications*, pages 4150–4156, 2023. doi: 10.1109/ICC45041.2023.10279775.
- 212 [3] Wendong Liu, Wenqiang Tian, Han Xiao, Shi Jin, Xiaofeng Liu, and Jia Shen. Evcsinet:
213 Eigenvector-based CSI feedback under 3GPP link-level channels. *IEEE Wireless Communications*
214 *Letters*, 10(12):2688–2692, Sept 2021. doi: 10.1109/LWC.2021.3112747.
- 215 [4] Lu Lu, Geoffrey Ye Li, A Lee Swindlehurst, Alexei Ashikhmin, and Rui Zhang. An overview of
216 massive MIMO: Benefits and challenges. *IEEE Journal of Selected Topics in Signal Processing*,
217 8(5):742–758, Apr. 2014.
- 218 [5] Shu Sun, Theodore S Rappaport, Sundeep Rangan, Timothy A Thomas, Amitava Ghosh, Istvan Z
219 Kovacs, Ignacio Rodriguez, Ozge Koymen, Andrzej Partyka, and Jan Jarvelainen. Propagation
220 path loss models for 5G urban micro-and macro-cellular scenarios. In *IEEE 83rd Vehicular*
221 *Technology Conference (VTC Spring)*, pages 1–6, Nanjing, China, May 2016.
- 222 [6] Chao-Kai Wen, Wan-Ting Shih, and Shi Jin. Deep learning for massive MIMO CSI feedback.
223 *IEEE Wireless Communications Letters*, 7(5):748–751, Mar. 2018.