# Turbo-AI: Iterative Machine Learning Based Channel Estimation for 2D Massive Arrays

Yejian Chen, Jafar Mohammadi, Stefan Wesemann and Thorsten Wild

Bell Laboratories, Nokia, Lorenzstraße 10, D-70435 Stuttgart, Germany

Yejian.Chen@Nokia-Bell-Labs.com

*Abstract*—Recently, Machine Learning (ML) is recognized as an effective tool for wireless communications and plays an evolutionary role to enhance Physical Layer (PHY) of 5th Generation (5G) and Beyond 5G (B5G) systems. In this paper, we focus on the ML-based channel estimation for 2-Dimensional (2D) massive antenna arrays. Due to the extremely high computational requirement for 2D arrays with *Ordinary Training*, we exploit 2D Kronecker covariance model to perform *Subspace Training* for vertical and horizontal spatial domain independently, which achieves a complexity cost saving factor $\mathcal{O}(M^4N^4)/\mathcal{O}(MN^4 + NM^4)$ for ML with an $M \times N$ 2D-array. Furthermore, we propose an iterative training approach, referred to as *Turbo-AI*. Along with *Subspace Training*, the new approach can monotonically reduce the effective variance of additive noise of the observation, and update the Neural Network (NN) models by re-training. Furthermore, we propose a concept, named *Universal Training*. It allows to use one NN for a wide range of Signal-to-Noise-Ratio (SNR) operation points and spatial angles, which can greatly simplify *Turbo-AI* usage. Numerical results exhibit that *Turbo-AI* can tightly approach the *genie-aided* channel estimation bound, especially at low SNR.

## I. Introduction

Witnessing so many Artificial Intelligence (AI) inspired inventions and applications, no one is going to doubt, whether the AI techniques are changing our world. The recent rapidly developed AI techniques exhibit high potential for many technical challenges. In 2016 and 2017, AlphaGo, the AI program developed by DeepMind, defeated the human professional world champions twice in the most challenging strategy board game *Go*, which can be recognized as a revolutionary landmark of deep learning [1]–[2]. In communications society, AI is attracting enormous attentions from both industry and academia for enhancing the new concepts and novel technologies within the context of 5th Generation (5G) New Radio (NR) [3] and even opening new perspectives for Beyond 5G (B5G) systems. In this paper, we focus on the exploitation of AI for PHYsical (PHY) layer of wireless communications. It is well known that the performance of the key components of PHY layer are mathematically bounded by derivable optimums. Among them, the fundamental and the famous one is the Shannon limit [4]. Obviously, it is impossible to outperform the optimum bound by means of deep learning. Instead, AI introduces an alternative path to approach the optimum bound with low complexity. Thus, AI provides us additional degrees of freedom to trade off the performance and complexity in PHY layer design for 5G and B5G system. There are potential AI applications for PHY layer, enumerated in [5], such as blind channel decoding and data detection, modulation recognition, channel estimation, and many others.

Hence, our paper is motivated in this context by introducing AI for channel estimation to enhance PHY layer performance. Precise channel estimation belongs to one of the key technical prerequisites to support PHY layer operations, such as equalization and data detection, power control, precise Channel State Information (CSI) for transmit precoding and feedback, user pairing for Multiple-Input Multiple-Output (MIMO) and Non-Orthogonal Multiple Access (NOMA) system, and so on. In [6], a Neural Network (NN) based channel estimation approach is proposed, in which the conventional Minimum Mean Square Error (MMSE) channel estimator is equivalently represented under certain array-specific property assumptions as a two-layer NN. The conventional channel estimators, e.g. [7], are usually more complex than the Machine Learning (ML) based solutions [8] for similar performance. Especially, the computational requirement will be significantly increased for massive Multiple-Input Multiple-Output (mMIMO) system [9] in these conventional channel estimators with a large number of antenna elements. A proof-of-concept is provided in [10] provides for NN-based channel estimation, by means of real-time experiment and measurement. In [11], the ML-based channel estimation is considered for mMIMO with different antenna configurations.

In this paper, we will focus on the low-cost ML-based channel estimation with massive 2-Dimensional (2D) antenna panels, by exploiting the approximation in [12]. With the Kronecker channel model, the spatial covariance matrix of a 2D-array can be decomposed as the Kronecker product of a vertical covariance matrix and a horizontal covariance matrix, respectively. The training for the full spatial covariance matrix of the 2D-array, denoted as highly complex *Ordinary Training*, can be replaced by low complex *Subspace Training* with two independent NNs in horizontal and vertical domains, respectively. Different combining strategies for channel estimates, obtained after the *Subspace Training* in vertical and horizontal subspaces, can improve the channel estimation with reduced complexity. Furthermore, after combining the estimates from independent subspaces, the estimation error, which can be interpreted as additive noise, is still a Gaussian random variable, but with reduced variance. Let the initial observation be replaced by the less noisy channel estimates, and repeat the *Subspace Training* and combining, to approach the optimum performance iteratively. Throughout this

paper, this procedure is referred to as *Turbo-AI*. It can be demonstrated through link level simulation that the *genie-aided* bound can be tightly reached. Furthermore, we propose the concept *Universal Training* to enhance the robustness of *Turbo-AI* for practical implementation.

This paper is organized as follows. In Section II we focus on the system model. An overview of *Subspace Training*, *Turbo-AI* and *Universal Training* will be presented in Section III. In Section IV, numerical results of link layer performance are provided. Finally, Section V renders some conclusions.

## II. SYSTEM MODEL

As depicted in Fig. 1, let us consider a 2D-array system with $M$ vertical elements and $N$ horizontal elements. The samples $k = 1, \cdots, K$ are collected from time domain as temporal symbols or from frequency domain as subcarriers. Hence, the Single-Input Multiple-Output (SIMO) system can be represented by following equation as

$$\mathbf{Y}_k = \mathbf{H}_k + \mathbf{Z}_k, \tag{1}$$

where the $M \times N$ matrix $\mathbf{Y}_k$ and $\mathbf{H}_k$ stands for the observations and wireless channel, respectively. The matrix $\mathbf{Z}_k$ denotes the Additive White Gaussian Noise (AWGN), element-wise characterized by noise power spectral density $N_0 = 2\sigma^2$. In [12], it is exhibited that the covariance matrix $\mathbf{R}$ of a 2D-
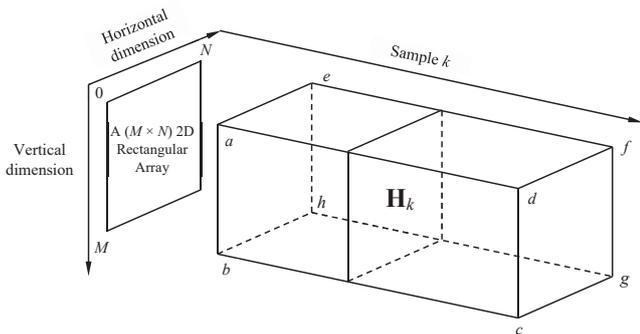


Fig. 1. Signal model with a 2D antenna panel

array can be approximated as

$$\mathbf{R} \approx \mathbf{R}_h \otimes \mathbf{R}_v, \tag{2}$$

where $\mathbf{R}_h$ and $\mathbf{R}_v$ denote the common horizontal and vertical covariance matrix in the sense of 1D-arrays in row and column, respectively. The covariance matrices $\mathbf{R}_v$, $\mathbf{R}_h$ and $\mathbf{R}$ are thus $M \times M$, $N \times N$ and $MN \times MN$ matrices. The spatial correlation between $i$-th and $j$-th antenna in $\mathbf{R}_h$ [13], denoted as $\mathbf{R}_h^{(i,j)}$, can be computed as

$$\mathbf{R}_h^{(i,j)} = \frac{1}{\sqrt{2}\varrho} \int_{-\pi}^{\pi} \exp\left[ -\frac{\sqrt{2}|\phi|}{\varrho} + j\frac{2\pi}{\lambda}\Delta d_{i,j}\sin(\theta + \phi)\right] \mathrm{d}\phi, \tag{3}$$

where $\Delta d_{i,j}$ denotes the distance between both antenna elements $i$ and $j$. Parameters $\varrho$ and $\theta$ stand for the angular spread and Direction of Arrival (DoA), respectively. The given

random variable $\phi$ obeys Laplacian power angular distribution concerning $\varrho$ and $\theta$. Equation (3) is also valid for vertical spatial domain $\mathbf{R}_v^{(i,j)}$. Considering the independent and identically distributed (i.i.d.) Rayleigh fading channel, it holds

$$\mathrm{vec}(\mathbf{H}_k) = \mathbf{R}^{0.5}\mathrm{vec}(\mathbf{H}_{w,k}). \tag{4}$$

The elements in the $M \times N$ matrix $\mathbf{H}_{w,k}$ follow i.i.d. Rayleigh distribution with $E[|\mathbf{H}_{w,k}|^2] = 1$, representing the sample $k$ randomly selected from time or frequency domain. Thus, the channel matrix $\mathbf{H}_k$ can be equivalently computed as

$$\mathbf{H}_k = \mathbf{R}_v^{0.5}\mathbf{H}_{w,k}(\mathbf{R}_h^{0.5})^T. \tag{5}$$

## III. CHANNEL ESTIMATION AND MACHINE LEARNING

In this section, let us focus on the channel estimation and training for machine learning. First of all, let us discuss about the benchmark solution of channel estimation. For a linear signal model in presence of Gaussian noise, the Linear Minimum Mean Square Error (LMMSE) channel estimator is regarded as the optimal linear estimator, in sense of being capable of minimizing the Mean Square Error (MSE). For a 2D-array system, depicted in Fig. 1, the MMSE weight matrix can be computed as

$$\mathbf{W}_{\text{genie}} = \mathbf{R}(\mathbf{R} + \sigma^2\mathbf{I}_{MN})^{-1}. \tag{6}$$

If the covariance matrix $\mathbf{R}$ is *a priori* known, it is referred to as *genie-aided* solution. Nevertheless, the drawback is obvious as well: if the number of antenna elements $M$ and $N$ become large, the matrix inversion in (6) can be computationally intensive.

### A. Estimation Through Subspaces

Let us re-consider the problem from the view point of horizontal and vertical subspaces as depicted in Fig. 2. We can stack the data as slices horizontally and vertically to create two observations. Let us define the $N \times M$ matrix $\mathbf{H}_h$ and $M \times N$ matrix $\mathbf{H}_v$ as

$$\mathbf{H}_h = \left[\mathbf{h}_1^{(h)} \cdots \mathbf{h}_M^{(h)}\right] \tag{7}$$

$$\mathbf{H}_v = \left[\mathbf{h}_1^{(v)} \cdots \mathbf{h}_N^{(v)}\right], \tag{8}$$

where $M$ samples of $N \times 1$ vector $\mathbf{h}_m^{(h)}$ and $N$ samples of $M \times 1$ vector $\mathbf{h}_n^{(v)}$ are obtained in horizontal and vertical subspaces, respectively. Consider MMSE channel estimation for both subspaces individually, it holds

$$\hat{\mathbf{H}}_h = \mathbf{W}_h\mathbf{Y}^T = \mathbf{R}_h(\mathbf{R}_h + \sigma^2\mathbf{I}_N)^{-1}\mathbf{Y}^T \tag{9}$$

$$\hat{\mathbf{H}}_v = \mathbf{W}_v\mathbf{Y} = \mathbf{R}_v(\mathbf{R}_v + \sigma^2\mathbf{I}_M)^{-1}\mathbf{Y} \tag{10}$$

where $\mathbf{W}_h$ and $\mathbf{W}_v$ denote $N \times N$ and $M \times M$ MMSE weighting matrices, and the observation is denoted as $M \times N$ matrix $\mathbf{Y}$. Let us simply consider the arithmetic mean to combine the estimates from vertical and horizontal subspaces. It holds

$$\hat{\mathbf{H}}_a = 0.5(\hat{\mathbf{H}}_h^T + \hat{\mathbf{H}}_v) = 0.5(\mathbf{Y}\mathbf{W}_h^T + \mathbf{W}_v\mathbf{Y}). \qquad (11)$$

The result in (11) can be further reformulated as a vectorized expression as

$$\mathrm{vec}(\hat{\mathbf{H}}_a) = 0.5(\mathbf{I}_N \otimes \mathbf{W}_v + \mathbf{W}_h \otimes \mathbf{I}_M)\mathrm{vec}(\mathbf{Y}), \qquad (12)$$

where the $MN \times MN$ matrix $\mathbf{W}_a = 0.5(\mathbf{I}_N \otimes \mathbf{W}_v + \mathbf{W}_h \otimes \mathbf{I}_M)$ is the effective weighting with respect to the subspace combining with arithmetic mean. Similarly, we can exploit the geometric mean to perform the subspace combining. It holds

$$\mathrm{vec}(\hat{\mathbf{H}}_g) = (\mathbf{I}_N \otimes \mathbf{W}_v)^{0.5}(\mathbf{W}_h \otimes \mathbf{I}_M)^{0.5}\mathrm{vec}(\mathbf{Y}) \qquad (13)$$
$$= (\mathbf{W}_h^{0.5} \otimes \mathbf{W}_v^{0.5})\mathrm{vec}(\mathbf{Y}),$$

where the $MN \times MN$ matrix $\mathbf{W}_g = \mathbf{W}_h^{0.5} \otimes \mathbf{W}_v^{0.5}$ is the effective weighting with geometric mean. With the properties of Kronecker product, it yields,

$$\hat{\mathbf{H}}_g = \mathbf{W}_v^{0.5}\mathbf{Y}(\mathbf{W}_h^{0.5})^T. \qquad (14)$$

Notice that the same data is reordered as the inputs for the horizontal NN and vertical NN, respectively. Nevertheless, the reshuffling patterns of corresponding AWGN components within both observations exhibit strong independence, represented as colorful spots in Fig. 2 for examples. The gain from
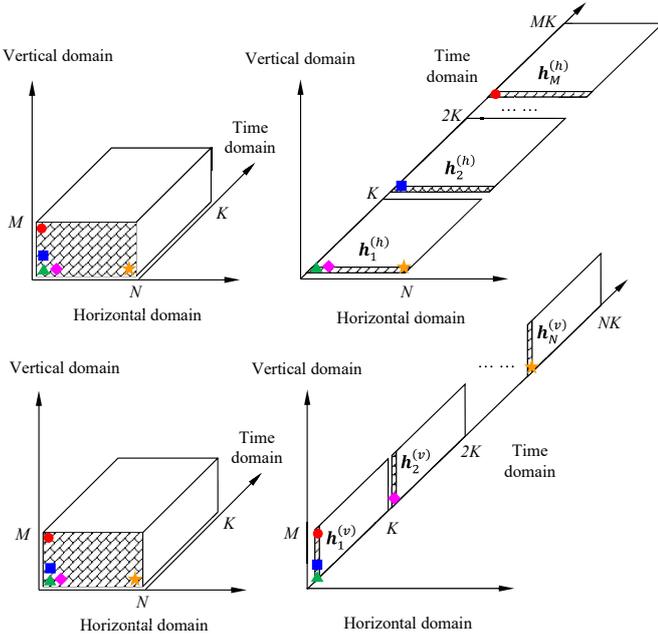


Fig. 2. Stacking $M$ vertical slices in horizontal spatial domain and stacking $N$ horizontal slices in vertical spatial domain

combining the outputs of two NNs after *Subspace Training* comes from the fact that the output randomness of each estimated element is originated from training these independent elements. It can be exhibited that the variance post-processing noise can be reduced after subspace combining, e.g. arithmetic mean. Let $\mathrm{VAR}[\mathbf{Z}_a^{(i,j)}]$ the variance of the AWGN component

on $i$-th row and $j$-th column in matrix $\mathbf{Z}_a$ after exploiting arithmetic mean. It holds,

$$\frac{\rho}{4}N_0 < \mathrm{VAR}[\mathbf{Z}_a^{(i,j)}] \le N_0 \qquad (15)$$

with $\rho = \sum_{m=1}^M E[|\mathbf{W}_v^{(i,m)}|^2] + \sum_{n=1}^N E[|\mathbf{W}_h^{(j,n)}|^2]$ and $0 \le \rho \le 2$. The variance of the original AWGN component is $\mathrm{VAR}[\mathbf{Z}^{(i,j)}] = N_0 = 2\sigma^2$. The proof of (15) will be provided in *Appendix A*.

In order to explore the property of arithmetic mean and geometric mean, let us introduce the notations $\mathbf{y} = \mathrm{vec}(\mathbf{Y})$, $\mathbf{A} = \mathbf{I}_N \otimes \mathbf{W}_v$ and $\mathbf{B} = \mathbf{W}_h \otimes \mathbf{I}_M$. With the arithmetic mean, the deviation of the channel estimates from the *genie-aided* MMSE channel estimator can be computed as

$$D_a = \frac{1}{MN}\mathbf{y}^H\Big[\frac{1}{2}(\mathbf{A} + \mathbf{B}) - \mathbf{W}_{\mathrm{genie}}\Big]^2\mathbf{y} \qquad (16)$$

and with geometric mean, the deviation can be similarly computed as

$$D_g = \frac{1}{MN}\mathbf{y}^H(\mathbf{A}^{0.5}\mathbf{B}^{0.5} - \mathbf{W}_{\mathrm{genie}})^2\mathbf{y}. \qquad (17)$$

Furthermore, it holds

$$D_a - D_g = \frac{1}{MN}\mathbf{y}^H\mathbf{L}(\mathbf{P} + \mathbf{Q})\mathbf{y} \qquad (18)$$

with $\mathbf{L} = (\mathbf{A}^{0.5} - \mathbf{B}^{0.5})^2$, $\mathbf{P} = \mathbf{A} + \mathbf{B} - 2\mathbf{W}_{\mathrm{genie}}$ and $\mathbf{Q} = 2(\mathbf{A}^{0.5}\mathbf{B}^{0.5} - \mathbf{W}_{\mathrm{genie}})$. According to our computational analysis for our dataset with both i.i.d. Rayleigh channel assumption and 3GPP compliant channel model, the following holds

$$D_a \ge D_g, \qquad (19)$$

which indicates that the geometric mean can provide the channel estimates with better quality. For computing the geometric mean $\hat{\mathbf{H}}_g$, the matrix square root operation $(\cdot)^{0.5}$ has to be invoked for matrices $\mathbf{W}_v$ and $\mathbf{W}_h$, which introduces additional complexity due to the non-linear operation.

### B. Neural Network and Complexity

It is exhibited in [6] that a MMSE channel estimator can be approximated as a two-layer NN network. Thus, in a $M \times N$ 2D-array system, with a $MN \times MN$ spatial covariance matrix $\mathbf{R}$, the number of input parameters is $2M^2N^2$ real values, if the sample covariance matrices are exploited for the training. Obviously, the complexity will be significantly increased, if a massive 2D-array is deployed. In [6], matrix transformation and decomposition are considered to reduce the number of NN parameters for antenna array with special structure. With the analysis in the previous subsection, the Kronecker model allows us to exploit the spatial decomposition naturally, and process the horizontal and vertical subspaces individually as Uniform Linear Array (ULA) antennas. As depicted in Fig. 4, we exploit two NNs in parallel, each of which has simple structure as a two-layer network with ReLu activation. The same $M \times N \times K$ data structure, see Fig. 2, will be reconstructed

and delivered to both NNs as the observations for vertical and horizontal subspaces. The input and the output of the NNs are spatial sample covariance matrices $\mathbf{R}_v$, $\mathbf{R}_h$ and the weight matrices $\mathbf{W}_v$, $\mathbf{W}_h$, respectively. Fig. 3 exhibits how to perform the combining for vertical and horizontal subspaces in sense of arithmetic mean and geometric mean, according to the analysis in the previous subsection. Exploiting the full MMSE
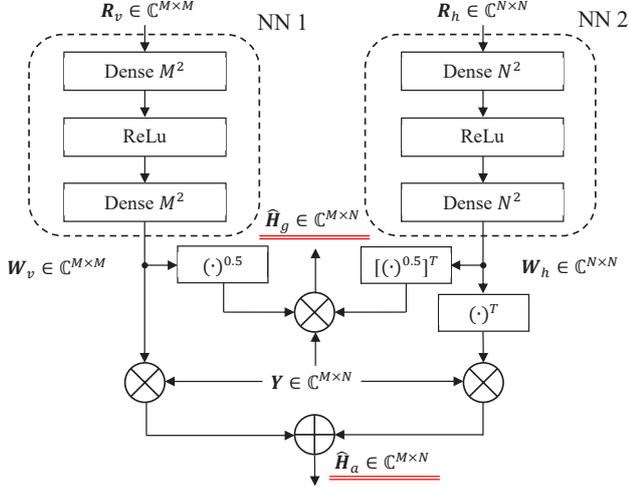


Fig. 3. Two parallel NNs for *Subspace Training*

*Ordinary Training* for the $M \times N$ 2D-array, two concatenated $(MN)^2$-In-$(MN)^2$-Out dense layers in one NN have to be considered. The number of real values as the input parameters is $2M^2N^2$. The number of the stored real values of both dense layers is $8M^4N^4$. Considering the cost of parameter passing, ReLu operation and real/imaginary part matrix multiplication, the complexity of the full MMSE *Ordinary Training* is $\mathcal{O}(M^4N^4)$. In Fig. 3, two parallel NNs are depicted, each of which has two concatenated dense layers. The number of real values as the input parameters of the network is $2(M^2 + N^2)$. The number of stored real values of both NNs is $4(M^4 + N^4)$. The total cost of *Subspace Training* to involve parameter passing, ReLu operation, real/imaginary part matrix multiplication and subspace combining is $\mathcal{O}(MN^4 + NM^4)$. Let us define a cost saving function $\eta(M, N)$ as

$$\eta(M, N) \approx \frac{\mathcal{O}(M^4N^4)}{\mathcal{O}(MN^4 + NM^4)}. \qquad (20)$$

For a $(M = 4, N = 8)$ and a $(M = 8, N = 16)$ 2D-array systems, the cost saving coefficients with *Subspace Training* are approximately 50 and 450, respectively.

### C. Turbo-AI

*Turbo-AI* is an iterative procedure, which consists of an iterative training stage and an iterative estimation stage. Let us explore the operation principle of *Turbo-AI* in this subsection. After being combined in horizontal and vertical subspaces, the channel estimates, obtained by the NN described in previous subsection, still exhibit residual additive noise, which is Gaussian, but with a reduced variance, compared to the noise in
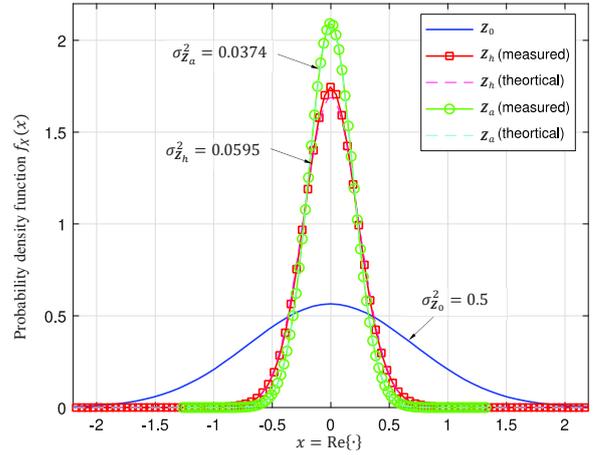


Fig. 4. Probability density function of post-processing noise after *Subspace Training* and combining at 0dB SNR

the channel observation input. This property can be illustrated in Fig. 4. Let us study the Probability Density Function (PDF) of the real part signal of the raw AWGN $\mathbf{Z}_0$, as well as that of the channel estimation error after horizontal *Subspace Training* $\mathbf{Z}_h$ and the channel estimation error after arithmetic mean $\mathbf{Z}_a$. Notice that we can plot the PDFs straightforwardly by means of measuring $\mathbf{Z}_h$ and $\mathbf{Z}_a$, or alternatively, plot the theoretical Gaussian distribution functions by the measured post-processing noise variance. The PDFs, obtained from both approaches, coincides to each other very well, which confirms that the channel estimation errors still obey Gaussian distribution with reduced variance. Let these noisy channel estimates replace old observations of the training data, and produce the new channel estimates repeatedly with the same NNs by means of *Subspace Training*. This iterative procedure is referred to as *Turbo-AI*, due to the iteratively achievable processing gain, which eventually makes the channel estimates get very close to the solution produced by *genie-aided* channel estimator.

Note that the *Turbo-AI* procedure introduces a new way to train the NN. As opposed to training the whole layers (from all iterations) at once, we train each layer to the same final labels, from bottom to top. This acts like a regularization by limiting the degrees of freedom that our network can take. It further boosts the convergence speed. Moreover, the chance of getting stuck in a local minima is a lot smaller this way, as in each stage we are dealing with a small 2-layered NN. This solution is not entirely heuristic since it is based on the observation that the dominant randomness at the input and the output of *Turbo-AI* is Gaussian with different variances.

In Fig. 5, the training phase for $i$-th iteration of *Turbo-AI* is presented. Notice that a dedicated set of NN weights, so-called NN-models, can be stored for the validation from other users with the same spatial characteristics. In the inference phase, depicted in Fig. 6, the user should estimate the initial SNR, decide to start with the best matched pre-trained NN-models $\mathbf{W}_{v,i}$ and $\mathbf{W}_{h,i}$, and replace both *Subspace Training* blocks in Fig. 5 to estimate the channel iteratively. Considering the potential complexity for hardware design, introduced by the

square root operations to calculate geometric mean for both $\mathbf{W}_{v,i}$ and $\mathbf{W}_{h,i}$ during the iterations, we adopt arithmetic mean for subspace combining to realize *Turbo-AI* as an overall relatively low complex solution throughout this paper. In *Appendix B*, the proof of the applicability of *Turbo-AI* will be provided.
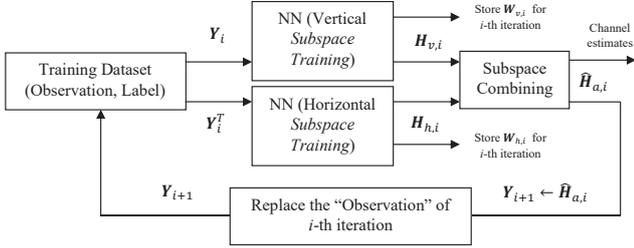
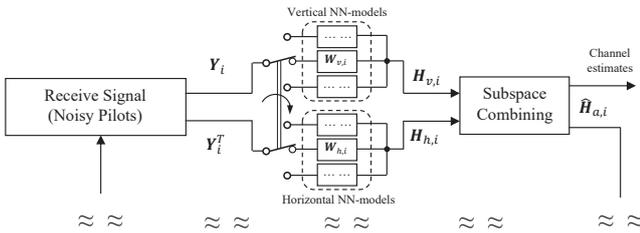Fig. 5. Training phase: Block diagram for $i$-th iteration training of *Turbo-AI*

Fig. 6. Inference phase: Block diagram for $i$-th iteration training of *Turbo-AI*

### D. Universal Training for Turbo-AI

*Turbo-AI* can iteratively improve the channel estimation. The training phase, depicted in Fig. 5, executes dedicated training for one direction with respect to vertical and horizontal spatial domain. For another direction, the vertical and horizontal NNs have to be retrained again. For the same reason, a set of vertical and horizontal NN-models have to be stored, because the effective SNR will increase during the iterations. Especially, the estimation of initial SNR to start the *Turbo-AI* in Fig. 6 should be precise. Otherwise, it can also cause the mismatch of the NN-models during the iterations. This motivates us to think about the question, whether it is possible to span the parameter space of the training, and make the NN-models of *Turbo-AI* become more universal and more robust against the parameter estimation error in the practice. Thus, we name this concept *Universal Training*. In Fig. 7, the
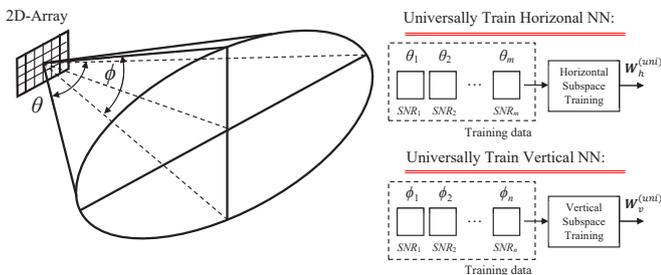
Fig. 7. *Universal Training* with respect to SNR and spatial domain

basic idea of *Universal Training* is illustrated. For training the horizontal NN, the training data consists of the data batches for different horizontal angle $\theta_m$ and SNR values randomly. During the training stage, the NN tries to find out a model to deliver sub-optimum solution for all possible combination of parameters $\theta_m$ and $\text{SNR}_m$. The randomization for the data combination is necessary to prevent ML from falling in a local optimum. Thus, the learning procedure takes longer than the conventional training.

After the training, a universal NN-model is ready for channel estimation within the parameter space spanned by data combination $\theta_m$ and $\text{SNR}_m$. Similarly, the vertical NN-model can be universally learned in another parameter space, spanned by parameter $\phi_n$ and $\text{SNR}_n$. As a matter of fact, the *Universal Training* is sub-optimal and can cause performance loss, if more parameters are introduced in training. On the other side, with *Universal Training*, *Turbo-AI* requires only a couple of vertical NN-models and horizontal NN-modes, or even one vertical NN-model and one horizontal NN-model to enable the iterations for arbitrary combination of $\theta_m$, $\phi_n$ and SNR value. This can introduce huge cost saving and flexibility in hardware design for a practical implementation.

## IV. SIMULATIONS AND NUMERICAL RESULTS

In this section, let us focus on the performance of *Turbo-AI* with link level simulations. Throughout our investigation, the spatial i.i.d. Rayleigh fading channel, modeled by (4) and (5), is adopted to generate channel responses. In Table I, the parameters are summarized.

TABLE I
LINK LEVEL SIMULATION AND TRAINING PARAMETERS

| | |
|---|---|
| Number of vertical elements | $M = 8$ |
| Number of horizontal elements | $N = 16$ |
| Antenna spacing | Half of wavelength |
| Vertical Angular spread | $\pi/180$ |
| Horizontal Angular spread | $\pi/90$ |
| Optimizer for ML | Adam |
| Learning rate | $\alpha = 0.009$ |
| Hyper-parameters | $\beta_1 = 0.9, \beta_2 = 0.999$ |
| Activation function | ReLu |
| Training data set size | $K = 500000$ (Time domain) |

In Fig. 8, the performance of *Turbo-AI* is presented for user validation with the stored NN-models through *Subspace Training*, as depicted in Fig. 6. With the cooperation between vertical and horizontal NN, the arithmetic mean based subspace combining can approach the *genie-aided* MMSE estimator very fast, especially at the first and the second iteration. Even at relatively low SNR, e.g. SNR at 0dB, the performance gap to the *genie-aided* bound is only 0.13dB after four iterations. In Fig. 9, the measured PDFs of the channel estimates from horizontal domain after the iterations are presented, which are a set of Gaussian distribution curves with reduced variance.

In Fig. 10, the *Turbo-AI* performance with *Universal Training* is summarized. As mentioned in previous section, the more parameters are introduced in the training, the more performance loss should be expected. In Fig. 10, we consider
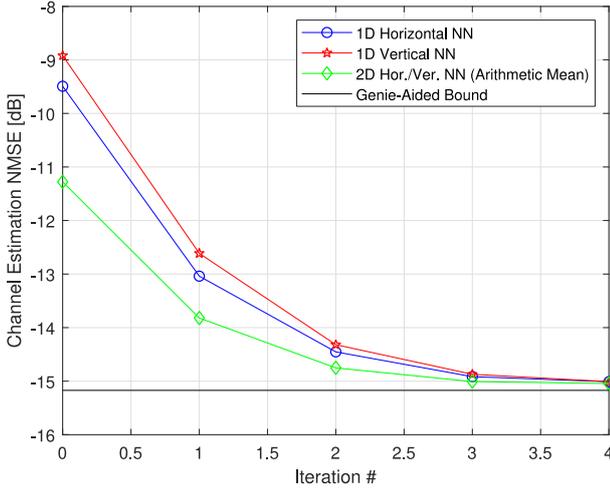
Fig. 8. *Turbo-AI* performance: channel estimation NMSE versus iterations, a 2D-array with $M = 8$ and $N = 16$, i.i.d. Rayleigh fading, SNR at 0dB
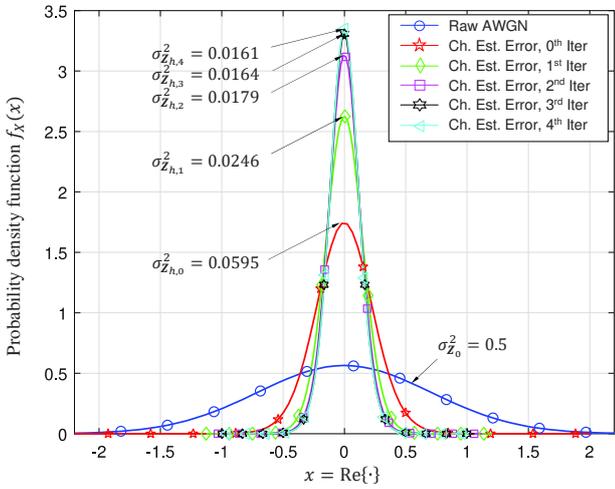


Fig. 10. *Turbo-AI* performance: dedicated training versus universal training, a 2D-array with $M = 8$ and $N = 16$, i.i.d. Rayleigh fading, SNR at 0dB

about 3.98dB. Nevertheless, comparing to the Least Square (LS) channel estimator, which achieves only 0dB channel estimation NMSE at 0dB SNR, *Turbo-AI* with one NN-model converges at approximately -11dB, which is quite satisfactory. In Fig. 11, the performance of diverse channel estimators is



Fig. 9. *Turbo-AI* performance: track the PDFs of horizontal estimates, a 2D-array with $M = 8$ and $N = 16$, i.i.d. Rayleigh fading, SNR at 0dB



Fig. 11. Channel estimation NMSE versus SNR, a 2D-array with $M = 8$ and $N = 16$, i.i.d. Rayleigh fading

three parameters, namely the range of SNR, horizontal and vertical spatial coverage $\theta$ and $\phi$, respectively. First of all, let us carry out the *Universal Training* only for SNR with the range from 0dB to 15dB, and assume that $\theta$ and $\phi$ are *a priori* known directions. Comparing to the dedicated training, the performance loss uniquely introduced by the universal SNR training is about 1.20dB. The universal SNR training is quite meaningful for *Turbo-AI*, since only one NN-model is required for the *Turbo-AI*, due to being able to adapt to all possible effective SNR values during the iterations. Then, let us increase the coverage horizontally and vertically with respect to both parameters $\theta$ and $\phi$. We can observe further obvious degradation. Finally, we carry out *Universal Training* for the whole sector, namely $-60° < \theta_m < 60°$, $30° < \phi_n < 90°$ and SNR $\in [0dB, 15dB]$. It is quite convenient to exploit *Turbo-AI* to estimate the channel with one vertical NN-model and one horizontal NN-model for the whole sector. Being the trade-off between complexity and performance, the performance gap to dedicated training is
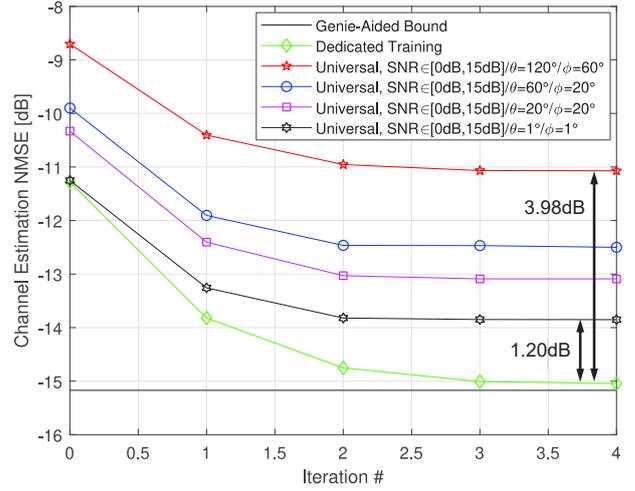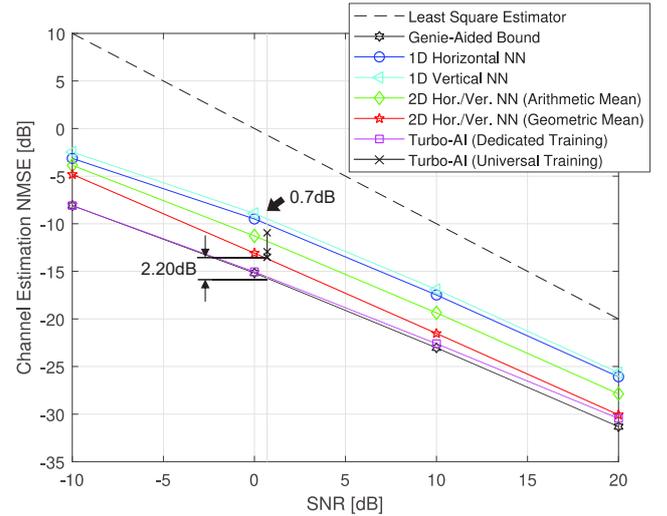
summarized. The ML-based channel estimators are bounded by LS estimator and *genie-aided* MMSE estimator. The 2D *Subspace Training* based channel estimators outperform the horizontal or vertical NNs. Especially, the 2D *Subspace Training* and combining with geometric mean outperforms that with arithmetic mean, which numerically verifies equation (15) to (17) in the previous section. *Turbo-AI* with dedicated training exhibits powerful performance. With four iterations, the *genie-aided* bound can be tightly approached, especially for low SNR region. Further, with the training parameters $10° < \theta_m < 30°$, $40° < \phi_n < 60°$ and SNR $\in [0dB, 15dB]$, we verify the pre-trained universal NN-model for a SNR value at 0.7dB, a strange value, which is not included in our selected

universal training data set at all. Simulation exhibits that this universal NN-model for *Turbo-AI* is still applicable. With three iterations, the performance gap to *genie-aided* bound is only 2.20dB.

## V. CONCLUSION

In this paper, we propose a set of ML-based channel estimation approaches for 2D massive arrays. Exploiting the 2D Kronecker channel covariance model, we can carry out vertical and horizontal *Subspace Training*, to estimate the channel in sense of arithmetic mean and geometric mean. *Subspace Training* provides a sub-optimal solution with much lower complexity than that of full array size inputs. Although the cost saving from *Subspace Training* with respect to hardware implementation is huge, a performance gap does still exist, considering the *genie-aided* MMSE channel estimator. Thus, we further trade off the complexity and performance, and introduce an iterative ML-based channel estimation approach, referred to as *Turbo-AI*, in which the observation for ML can be iteratively de-noised by means of the cooperation of horizontal and vertical *Subspace Training* and combining. Numerical results exhibit that the *genie-aided* bound performance can be tightly approached by *Turbo-AI*, even at very low SNR. From the view point of system design, *Subspace Training* and *Turbo-AI* provide diverse application possibilities to support both offline learning and online learning with scalable complexity. Furthermore, in order to enhance the flexibility of hardware design, *Universal Training* is introduced as a counterpart for training with dedicated parameters. In *Universal Training*, we span the parameter space in training phase and make the NN-model adapt to a wide range of parameters. For instance, the parameter space can be spanned by vertical spatial coverage $\phi$, horizontal spatial coverage $\theta$ and SNR. It is demonstrated that flexible performance can be achieved by adjusting the range of the parameter space. Thus, *Universal Training* provides additional degrees of freedom to facilitate practical implementation. For the future perspective, we will validate the concepts, e.g. *Subspace Training*, *Turbo-AI* and *Universal Training*, in 3GPP spatial channel model for a multicarrier system. The channel correlation in space, time and frequency can be jointly considered. These additional dimensions are assumed to be able to provide further enhancement for ML-based channel estimation.

## APPENDIX

### A. Proof of (15)

Before proving equation (15), let us introduce several *Lemmas* to define the NN, which is exploited for channel estimation. We inherit the notations from previous sections. The observations and labels of a 2D-array with $M$ vertical elements and $N$ horizontal elements are denoted as $M \times N$ matrices $\mathbf{Y}$ and $\mathbf{H}$, respectively, with $\mathbf{H}_v = \mathbf{H}$ and $\mathbf{H}_h = \mathbf{H}^T$. Let $\mathbf{Z}$ denote the $M \times N$ AWGN matrix, with covariance matrix $\mathbf{R}_Z = N_0 \mathbf{I}_M = 2\sigma^2 \mathbf{I}_M$. Without loss of generality,

we adopt the vertical subspace for defining the *Lemmas*. The same *Lemmas* hold for horizontal domain as well.

*Lemma 1:* When a neural network is trained over adequate amount of channel observation, sampled from one given distribution, this neural network is called converged, if following properties hold,

$$E[\|\mathbf{W}_v\mathbf{Y} - \mathbf{H}\|_F^2] \le E[\|\mathbf{Z}\|_F^2] \tag{21}$$

$$\|\mathbf{W}_v - \mathbf{W}_{v,\text{genie}}\|_F < \epsilon \tag{22}$$

where $\mathbf{W}_v$ denotes the $M \times M$ weighting matrix and $\mathbf{W}_v\mathbf{Y} - \mathbf{H}$ represents the zero-mean estimation error. The operator $\|\mathbf{A}\|_F$ denotes the Frobenius norm of matrix $\mathbf{A}$. Matrix $\mathbf{W}_{v,\text{genie}}$ denotes the optimal MMSE weighting matrix in vertical subspace, defined in (6), and $\epsilon$ is small non-negative value. *Lemma 1* follows from the the MSE loss function.

Equation (22) in *Lemma 1* is difficult to be proven mathematically, as the proof of convergence to the optimum point generally does not exist for an arbitrary architecture and loss function. The conditions in *Lemma 1* are based on our observations of training the NNs on many different scenarios and channel data distributions, corrupted by Gaussian noise. We do not claim these conditions and properties for an arbitrary type of data. Nevertheless, the procedure for the core NN design based on [6] and our simulations indicates that it holds for Gaussian inputs. Furthermore, we treat $\mathbf{W}_v$ as a deterministic matrix in proving (15), if the NN is converged and the matrix $\mathbf{W}_{v,\text{genie}}$ itself is deterministic as well.

*Lemma 2:* From *Lemma 1*, it straightforwardly holds

$$E[\|\mathbf{W}_v\mathbf{Z}\|_F^2] \le E[\|\mathbf{Z}\|_F^2].$$

Furthermore, for an arbitrary vertical array element $i$, with $i \in \{1, 2, \cdots, M\}$, through the spatial filtering with $\mathbf{W}_v$, the post-processing noise will be reduced, if the convergence of NN is reached, namely

$$E[\|\mathbf{W}_v^{(i,:)}\mathbf{Z}\|_F^2] \le E[\|\mathbf{Z}^{(i,:)}\|_F^2], \tag{23}$$

where $\mathbf{W}_v^{(i,:)}$ and $\mathbf{Z}^{(i,:)}$ denote the $i$-th row of matrices $\mathbf{W}_v$ and $\mathbf{Z}$, respectively.

With the *Lemma 1* and *Lemma 2* introduced above, now we can start to prove equation (15).

*Proof:* With *Lemma 2*, we obtain

$$E[\|\mathbf{W}_v\mathbf{Z}\|_F^2] = E[\text{Tr}(\mathbf{W}_v\mathbf{Z}\mathbf{Z}^H\mathbf{W}_v^H)] = 2\sigma^2\text{Tr}(\mathbf{W}_v\mathbf{W}_v^H)$$

$$\le E[\|\mathbf{Z}\|_F^2] = E[\text{Tr}(\mathbf{Z}\mathbf{Z}^H)] = 2\sigma^2 M. \tag{24}$$

Thus, it yields,

$$\mathrm{Tr}(\mathbf{W}_v \mathbf{W}_v^H) \le M \Leftrightarrow \sum_{i=1}^{M} \sum_{m=1}^{M} |\mathbf{W}_v^{(i,m)}|^2 \le M, \qquad (25)$$

where $\mathbf{W}_v^{(i,m)}$ denotes the weighting coefficient on $i$-th row and $m$-th column of $\mathbf{W}_v$. Exploiting the same operations for (23), considering (24) and (25), we obtain,

$$0 \le \sum_{m=1}^{M} |\mathbf{W}_v^{(i,m)}|^2 \le 1, \; \forall i \in \{1, 2, \cdots, M\}. \qquad (26)$$

Similarly, for horizontal case, we have

$$0 \le \sum_{n=1}^{N} |\mathbf{W}_h^{(j,n)}|^2 \le 1, \; \forall j \in \{1, 2, \cdots, N\}. \qquad (27)$$

Considering arithmetic mean to combine the estimates from vertical and horizontal subspaces, we have

$$\hat{\mathbf{H}}_a^{(i,j)} = \frac{1}{2}\big[\hat{\mathbf{H}}_v^{(i,j)} + \hat{\mathbf{H}}_h^{(j,i)}\big]. \qquad (28)$$

Furthermore, it holds,

$$E\big[|\hat{\mathbf{H}}_a^{(i,j)}|^2\big] = \frac{1}{4} E\big[|\hat{\mathbf{H}}_v^{(i,j)} + \hat{\mathbf{H}}_h^{(j,i)}|^2\big]$$
$$= \underbrace{\frac{1}{4} E\big[|\hat{\mathbf{H}}_v^{(i,j)}|^2\big]}_{A} + \underbrace{\frac{1}{4} E\big[|\hat{\mathbf{H}}_h^{(j,i)}|^2\big]}_{B} + \underbrace{\frac{1}{2}\mathrm{Re}\big\{E\big[\hat{\mathbf{H}}_v^{(i,j)}\hat{\mathbf{H}}_h^{(j,i)^*}\big]\big\}}_{C}.$$
$$(29)$$

Now, let us resolve the terms $A$, $B$ and $C$ in the equation above. It holds,

$$A = \frac{1}{4} E\Big[\big|\sum_{m=1}^{M} \mathbf{W}_v^{(i,m)}\mathbf{H}_v^{(m,j)}\big|^2\Big] + \frac{\sigma^2}{2} \sum_{m=1}^{M} |\mathbf{W}_v^{(i,m)}|^2 \quad (30)$$

$$B = \frac{1}{4} E\Big[\big|\sum_{n=1}^{N} \mathbf{W}_h^{(j,n)}\mathbf{H}_h^{(n,i)}\big|^2\Big] + \frac{\sigma^2}{2} \sum_{n=1}^{N} |\mathbf{W}_h^{(j,n)}|^2 \qquad (31)$$

$$C = \frac{1}{2}\mathrm{Re}\Big\{E\Big[\sum_{m=1}^{M} \mathbf{W}_v^{(j,m)}\mathbf{H}_v^{(m,i)} \sum_{n=1}^{N} \mathbf{W}_h^{(j,n)^*}\mathbf{H}_h^{(n,i)^*}\Big]\Big\}$$
$$+ \sigma^2 \mathbf{W}_v^{(i,i)}\mathbf{W}_h^{(j,j)}. \qquad (32)$$

Let us pick out the post-processing AWGN related terms from $A$, $B$ and $C$. It holds,

$$\mathrm{VAR}\big[\mathbf{Z}_a^{(i,j)}\big] = \frac{\sigma^2}{2} \sum_{m=1}^{M} |\mathbf{W}_v^{(i,m)}|^2 + \frac{\sigma^2}{2} \sum_{n=1}^{N} |\mathbf{W}_h^{(j,n)}|^2$$
$$+ \sigma^2 \mathbf{W}_v^{(i,i)}\mathbf{W}_h^{(j,j)}, \qquad (33)$$

where the $M \times N$ matrix $\mathbf{Z}_a$ denotes the post-processing noise after deploying arithmetic mean. Consider the following inequality

$$\sum_{m=1}^{M} |\mathbf{W}_v^{(i,m)}|^2 + \sum_{n=1}^{N} |\mathbf{W}_h^{(j,n)}|^2$$
$$\ge \mathbf{W}_v^{(i,i)^2} + \mathbf{W}_h^{(j,j)^2} \ge 2\mathbf{W}_v^{(i,i)}\mathbf{W}_h^{(j,j)}. \qquad (34)$$

With (33) and (34), the upper bound of $\mathrm{VAR}\big[\mathbf{Z}_a^{(i,j)}\big]$ can be represented as

$$\mathrm{VAR}\big[\mathbf{Z}_a^{(i,j)}\big] \le 2\big\{ \sum_{m=1}^{M} |\mathbf{W}_v^{(i,m)}|^2 + \sum_{n=1}^{N} |\mathbf{W}_h^{(j,n)}|^2 \big\} \frac{\sigma^2}{2}$$
$$\le 2\sigma^2 = N_0. \qquad (35)$$

With (33), we find the lower bound of $\mathrm{VAR}\big[\mathbf{Z}_a^{(i,j)}\big]$ as

$$\mathrm{VAR}\big[\mathbf{Z}_a^{(i,j)}\big] > \frac{\sigma^2}{2} \sum_{m=1}^{M} |\mathbf{W}_v^{(i,m)}|^2 + \frac{\sigma^2}{2} \sum_{n=1}^{N} |\mathbf{W}_h^{(j,n)}|^2$$
$$= \frac{\sigma^2}{2}\rho = \frac{\rho}{4}N_0. \qquad (36)$$

with $\rho = \sum_{m=1}^{M} |\mathbf{W}_v^{(i,m)}|^2 + \sum_{n=1}^{N} |\mathbf{W}_h^{(j,n)}|^2$ and $0 \le \rho \le 2$. Finally, let us make a summary, by jointly considering (35) and (36). It holds,

$$\frac{\rho}{4}N_0 < \mathrm{VAR}\big[\mathbf{Z}_a^{(i,j)}\big] \le N_0. \qquad (37)$$

(q.e.d.) ∎

*Further Discussion*

In Fig. 12, numerical demonstration is provided to exhibit the contour of $|\mathbf{W}_v^{(i,j)}|^2$ with $M = 8$. The presented weight matrix is the average of a large number of weight matrices, generated by the pre-trained NN-model for certain stationary statistics. Thus, the weight matrices in Fig. 12 can be regarded as the approximation of the optimum $|\mathbf{W}_{v,\mathrm{genie}}^{(i,j)}|^2$ within different batch stationary inference stages. In i.i.d. data batch, no spatial correlation can be exploited, and in the data batch involving Kronecker model, rich spatial correlation exists. Fig. 12 not only verifies the equations (25)-(27) numerically, but also reveals the fact, that the off-diagonal elements in $\mathbf{W}_v$ and $\mathbf{W}_h$ will become dominant, if the spatial correlation increases. This makes the term $\mathbf{W}_v^{(i,i)}\mathbf{W}_h^{(j,j)}$ in (33) reduce and become negligible, and helps to reduce the post-processing noise variance and improve the quality of channel estimates. Let us again focus on post-processing noise variance in (33), and introduce following notations for further simplification as

$$z = \frac{1}{\sigma^2}\mathrm{VAR}\big[\mathbf{Z}_a^{(i,j)}\big], \qquad (38)$$

$$x = \sum_{m=1}^{M} \big|\mathbf{W}_v^{(i,m)}\big|^2, \qquad (39)$$

$$y = \sum_{n=1}^{N} \big|\mathbf{W}_h^{(j,n)}\big|^2, \qquad (40)$$

$$d = \mathbf{W}_v^{(i,i)}\mathbf{W}_h^{(j,j)}. \qquad (41)$$

Obviously, equation (33) can be thus represented as

$$z = \frac{1}{2}(x + y) + d. \qquad (42)$$

Furthermore, considering the post-processing noise variance $\text{VAR}\big[\mathbf{Z}_v^{(i,j)}\big]$ and $\text{VAR}\big[\mathbf{Z}_h^{(i,j)}\big]$, which are uniquely obtained from vertical and horizontal *Subspace Training*, respectively. With equation (29) to equation (31), let us additionally introduce,
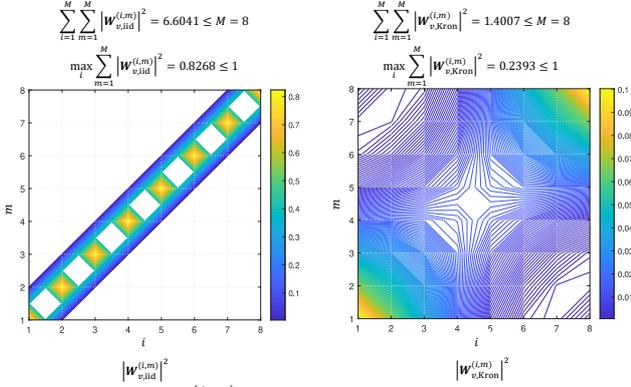


Fig. 12. Contour of $|\mathbf{W}_v^{(i,m)}|^2$ with and without spatial correlation at 10dB SNR

$$\frac{1}{\sigma^2}\text{VAR}\big[\mathbf{Z}_v^{(i,j)}\big] = 2\sum_{i=0}^{M}\big|\mathbf{W}_v^{(i,m)}\big|^2 = 2x, \qquad (43)$$

$$\frac{1}{\sigma^2}\text{VAR}\big[\mathbf{Z}_h^{(i,j)}\big] = 2\sum_{j=0}^{N}\big|\mathbf{W}_h^{(j,n)}\big|^2 = 2y. \qquad (44)$$

Thus, we are able to find the effective subspace combining region for $(x,y)$-pairs, which realizes both inequalities $\text{VAR}\big[\mathbf{Z}_a^{(i,j)}\big] < \text{VAR}\big[\mathbf{Z}_v^{(i,j)}\big]$ and $\text{VAR}\big[\mathbf{Z}_a^{(i,j)}\big] < \text{VAR}\big[\mathbf{Z}_h^{(i,j)}\big]$ for the post-processing noise variance in vertical and horizontal domains individually, so that the boosting effect in Fig. 8 by means of subspace combining can be achieved, if *Turbo-AI* is additionally deployed. With the introduced simplified notations, following equation system can be established as,

$$\begin{cases} z = \dfrac{1}{2}(x+y) + d, \\ z < 2x, \\ z < 2y, \end{cases} \qquad (45)$$

with $0 \le x \le 1$, $0 \le y \le 1$ and $0 \le d \le 1$. In Fig. 13, the equation system (45) is sketched in Cartesian coordinates. The blue intersection part denotes the $(x,y,z)$-pairs, that satisfy the equation (45), and its corresponding projection on plane $z = d$ is illustrated as the red shadow. With Fig. 13, we obtain following inequality

$$\frac{1}{3}\sum_{m=1}^{M}\big|\mathbf{W}_v^{(i,m)}\big|^2 + \frac{2}{3}\mathbf{W}_v^{(i,i)}\mathbf{W}_h^{(j,j)} < \sum_{n=1}^{N}\big|\mathbf{W}_h^{(j,n)}\big|^2$$
$$< 3\sum_{m=1}^{M}\big|\mathbf{W}_v^{(i,m)}\big|^2 - 2\mathbf{W}_v^{(i,i)}\mathbf{W}_h^{(j,j)}. \qquad (46)$$

Notice that the equation (46) provides the range for both $\sum_{m=1}^{M}|\mathbf{W}_v^{(i,m)}|^2$ and $\sum_{n=1}^{N}|\mathbf{W}_h^{(j,n)}|^2$, in order to effectively support *Turbo-AI*. First of all, for a 2D massive array system,

$\sum_{m=1}^{M}|\mathbf{W}_v^{(i,m)}|^2$ and $\sum_{n=1}^{N}|\mathbf{W}_h^{(j,n)}|^2$ should be adjusted by selecting the number of antenna elements for both vertical and horizontal subspaces, or by configuring the antenna spacing to realize spatial correlation. Secondly, notice that the red shadowing part in Fig. 13 will be reduced, if $d = \mathbf{W}_v^{(i,i)}\mathbf{W}_h^{(j,j)}$ increases. Once $d = 1$ holds, the shadowing part will disappear completely. This indicates that the subspace combining will not be effective any more, if the diagonal elements of matrices $\mathbf{W}_v$ and $\mathbf{W}_h$ become dominant, due to loss of spatial correlation.
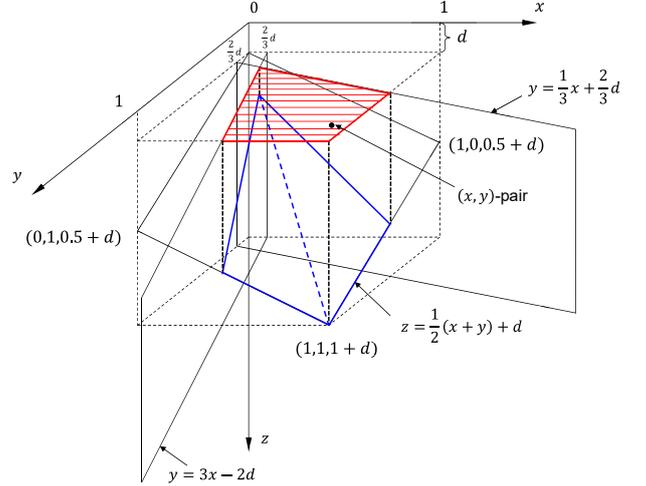


Fig. 13. Effective combining region for $(x,y)$-pairs to boost the performance with *Turbo-AI*

### B. Proof of Turbo-AI

*Turbo-AI* is an iterative procedure. In training stage, *Turbo-AI* exploits the new estimates to replace the observations and repeats the training. The applicability of *Turbo-AI* comes from the fact that the additive noise variance after each training will be smaller than that of the previous observation, which can be represented by following equality,

$$E\big[\|\mathbf{W}_{I+1}\big(\prod_{i=0}^{I}\mathbf{W}_i\mathbf{Y}\big) - \mathbf{H}\|_F^2\big] \le E\big[\|\prod_{i=0}^{I}\mathbf{W}_i\mathbf{Y} - \mathbf{H}\|_F^2\big], \qquad (47)$$

where $\mathbf{W}_i$ denotes the weighting matrix of $i$-th iteration, $\mathbf{Y}$ and $\mathbf{H}$ denotes the observations and labels, respectively. Furthermore, the noise distribution at the input and output of the NN remains Gaussian due to the fact that our NN in essence is a linear operator on the input data $\mathbf{Y}$.

*Proof:* In *Appendix A*, the convergence of a neural network is defined. Thus, after initializing *Turbo-AI* at 0-th iteration, if the convergence is reached after the learning stage, it holds,

$$E\big[\|\mathbf{W}_0\mathbf{Y} - \mathbf{H}\|_F^2\big] \le E\big[\|\underbrace{\mathbf{Y} - \mathbf{H}}_{\mathbf{Z}}\|_F^2\big]. \qquad (48)$$

And at first iteration, as long as the convergence is reached with the updated observations $\mathbf{W}_0\mathbf{Y}$, a weighting matrix $\mathbf{W}_1$

exists, which satisfies the following inequality,

$$E\big[\|\mathbf{W}_1\mathbf{W}_0\mathbf{Y} - \mathbf{H}\|_F^2\big] \le E\big[\|\underbrace{\mathbf{W}_0\mathbf{Y} - \mathbf{H}}_{\mathbf{Z}_0}\|_F^2\big]. \qquad (49)$$

In order to prove (37) in sense of mathematical induction, let us assume that the following inequality holds at $I$-th iteration,

$$E\big[\|\mathbf{W}_I\big(\prod_{i=0}^{I-1}\mathbf{W}_i\mathbf{Y}\big) - \mathbf{H}\|_F^2\big] \le E\big[\|\underbrace{\prod_{i=0}^{I-1}\mathbf{W}_i\mathbf{Y} - \mathbf{H}}_{\mathbf{Z}_{I-1}}\|_F^2\big].$$

$$(50)$$

Then, at $(I+1)$-th iteration, with the updated observations $\prod_{i=0}^{I}\mathbf{W}_i\mathbf{Y}$, the training can be carried out in the same neural network again. The weighting matrix $\mathbf{W}_{I+1}$ will be obtained, as the convergence is reached. It holds,

$$E\big[\|\mathbf{W}_{I+1}\big(\prod_{i=0}^{I}\mathbf{W}_i\mathbf{Y}\big) - \mathbf{H}\|_F^2\big] \le E\big[\|\underbrace{\prod_{i=0}^{I}\mathbf{W}_i\mathbf{Y} - \mathbf{H}}_{\mathbf{Z}_I}\|_F^2\big].$$

$$(51)$$

(q.e.d.) ∎

## REFERENCES

[1] Silver, D. *et. al.*; "Mastering the Game of Go with Deep Neural Networks and Tree Search," *Nature*, Vol. 529, pp. 484–503, Jan. 2016.

[2] Silver, D. *et. al.*; "Mastering the Game of Go without Human Knowledge," *Nature*, Vol. 550, pp. 354–371, Oct. 2017.

[3] Boccardi, F; Heath, R.W.; Lozano, A.; Marzetta, T.L.; Popovski, P.; "Five Disruptive Technology Directions for 5G," *IEEE Commun. Mag.*, Vol. 52, No. 2, pp. 74–80, Feb. 2014.

[4] Shannon, C.E.; "A Mathematical Theory of Communication," *Bell System Technical Journal 27 (3)*, pp. 379–423, Jul. 1948.

[5] O'Shea, T.; Hoydis, J. "An Introduction to Deep Learning for the Physical Layer," *IEEE Trans. Cognitive Commun. Network*, Vol. 3, No. 4, pp. 563–575, Dec. 2017.

[6] Neumann, D.; Wiese, T.; Utschick, W.; "Learning the MMSE channel estimator," *IEEE Trans. Signal Process.*, Vol. 66, No. 11, pp. 2905–2917, Jun. 2018.

[7] Chen, Yejian; ten Brink, S.; "Pilot Strategies for Trellis-Based MIMO Channel Tracking and Data Detection," in *Proc. 2013 IEEE Global Telecommun. Conf. (Globecom'13)*, pp. 4313–4318, Atlanta, USA, Dec. 2013.

[8] Dong, Peihao; Zhang, Hua; Li, Geoffrey Ye; Gaspar, I.; Naderi-Alizadeh, N.; "Deep CNN-Based Channel Estimation for mmWave Massive MIMO Systems," *IEEE J. Select. Topics Signal Process.*, Vol. 13, No. 5, pp. 989–1000, Sep. 2019.

[9] Marzetta, T.L.; "Noncooperative Cellular Wireless with Unlimited Numbers of Base Station Antennas," *IEEE Trans. Wireless Commun.*, Vol. 9, No. 11, pp. 3590–3600, Nov. 2010.

[10] Hellings, C.; Dehmani, A.; Wesemann, S.; Koller, M.; Utschick, W.; "Evaluation of Neural-Network-Based Channel Estimators Using Measurement Data," in *Proc. ITG Workshop on Smart Antennas (WSA'19)*, Apr. 2019.

[11] Koller, M.; Hellings, C.; Utschick, W.; "Learning-Based Channel Estimation for Various Antenna Array Configurations," in *Proc. IEEE 20th Int. Workshop Signal Process. Advances Wireless Commun. (SPAWC'19)*, Jul. 2019.

[12] Ying, Dawei; Vook, F.W.; Thomas, T.A.; Love, D.J.; Ghosh, A.; "Kronecker product correlation model and limited feedback codebook design in a 3D channel model," in *Proc. IEEE Int. Conf. Commun. (ICC)*, Jun. 2014.

[13] 3GPP Technical Specification Group Radio Access Network; "Spatial Channel Model for Multiple Input Multiple Output (MIMO) Simulations," *3GPP TS 25.996 v.15.0.0*, Jun. 2018.