

Generative and Explainable AI for High-Dimensional Channel Estimation

Nghia Thinh Nguyen and Tri Nhu Do

Department of Electrical Engineering, Polytechnique Montréal, Montréal, Québec, Canada.

Emails: {nghia-thinh.nguyen, tri-nhu.do}@polymtl.ca

Abstract—In this paper, we propose a new adversarial training framework to address high-dimensional instantaneous channel estimation in wireless communications. Specifically, we train a generative adversarial network to predict a channel realization in the time-frequency-space domain, in which the generator exploits the third-order moment of the input in its loss function and applies a new reparameterization method for latent distribution learning to minimize the Wasserstein distance between the true and estimated channel distributions. Next, we propose an explainable artificial intelligence mechanism to examine how the critic discriminates the generated channel. We demonstrate that our proposed framework is superior to existing methods in terms of minimizing estimation errors. Additionally, we find that the critic’s attention focuses on the high-power portion of the channel’s time-frequency representation.

Index Terms—5G NR, 3GPP, Channel estimation, MIMO, VAE, WGAN-GP, Third-order moment, Explainable AI

I. INTRODUCTION

Channel state information (CSI) acquisition, particularly channel estimation in 5G New Radio (NR), presents unprecedented challenges, especially in handling high-dimensional data [1]. The increased degrees of freedom (DoF) related to the dimension of transmission pilots across the time, frequency, and space domains lead to higher computational complexity with traditional methods [2]. In this context, generative artificial intelligence (AI) emerges as a potential solution to these pressing challenges [3]. Recently, researchers have focused on enhancing realistic 5G NR pilot-based channel estimation techniques through advanced simulation frameworks like Sionna [4], following standards such as 3GPP 38.901. However, the shift from conventional methods like Least Square (LS) and Linear Minimum Mean Square Error (LMMSE) to AI-driven approaches in high-dimensional scenarios introduces its own set of research challenges [5].

Generative adversarial network (GAN)-based channel estimation has been proposed as an alternative solution. [6] employs a GAN-based approach as the primary method to estimate the channel for high degrees of freedom in terms of channel components. [7] uses a GAN-based method combined with pilot information to estimate the channel. [8] utilized a Wasserstein GAN with gradient penalty (WGAN-GP) as a component to estimate the channel based on the LMMSE method. However, a fundamental issue with the GAN model is the imbalance between the critic and the generator, as the critic converges too quickly, leading the generator to produce a limited variety of samples or experience vanishing gradients. This results in the generator failing to learn effectively from

the critic’s feedback, manifesting as inaccurate or incomplete reconstruction of channel state information [9].

In this paper, we propose a novel framework for adversarial training aimed at instantaneous channel estimation. Specifically, we focus on enhancing the generator’s capability by leveraging additional probabilistic characteristics of the input features. To achieve this, we introduce an innovative loss function that integrates the third-order moment of the channel impulse response (CIR) with reconstruction loss, thereby capturing higher-order statistics of the channel characteristics. Additionally, we aim to strengthen the latent distribution learning capability of semi-supervised methods, such as Variational Autoencoder (VAE). Based on the proposed framework, we pose a *research question*: How can the AI model, particularly the critic in our framework, interpret the channel gain across the time, frequency, and space domains to effectively discriminate the generated channel?

The contributions of this work are twofold¹. (i) We propose the use of the third-order moment of high DoF input in the generator’s loss function and introduce a novel method for VAE’s latent distribution learning, resulting in superior performance, as demonstrated by lower normalized mean square error (NMSE) in channel estimation. Notably, our estimation framework does not require pilot knowledge and, in some cases, achieves better performance than traditional methods. (ii) We introduce a new explainable AI method using activation mapping for the critic, termed AM4C. We observe a significant alignment between the critic’s attention and the high-power gain regions in the time-frequency representation (TFR) of the CIR, providing insights into the design of the generator’s loss, pilot patterns, and other loss functions for future research.

II. SYSTEM MODEL

Notations. Underlined notations denote random variables, such as a single scalar $\underline{\alpha}$, a vector $\underline{\mathbf{x}}$, or a matrix $\underline{\mathbf{X}}$, while non-underlined notations represent realizations of these random variables, e.g., α , \mathbf{x} , and \mathbf{X} .

1) *System Components and Network Topology*: We consider a random network topology to analyze our proposed channel estimation framework comprehensively. Uplink (UL) estimation is performed by a single-antenna user terminal (UT) communicating with a base station (BS) equipped with K_{ant} antennas. The coordinates of the UT and BS in spherical

¹The source code developed for this paper is available at [tnd-lab.work](https://tnd-lab.github.io)

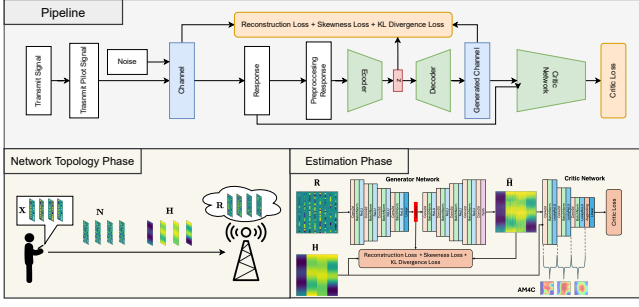


Fig. 1. Proposed WGAN-GP framework with explainable AI mechanism.

coordinates are denoted as $\mathbf{u}_{\text{ut,bs}} = [d_{\text{ut,bs}}, \varphi_{\text{ut,bs}}, \phi_{\text{ut,bs}}]$, where $d_{\text{ut,bs}}$ is the radial distance, $\varphi_{\text{ut,bs}}$ is the elevation angle of arrival along the z -axis, and $\phi_{\text{ut,bs}}$ is the azimuth angle in the horizontal x - y plane. The UT antenna follows an omnidirectional pattern, so the antenna response vector for the UT is represented as $a_{\text{ut}}(\varphi, \phi) = 1$. The BS is equipped with a uniform linear array (ULA) of K_{ant} cross-polarized antennas. The UL array response vector is: $\mathbf{a}_{\text{bs}}(\varphi, \phi) = [a_1(\varphi, \phi)e^{j\psi_1}, \dots, a_{K_{\text{ant}}}(\varphi, \phi)e^{j\psi_{K_{\text{ant}}}}]$, where $a_n(\varphi, \phi)$ represents the individual gain and phase response of the n -th antenna, λ_{cf} is the wavelength of the carrier frequency, and $\psi_n = -j\frac{2\pi}{\lambda_{cf}}(n-1)d_y \sin(\varphi) \sin(\phi)$ denotes the phase shift of the n -th antenna, d_y is the distance between two antennas.

2) *Signal Modeling*: Let K_{sym} be the number of OFDM symbols, i.e., time domain, K_{sc} be the number of effective subcarriers, i.e., frequency domain. Recall that K_{ant} is the number of antennas, representing the space domain. A pilot mask, $\mathbf{M} \in \{0, 1\}^{K_{\text{sc}} \times K_{\text{sym}}}$, is created using the orthogonal pilot pattern with Kronecker structure [4]. Let $\mathcal{T} \subseteq \{0, 1, \dots, N_{\text{sym}} - 1\}$ be the set of OFDM symbol indices reserved for pilot transmission, $\mathcal{K}_{i,j} \subseteq \{0, 1, \dots, N_{\text{sc}} - 1\}$ be the set of subcarrier indices allocated to transmitter i and stream j for pilot transmission. Considering *sparse* pilot pattern, the mask exhibits sparseness into distinct time and frequency components, i.e., $\mathbf{M} = \mathbf{P}_t \otimes \mathbf{P}_f^{i,j}$, where $\mathbf{P}_t(s)$ is a binary vector indicating pilot OFDM symbol indices, $\mathbf{P}_f^{i,j}(k)$ is a binary vector indicating the subcarrier indices allocated to transmitter i and stream j [4]. Here, $\mathbf{X}_p \in \mathbb{C}^{K_{\text{sc}} \times K_{\text{sym}}}$ can be represented as $\mathbf{X}_p^{i,j}(k, s) = p_{i,j}$ if $\mathbf{M}^{i,j}(k, s) = 1$, and 0 otherwise, where $p_{i,j}$ are the QAM pilot symbols for transmitter i , stream j . \mathbf{M} is designed as $\mathbf{M}^{i,j}(k, s) = 1$ if $s \in \mathcal{T}, k \equiv iK_{\text{str}} + j \pmod{K_{\text{seq}}}$, and 0 otherwise. Combining the mask and symbol assignment, the complete pilot pattern $\mathbf{X} \in \mathbb{C}^{K_{\text{sc}} \times K_{\text{sym}}}$ within the resource grid (RG) can be expressed as $\mathbf{X}^{i,j}(k, s) = \mathbf{M}^{i,j}(k, s) \times \mathbf{X}_p^{i,j}(k, s)$.

3) *Channel Modeling*: We consider communication between the UT and BS using the Urban Micro (UMi) channel model [4]. Let $h_{n,s,k}$ denote the channel for the n -th antenna and the k -th subcarrier at time step s . With the applied network topology, it can be represented as

$$h_{n,k,s} = \sum_{\ell=0}^{K_{\text{path}}-1} h_{\ell}(s) e^{-j2\pi k \Delta_f \tau_{\ell}} a_n(\varphi_{\ell}, \phi_{\ell}) e^{j\psi_{n,\ell}}, \quad (1)$$

where K_{path} denotes the number of multi-path components, and Δ_f represents the subcarrier spacing. Here, $h_{\ell}(s) \sim \mathcal{CN}(0, 1)$ denotes the complex amplitude, τ_{ℓ} the delay, and $\psi_{n,\ell}$ the phase shift, all provided by the UMi model [4]. According to (1), the instantaneous TFR of the channel $\mathbf{H} \in \mathbb{C}^{K_{\text{ant}} \times K_{\text{sym}} \times K_{\text{sc}}}$ represents a single realization of \mathbf{H} .

4) *Received Observation Modeling*: Based on the channel model and signal modeling, assuming that the cyclic prefix (CP) is added to the transmit waveform and removed from the received waveform in the time domain, the TFR of the received signal, \mathbf{R} , at the BS for the UT is expressed as

$$\mathbf{R} = \mathbf{H} \odot \mathbf{X}_{\text{exd}} + \mathbf{N}, \quad (2)$$

where $\mathbf{R} \in \mathbb{C}^{K_{\text{ant}} \times K_{\text{sym}} \times K_{\text{sc}}}$ represents the received signal, and $\mathbf{N} \in \mathbb{C}^{K_{\text{ant}} \times K_{\text{sym}} \times K_{\text{sc}}}$ is an i.i.d. matrix of additive white Gaussian noise (AWGN) with entries $[\mathbf{N}]_{ij} \sim \mathcal{CN}(0, 1)$. The expanded version of \mathbf{X} is denoted by $\mathbf{X}_{\text{exd}} = \text{repmat}(\mathbf{X}, K_{\text{ant}}, 1, 1)$, where repmat replicates \mathbf{X} along the antenna dimension.

III. PROPOSED ADVERSARIAL TRAINING FRAMEWORK

A. Adversarial Training for Channel Estimation

As aforementioned, we address the channel estimation problem using an adversarial training framework consisting of a generator G that generates synthetic channel estimates and a critic D that discriminates between the generated channel and the true channel. The generator G receives the observation \mathbf{R} as input, mapping it as $G(\mathbf{R})$ such that $G(\mathbf{R}) : \mathbb{R}^{K_{\text{ant}} \times K_{\text{sym}} \times K_{\text{sc}}} \rightarrow \mathbb{R}^{K_{\text{ant}} \times K_{\text{sym}} \times K_{\text{sc}}}$, where $\mathbf{R} \mapsto \hat{\mathbf{H}}$. Thus, the instantaneous channel estimate is $\hat{\mathbf{H}} = G(\mathbf{R})$. Given the true random instantaneous channel \mathbf{H} , the input of the critic can be either \mathbf{H} or $\hat{\mathbf{H}}$. The critic can be represented as $D(\mathbf{H}^*) : \mathbb{R}^{K_{\text{ant}} \times K_{\text{sym}} \times K_{\text{sc}}} \rightarrow [0, 1]$ where $\mathbf{H}^* \mapsto \{p(\mathbf{H}^*|H_0), p(\mathbf{H}^*|H_1)\}$. Here, \mathbf{H}^* can represent either \mathbf{H} or $\hat{\mathbf{H}}$ channel model. The likelihood probability $p(\mathbf{H}^*|H_0)$ and $p(\mathbf{H}^*|H_1)$ corresponding to the hypothesis H_0 for the generated channel and H_1 for the true channel.

The overall objective of the GAN framework is a *two-player minimax game*. The critic tries to maximize the value function $V(D, G)$ to distinguish between true and generated channels, while the generator tries to minimize this same value function to deceive the critic. Let $p_{\mathbf{H}}(\mathbf{H})$, $p_{\hat{\mathbf{H}}}(\hat{\mathbf{H}})$, $p_{\mathbf{R}}(\mathbf{R})$ denote the distribution of the true channel, the generated channel, and the received observation, respectively. The minimax game of our channel estimation problem is expressed as

$$\begin{aligned} \min_G \max_D V(D, G) &= \\ \mathbb{E}_{\mathbf{H} \sim p_{\mathbf{H}}(\mathbf{H})} [\log D(\mathbf{H})] + \mathbb{E}_{\mathbf{R} \sim p_{\mathbf{R}}(\mathbf{R})} [\log(1 - D(G(\mathbf{R})))] \\ &= \mathbb{E}_{\mathbf{H} \sim p_{\mathbf{H}}(\mathbf{H})} [\log D(\mathbf{H})] + \mathbb{E}_{\hat{\mathbf{H}} \sim p_{\hat{\mathbf{H}}}(\hat{\mathbf{H}})} [\log(1 - D(\hat{\mathbf{H}}))]. \end{aligned} \quad (3)$$

The technical challenge is the imbalanced optimization between the critic and generator networks during the early training phase. The critic D output has a more direct objective, discriminating tasks constrained to $[0, 1]$, while the generator G learns to produce channel distributions over a wider and

more complex domain, i.e., $[\min \mathbf{H}, \max \mathbf{H}]$. In addition, \mathbf{R} and \mathbf{H} are high-dimensional data, which makes it challenging for G to learn the complex underlying distribution.

B. Problem Formulation for Optimizing the Generator

The training objective in our paper is to enable the generator to converge quickly from one distribution to another, i.e., from $p_{\hat{\mathbf{H}}}(\hat{\mathbf{H}})$ to $p_{\mathbf{H}}(\mathbf{H})$. To achieve this, the training process aims to minimize their Wasserstein distance, i.e., $W(p_{\mathbf{H}}(\mathbf{H}), p_{\hat{\mathbf{H}}}(\hat{\mathbf{H}}))$. Since the set of all joint distributions $\prod(p_{\mathbf{H}}(\mathbf{H}), p_{\hat{\mathbf{H}}}(\hat{\mathbf{H}}))$ is not available, Kantorovich-Rubinstein's dual formulation is used to simplify the Wasserstein distance [10].

$$W(p_{\mathbf{H}}(\mathbf{H}), p_{\hat{\mathbf{H}}}(\hat{\mathbf{H}})) = \sup_{\|f\|_c \leq 1} \{ \mathbb{E}_{\mathbf{H} \sim p_{\mathbf{H}}(\mathbf{H})} [f(\mathbf{H})] - \mathbb{E}_{\hat{\mathbf{H}} \sim p_{\hat{\mathbf{H}}}(\hat{\mathbf{H}})} [f(\hat{\mathbf{H}})] \}, \quad (4)$$

where $f(\cdot)$ is a function parameterized by a neural network. Even if the probability density functions $p_{\hat{\mathbf{H}}}(\hat{\mathbf{H}})$ and $p_{\mathbf{H}}(\mathbf{H})$ do not perfectly align, the Wasserstein distance between them can remain small as long as their supports are close and the mass (probability density) is similarly distributed within those supports. In other words, if the generated channel estimates $\hat{\mathbf{H}}$ have a distribution whose support largely overlaps with that of the true channel \mathbf{H} and the mass is distributed in a similar fashion, then the Wasserstein distance will be small. Therefore, our goal is to ensure that the distribution $p_{\hat{\mathbf{H}}}(\hat{\mathbf{H}})$ of the estimated channels closely matches the distribution $p_{\mathbf{H}}(\mathbf{H})$ of the true channels, both in terms of support and mass distribution. To this end, first we let $\hat{\mathbf{H}}$ and \mathbf{H} be normalized to the same scale, i.e., the same range, such as \mathbf{H} and $\hat{\mathbf{H}}$ being in $[-1, 1]$. Considering pilot sparsity and adversarial training, the proposed channel estimation problem is formulated as

$$G^* = \arg \min_{(G|D)} W(p_{\mathbf{H}}(\mathbf{H}), p_{\hat{\mathbf{H}}}(\hat{\mathbf{H}})) \quad (5a)$$

$$\text{subject to } \Re\{[\mathbf{H}]_{ij}\}, \Im\{[\mathbf{H}]_{ij}\}, \Re\{[\hat{\mathbf{H}}]_{ij}\}, \Im\{[\hat{\mathbf{H}}]_{ij}\} \in [-1, 1], \forall i, \forall j, \quad (5b)$$

$$|p_{i,j}|^2 \leq 1, i \leq |\mathcal{T}|, j \leq K_{sc}, \quad (5c)$$

where $(G|D)$ indicates that we focus on optimizing G given a specific D . Constraint (5b) represents data normalization, while constraint (5c) enforces sparsity in the pilot pattern. We set $p_{i,j} = \pm \frac{1}{\sqrt{2}} \pm j \frac{1}{\sqrt{2}}$, ensuring that $|p_{i,j}|^2 = 1$. Note that while the realization values are scaled to $[-1, 1]$, the distribution of the true channel $p_{\mathbf{H}}(\mathbf{H})$ and the generated channel $p_{\hat{\mathbf{H}}}(\hat{\mathbf{H}})$ differ. Achieving alignment between these distributions remains challenging.

IV. THIRD-ORDER MOMENT-BASED VAE-WGAN-GP

A. Mathematical Description of the Proposed G and D

1) *NN Architecture*: Our goal is to upgrade the generator to a more refined version. Based on the work of [10], further investigated in [11], we define the generator function f_g with the input being the observation data \mathbf{R} , and the set of parameters θ is updated accordingly. Note that for a composite function,

$f \circ g = f(g(x))$, and $\mathbf{H}^* = f_g(\mathbf{R}; \theta)$. Let $f_g^{[bl_i]}$ denote the i -th feature extraction block, which can be represented as $f_g^{[bl_i]} = f_g^{[a_i]} \circ f_g^{[b_i]} \circ f_g^{[c_i]}$, comprising activation, batch normalization, and convolution layers. Drawing from the VAE architecture [12], the G function can be divided into the encoder, latent space, and decoder parts. Let $\mathbf{a}_{g,enc} = f_{g,enc}(\mathbf{R}; \theta_{\mathbf{a}_{g,enc}}) = f_{g,enc}^{[bl_1]} \circ f_{g,enc}^{[bl_2]} \circ f_{g,enc}^{[bl_3]} \circ f_{g,enc}^{[bl_4]} \circ f_{g,enc}^{[l]}$ represent the output vector of the encoder function. The latent vector \mathbf{z} is the output of the transformation function $\mathbf{z} = f_{g,trans}(\mathbf{a}_{g,enc})$ and serves as the input to the decoder function. The decoder output is given by $\hat{\mathbf{H}} = f_{g,dec}(\mathbf{z}; \theta_{g,dec}) = f_{g,dec}^{[bl_1]} \circ f_{g,dec}^{[bl_2]} \circ f_{g,dec}^{[bl_3]} \circ f_{g,dec}^{[bl_4]} \circ f_{g,dec}^{[tanh]}$. The function $f_{g,dec}^{[tanh]}$ ensures that the output is normalized within $[-1, 1]$, aligning with the scaling of the true channel \mathbf{H} . Thus, G is mathematically described as

$$f_g(\mathbf{R}; \theta) = f_{g,dec}(f_{g,trans}(f_{g,enc}(\mathbf{R}; \theta_{\mathbf{a}_{g,enc}})); \theta_{g,dec}). \quad (6)$$

In contrast, the critic network, denoted as f_c , has a set of parameters ω and includes a feature extraction block defined as $f_{c,ft}^{[bl_i]} = f_c^{[a_i]} \circ f_c^{[b_i]} \circ f_c^{[c_i]}$ at the i -th block. Thus, the critic D can be mathematically expressed as

$$f_c(\mathbf{H}^*; \omega) = f_c^{[l_4]}(f_c^{[avg_4]}(f_{c,ft}^{[bl_3]}(f_{c,ft}^{[bl_2]}(f_{c,ft}^{[bl_1]}(\mathbf{H}^*; \omega_{c,ft}^{[bl_1]}); \omega_{c,ft}^{[bl_2]}); \omega_{c,ft}^{[bl_3]}); \omega_{c,ft}^{[avg_4]}), \omega_{c,ft}^{[l_4]}). \quad (7)$$

Before training the critic network, the function $f_c^{[dim]}(\mathbf{H}) : \mathbb{R}^{K_{ant} \times K_{sym} \times K_{sc}} \mapsto \mathbb{R}^{(K_{ant} \times K_{sym}) \times K_{sc}}$ is defined to reduce the dimensionality of the channel input \mathbf{H} , making the model easier to train and reducing computational costs for high-dimensional channels. As the critic's target, the output of the linear function $f_c^{l_4}$ is scalar, with $f_c^{l_4}(X_{\mathbf{H}}^{[l_4]}, \omega_{c,ft}^{[l_4]}) > f_c^{l_4}(X_{\hat{\mathbf{H}}}^{[l_4]}, \omega_{c,ft}^{[l_4]})$, where $X_{\mathbf{H}}^{[l_4]}$ and $X_{\hat{\mathbf{H}}}^{[l_4]}$ are the inputs to the linear function for $D(\mathbf{H})$ and $D(\hat{\mathbf{H}})$, respectively.

2) *Proposed Latent Distribution Learning in VAE's Probabilistic Framework*: $q(\mathbf{z}|\mathbf{H})$ represent the true posterior probability given the real channel \mathbf{H} . Since \mathbf{H} cannot be observed directly, $q(\mathbf{z}|\mathbf{H})$ cannot be computed directly. We approximate it by sampling the approximate posterior probability of the encoder output vector, where $q(\mathbf{z}|\mathbf{H}) \approx p(\mathbf{z}|\mathbf{a}_{g,enc})$. Using the reparameterization method, the latent vector \mathbf{z} is defined as $\mathbf{z} = f_{g,trans}(\mathbf{a}_{g,enc}) = \mu_{\mathbf{a}_{g,enc}} + \epsilon \sigma_{\mathbf{a}_{g,enc}}$, where $\mathbf{z} \sim p(\mathbf{z}|\mathbf{a}_{g,enc}) = \mathcal{N}(\mu_{\mathbf{z},\theta}, \sigma_{\mathbf{z},\theta})$. In most cases, the estimated latent vector $\mathbf{z} \sim \mathcal{N}(0, 1)$ is used for the inference phase. However, in our scenario, $\hat{\mathbf{H}} \sim q(\hat{\mathbf{H}}|\mathbf{z})$ cannot be generated randomly by the decoder using a Gaussian random vector as input. Moreover, $\mathbf{a}_{g,enc}$ is a compressed vector of the observation \mathbf{R} rather than \mathbf{H} , making $q(\mathbf{z}|\mathbf{H}) \approx p(\mathbf{z}|\mathbf{a}_{g,enc})$ unreasonable, as $p_{\mathbf{R}}(\mathbf{R})$ and $p_{\mathbf{H}}(\mathbf{H})$ differ, as evidenced by the results. The proposed objective is to optimize the latent vector $\mathbf{z} \sim p(\mathbf{z}|\mathbf{a}_{g,enc})$ to closely approximate $p(\mathbf{z}|\mathbf{H})$, with $p_{\mathbf{H}}(\mathbf{H}) \sim \mathcal{N}(\mu_{\mathbf{H}}, \sigma_{\mathbf{H}}^2)$. It is assumed that the number of training sets of real channel $|\mathcal{D}_{train}|$ is large enough to compute $\mu_{\mathbf{H}}$, and $\sigma_{\mathbf{H}}$. The latent vector $\mathbf{z} \sim \mathcal{N}(\mu_{\mathbf{z},\theta}, \sigma_{\mathbf{z},\theta}^2)$ is estimated directly from compressed vector $\mathbf{a}_{g,enc}$ and it is transformed close to $\mathcal{N}(\mu_{\mathbf{H}}, \sigma_{\mathbf{H}}^2)$. Applying trick parameterization, $\mathbf{z} = f_{g,trans}(\mathbf{a}_{g,enc}) = \mu_{\mathbf{H}} + \mathbf{a}_{g,enc} \sigma_{\mathbf{H}}$.

B. Proposed Loss Function

The target critic function D tries to maximize the distance of two distributions, while the generator model G is the opposite.

1) *Loss of D* : The loss function of critic network D is minimized by updating ω . Its formulation is represented as

$$L_D = \underbrace{\mathbb{E}_{\hat{\mathbf{H}} \sim p_{\hat{\mathbf{H}}}(\hat{\mathbf{H}})}[D\omega(\hat{\mathbf{H}})]}_{\triangleq L_f: \text{generative loss}} - \underbrace{\mathbb{E}_{\mathbf{H} \sim p_{\mathbf{H}}(\mathbf{H})}[D\omega(\mathbf{H})]}_{\triangleq L_r: \text{true loss}}. \quad (8)$$

Consider $p_{\hat{\mathbf{H}}}(\hat{\mathbf{H}})$ as the mixture distribution of the $p_{\mathbf{H}}(\mathbf{H})$ and $p_{\hat{\mathbf{H}}}(\hat{\mathbf{H}})$, is given by: $p_{\hat{\mathbf{H}}}(\hat{\mathbf{H}}) = \epsilon p_{\mathbf{H}}(\mathbf{H}) + (1 - \epsilon)p_{\hat{\mathbf{H}}}(\hat{\mathbf{H}})$ where $\epsilon \sim \mathcal{U}(0, 1)$ is a random scalar, drawn from a uniform distribution between 0 and 1. To constrain the critic network's update speed and preserve the 1-Lipschitz property. The gradient penalty loss can be defined as

$$L_{gp} = \lambda_{gp} \mathbb{E}_{\hat{\mathbf{H}} \sim p_{\hat{\mathbf{H}}}(\hat{\mathbf{H}})}[(\|\nabla D\omega(\hat{\mathbf{H}})\|_2 - 1)^2]. \quad (9)$$

From (8) and (9), the adopted L_D is expressed as

$$L_D = \text{generative loss} - \text{true loss} + \text{gradient penalty} \\ = L_f - L_r + L_{gp}. \quad (10)$$

2) *Proposed Loss of G* : In the minimax game framework, the discriminator D function seeks to maximize the Lipschitz continuity of its output by controlling the gradients of $\hat{\mathbf{H}}$ through its parameters ω . Concurrently, the generator G updates its parameters θ to minimize the discrepancy between the generated data and real data. Initially, the loss of G can be designed based on the negative fake loss [10] as

$$L_G = -L_f = -\mathbb{E}_{\mathbf{R} \sim p_{\mathbf{R}}(\mathbf{R})}[D(G\theta(\mathbf{R}))]. \quad (11)$$

Taking into account the evidence lower bound of the VAE generator, we consider the KL divergence loss and the reconstruction loss, respectively, as

$$L_{KL} = \lambda_{KL} \frac{1}{2} \sum_{d=1}^{z_{\text{dim}}} \left(\log \frac{\sigma_{H_d}^2}{\sigma_{z_d, \theta}^2} + \frac{\sigma_{z_d, \theta}^2 + (\mu_{z_d, \theta} - \mu_{H_d})^2}{\sigma_{H_d}^2} - 1 \right), \quad (12)$$

$$L_{\text{rec}} = \mathbb{E}_{\mathbf{H} \sim p_{\mathbf{H}}(\mathbf{H}), \mathbf{R} \sim p_{\mathbf{R}}(\mathbf{R})}[\|\mathbf{H} - G\theta(\mathbf{R})\|^2]. \quad (13)$$

Next, we aim to make the two distributions $p_{\mathbf{H}}(\mathbf{H})$ and $p_{\hat{\mathbf{H}}}(\hat{\mathbf{H}})$ as close as possible. Let γ_r and γ_g represent the *coefficients of skewness* or the third-order moments of the true and estimated channel distributions, respectively, which can be expressed as

$$\gamma_r = (\mathbb{E}_{\mathbf{H} \sim p_{\mathbf{H}}(\mathbf{H})}[(\mathbf{H} - \mu_r)^3]) / \sigma_r^3, \quad (14)$$

$$\gamma_g = (\mathbb{E}_{\hat{\mathbf{H}} \sim p_{\hat{\mathbf{H}}}(\hat{\mathbf{H}})}[(\hat{\mathbf{H}} - \mu_g)^3]) / \sigma_g^3, \quad (15)$$

where μ_r , μ_g , σ_r , and σ_g are the mean and the standard deviation of the true and generative distribution respectively. It is assumed that the distance of the skewness of true and generative distributions is close and continuous if the distance of two distributions is close. Thus, from (14) and (15), we can define our third order moment-based loss as

$$L_\gamma = |\gamma_r - \gamma_g|. \quad (16)$$

From (11)-(13), and (16), the proposed L_G is expressed as

$$L_G = -\text{fake loss} + \text{rec loss} + \text{kl loss} + \text{skewness loss} \\ = -L_f + L_{\text{rec}} + L_{KL} + L_\gamma. \quad (17)$$

The proposed VAE-WGAN-GP-based channel estimation method, with enhanced G loss and latent distribution learning, is outlined in the pseudo-code of Algorithm 1.

Algorithm 1: Training Procedure of the Proposed VAE-WGAN-GP Generator and Critic Models

Input: Noisy received signal $\mathbf{R} \sim p_{\mathbf{R}}(\mathbf{R})$. The true channel data $\mathbf{H} \sim p_{\mathbf{H}}(\mathbf{H})$

Output: Trained Generator G for channel estimation $\hat{\mathbf{H}} \sim p_{\hat{\mathbf{H}}}(\hat{\mathbf{H}})$

```

1 Compute  $\mu_{\mathbf{H}}, \sigma_{\mathbf{H}}$  from training set  $\mathcal{D}_{\text{train}}$ ;
2 for  $e = 0, \dots, K_{\text{epoch}}$  do
3   Sample  $\{\mathbf{H}_i^{[e]}\}_{i=1}^b \sim p_{\mathbf{H}}(\mathbf{H})$  a batch from true channel;
4   Sample  $\{\mathbf{R}_i^{[e]}\}_{i=1}^b \sim p_{\mathbf{R}}(\mathbf{R})$  a batch from received signal;
5   for  $t = 0, \dots, n_c$  do
6      $\{\hat{\mathbf{H}}_i^{[t],[e]}\}_{i=1}^b \leftarrow G(\{\mathbf{R}_i^{[e]}\}_{i=1}^b) \sim p_{\hat{\mathbf{H}}}(\hat{\mathbf{H}})$ 
       predict a batch from generated channel;
7      $\{\tilde{\mathbf{H}}_i^{[t],[e]}\}_{i=1}^b \leftarrow \epsilon \{\mathbf{H}_i^{[e]}\}_{i=1}^b + (1 - \epsilon) \{\hat{\mathbf{H}}_i^{[t],[e]}\}_{i=1}^b$  a batch of mixture channel;
8     Training  $a_c^{[t],[e]} \leftarrow D(\{\mathbf{H}_i^{*[t],[e]}\}_{i=1}^b)$ ;
9     Compute  $f_{L_D}(a_c^{[t],[e]})$ ;
10    Compute  $f_{L_{gp}}(\{\tilde{\mathbf{H}}_i^{[t],[e]}\}_{i=1}^b)$ ;
11     $f_{L_D}^{[t],[e]} \leftarrow f_{L_D}(a_c^{[t],[e]}) + f_{L_{gp}}(\{\tilde{\mathbf{H}}_i^{[t],[e]}\}_{i=1}^b)$ ;
12     $\omega^{[t+1],[e]} \leftarrow \text{Opt}(f_{L_D}^{[t],[e]}, \omega^{[t],[e]}, \eta, \beta_1, \beta_2)$ ;
13  end
14   $\mathbf{a}_{g,\text{enc}}^{[e]} \leftarrow f_{g,\text{enc}}(\{\mathbf{R}_i^{[e]}\}_{i=1}^b) \sim p(\mathbf{a}_{g,\text{enc}})$ ;
15  Reparameterize
     $\mathbf{z}^{[e]} \leftarrow \mu_{\mathbf{H}} + \mathbf{a}_{g,\text{enc}}^{[e]} \sigma_{\mathbf{H}} \sim p(\mathbf{z} | \mathbf{a}_{g,\text{enc}})$ ;
16   $\{\hat{\mathbf{H}}_i^{[e]}\}_{i=1}^b \leftarrow f_{g,\text{dec}}(\mathbf{z}^{[e]})$ ;
17  Compute  $f_{L_f}(\{\hat{\mathbf{H}}_i^{[e]}\}_{i=1}^b)$ ;
18  Compute  $f_{L_{\text{rec}}}(\{\hat{\mathbf{H}}_i^{[e]}\}_{i=1}^b, \{\mathbf{H}_i^{[e]}\}_{i=1}^b)$ ;
19  Compute  $f_{L_{KL}}(\mathbf{z}^{[e]}, \mu_{\mathbf{H}}, \sigma_{\mathbf{H}})$ ;
20   $f_{L_G}^{[e]} \leftarrow f_{L_f}(\{\hat{\mathbf{H}}_i^{[e]}\}_{i=1}^b) + f_{L_{\text{rec}}}(\{\hat{\mathbf{H}}_i^{[e]}\}_{i=1}^b, \{\mathbf{H}_i^{[e]}\}_{i=1}^b) + f_{L_{KL}}(\mathbf{z}^{[e]}, \mu_{\mathbf{H}}, \sigma_{\mathbf{H}})$ ;
21   $\theta^{[e+1]} \leftarrow \text{Opt}(f_{L_G}^{[e]}, \theta^{[e]}, \eta, \beta_1, \beta_2)$ ;
22 end
23 Initialization: Initialize parameters of  $G$  and  $D$ , set learning rates  $\eta$ , the batch size  $b$ , the number of iterations of the critic per generator iteration  $n_c$ , other Optimizer's parameters  $\beta_1$  and  $\beta_2$ 

```

Remark 1. Numerical results demonstrate that all components of the Generator loss L_G , including the fake loss L_f , reconstruction loss L_{rec} , KL divergence loss L_{KL} , and the third-order moment-based skewness loss L_γ , successfully converge

during the training process. This collective convergence of the individual loss components ensures the overall convergence of the proposed Generator loss L_G , as illustrated in Fig. 2.

C. Inference and Performance Analysis

To infer the model, our strategy involves simulating $|\mathcal{D}_{\text{test}}|$ samples from the Sionna framework, with each sample containing instantaneous channels with added noise as test sets. To adapt the model architecture, we reshape the actual channel based on $f_c^{[\text{dim}]}(\mathbf{H}^*)$. The evaluation metric is the Normalized Mean Square Error (NMSE), which measures the difference between the true and generated channels, defined as

$$\text{NMSE} = \frac{\sum_{i=1}^{|\mathcal{D}_{\text{test}}|} (\mathbf{H}_i - \hat{\mathbf{H}}_i)^2}{\sum_{i=1}^{|\mathcal{D}_{\text{test}}|} (\mathbf{H}_i)^2}. \quad (18)$$

D. Explainable AI: Activation Mapping for the Critic (AM4C)

We propose an Explainable AI method to study the interpretability of the critic D during the training loop. Specifically, utilizing the activation mapping method, we propose the AM4C algorithm to identify and visualize the regions in the TRF of the channel that significantly influence the decision of D . The activation map of the k feature map at the spatial location (i, j) from the convolution function $f_c^{[cn]}$, where $n \in \{1, 2, 3\}$ is expressed as

$$f_c^{[cn]}(\mathbf{X}_c^{[cn]}; \omega_c^{[cn]}) = \{\mathbf{A}_{k,i,j}^{cn} | k \in \{1, \dots, K\}\} \quad (19)$$

where K is the number of spatial locations, $\mathbf{X}_c^{[cn]}$ is input of convolution layer. Next, we compute the gradient of the critic output of $f_c(\mathbf{H}^*; \omega$ with respect to each feature map's activation $\frac{\partial f_c(\mathbf{H}^*; \omega)}{\partial f_c^{[cn]}(\mathbf{X}_c^{[cn]}; \omega_c^{[cn]})}$. Let a_c^k be the weight representing the overall influence of the critic decision value as follows

$$a_c^k = \frac{1}{K} \sum_{i,j} \frac{\partial f_c(\mathbf{H}^*; \omega)}{\partial f_c^{[cn]}(\mathbf{X}_c^{[cn]}; \omega_c^{[cn]})}. \quad (20)$$

Consider ReLU activation function to remove negative values. The activation map $M_D^{cn}(i, j)$ can be defines as

$$M_D^{cn}(i, j) = \text{ReLU}\left(\frac{1}{K} \sum_k \sum_{i,j} \frac{\partial f_c(\mathbf{H}^*; \omega)}{\partial f_c^{[cn]}(\mathbf{X}_c^{[cn]}; \omega_c^{[cn]})} \times f_c^{[cn]}(\mathbf{X}_c^{[cn]}; \omega_c^{[cn]})\right). \quad (21)$$

The resulting activation map $M_D^{cn}(i, j)$ provides a visual and quantitative representation of which spatial regions in the input channel data are most influential in the critic D . High values in the activation map indicate areas with substantial impact, thereby serving as the *attention regions* that the model focuses on during channel estimation.

Remark 2. After training, the attention areas identified by the critic D consistently align with regions of high magnitude in the channel's TRF, as demonstrated in Fig. 6.

V. NUMERICAL RESULTS AND VISUALIZATION

Key simulation settings are presented in Table I; additional settings for simulation, training, and inference can be found in our published source code.

TABLE I
NETWORK, SIMULATION, AND HYPER-PARAMETERS SETTINGS

Parameter	Value	Parameter	Value
K_{ant}	4	λ_{gp}	10
K_{sym}	8	λ_{KL}	0.1
K_{sc}	32	\mathcal{T}	$\{2, 6\}$

1) *Training Evaluation:* To demonstrate the advantage of our proposed method, Fig. 2 shows that the generator loss and the critic loss tend to decrease shapely after K_{epoch} . It means that the proposed loss function convergence.

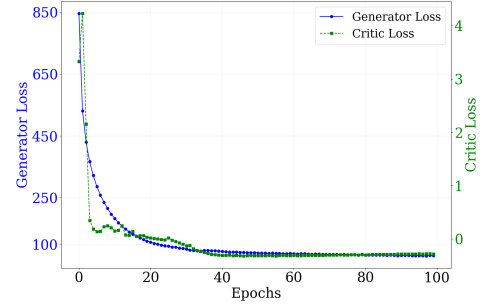


Fig. 2. Loss function of G and D

2) *Inference:* For testing, we use the same topology as the training set. To ensure fidelity, we also generate instantaneous channels with random noise as a test set and then evaluate the dataset. The comparison is shown in Fig. 3; the magnitude of the channel on each antenna is reconstructed similarly to the true channel. The results also indicate that the normalization process impacts the reconstruction of channel characteristics.

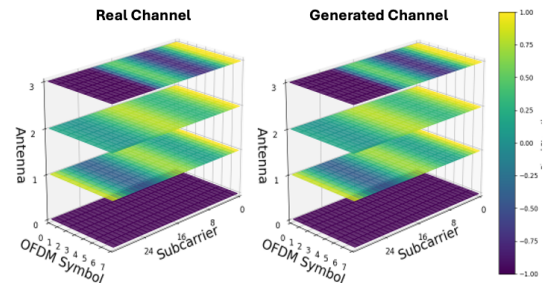


Fig. 3. True 3D channel realization and corresponding generative channel.

3) *Performance Evaluation in terms of NMSE:* We compare the NMSE of the proposed method with conventional methods using LS with interpolation techniques, i.e., nearest neighbor, linear, and LMMSE across different data domains. Fig. 4 shows that our methods obtain comparable outcomes with LMMSE methods in the lower SNRs. Besides it is a note that learning to generate samples without explicit distribution of WGAN-GP gives a poor result if compared with VAE methods. However, combining VAE and WGAN-GP can be a better solution, even with the proposed loss function which is the best result in our experiment.

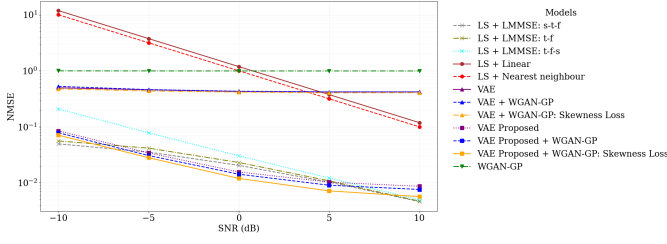


Fig. 4. NMSE of the proposed VAE-WGAN-GP and conventional methods.

It is true that the encoder input \mathbf{R} estimate from the function of Gaussian noise which is slightly shifted by the operation of \mathbf{X} containing mostly zero. It means that $\mathbf{R} \sim \mathcal{N}(\mu_{\mathbf{R}}, \sigma_{\mathbf{R}})$, where $\mu_{\mathbf{R}} \approx 0$, and $\sigma_{\mathbf{R}} \approx 1$ as can be seen from the first image of Fig. 5. The distribution of the encoder output $\mathbf{a}_{g,enc}$ is estimated as the same shape as the decoder input \mathbf{z} and they are a right-skewed distribution. As a consequence, The distribution of the generated $\hat{\mathbf{H}}$ and the true \mathbf{H} is similar.

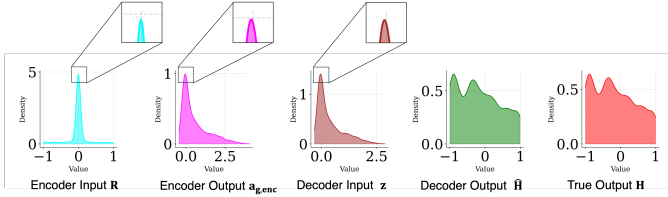


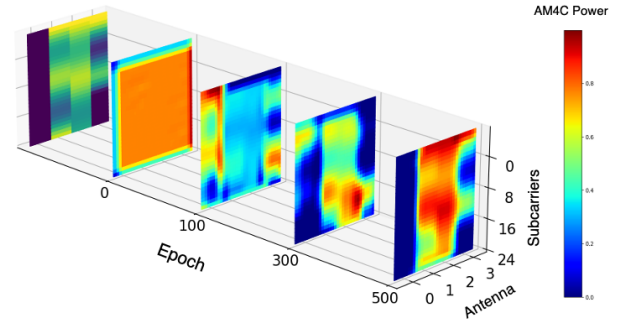
Fig. 5. The distribution of input and output of the model, $p_{\mathbf{R}}(\mathbf{R})$, $p(\mathbf{a}_{g,enc})$, approximated posterior $p(\mathbf{z}|\mathbf{a}_{g,enc})$, $p_{\hat{\mathbf{H}}}(\hat{\mathbf{H}})$, $p_{\mathbf{H}}(\mathbf{H})$

4) Explainable AI for Attention Area of the Channel Matrix:

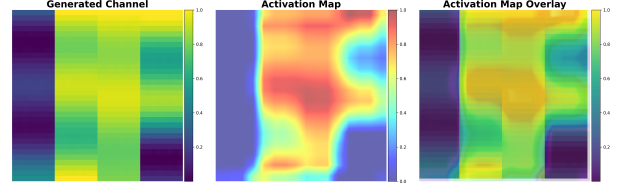
In Fig. 6a, we illustrate how the generator model updates features during training and how the critic model utilizes these features to determine whether they are generative. At epoch zero, it is observed that the critic focuses on the entire channel due to the initial noise distribution. By epoch 100, the critic tends to focus on lower magnitudes, as the generator easily produces these values and the critic seeks to avoid being misled by the generator. However, from epoch 300 to 500, the generator nearly surpasses the critic, prompting the critic to shift its focus to higher magnitudes, which contain more valuable information for classification. The activation map, in Fig. 6b, originates from the first convolutional layer of the critic, indicating that this layer is the closest to extracting information from the generator. This suggests that the critic model focuses on higher magnitudes of resource elements in the TFR to for decision-making.

VI. CONCLUSIONS

In this paper, we propose a new WGAN-GP model for channel estimation, in which the generator exploits the third-order moment of the data to improve its performance. We also propose AM4C as an explainable AI mechanism to obtain the attention area of the critic. We show that the proposed WGAN-GP performs superiorly, achieving lower NMSE than conventional methods, particularly in favorable training settings. Using the proposed AM4C, we demonstrate



(a)



(b)

Fig. 6. (a) The evolution of the attention area of the critic over training epochs. (b) Attention area of the trained critic overlaid over generative channel.

that the critic relies on the high power gain regions of the time-frequency representation (TFR) of the channel to learn and discriminate the generated channel from the true channel.

REFERENCES

- [1] M. Hirzallah, M. Krunz, B. Keciciglu, and B. Hamzeh, "5g new radio unlicensed: Challenges and evaluation," *IEEE Transactions on Cognitive Communications and Networking*, vol. 7, no. 3, p. 689–701, Sep. 2021.
- [2] X. Ma, Z. Gao, F. Gao, and M. Di Renzo, "Model-driven deep learning based channel estimation and feedback for millimeter-wave massive hybrid mimo systems," *IEEE Journal on Selected Areas in Communications*, vol. 39, no. 8, p. 2388–2406, Aug. 2021.
- [3] H. Ye, G. Y. Li, and B.-H. Juang, "Power of deep learning for channel estimation and signal detection in ofdm systems," *IEEE Wireless Communications Letters*, vol. 7, no. 1, p. 114–117, Feb. 2018.
- [4] Hoydis et al., "Sionna: An open-source library for next-generation physical layer research," Mar. 2022. [Online]. Available: <https://nvlabs.github.io/sionna>
- [5] Jiao et al., "Advanced deep learning models for 6g: Overview, opportunities, and challenges," *IEEE Access*, vol. 12, p. 133245–133314, 2024.
- [6] Haider et al., "Gan-based channel estimation for irs-aided communication systems," *IEEE Transactions on Vehicular Technology*, vol. 73, no. 4, p. 6012–6017, Apr. 2024.
- [7] Y. Du, Y. Li, M. Xu, J. Jiang, and W. Wang, "A joint channel estimation and compression method based on GAN in 6G communication systems," *Applied Sciences*, vol. 13, no. 4, p. 2319, Feb. 2023.
- [8] E. Balevi and J. G. Andrews, "Wideband channel estimation with a generative adversarial network," *IEEE Transactions on Wireless Communications*, vol. 20, no. 5, p. 3049–3060, May 2021.
- [9] D. Bau, J.-Y. Zhu, J. Wulff, W. Peebles, B. Zhou, H. Strobelt, and A. Torralba, "Seeing what a gan cannot generate," in *2019 IEEE/CVF International Conference on Computer Vision (ICCV)*. IEEE, Oct. 2019.
- [10] I. Gulrajani, F. Ahmed, M. Arjovsky, V. Dumoulin, and A. Courville, "Improved training of wasserstein gans," in *Proceedings of the 31st International Conference on Neural Information Processing Systems*, 2017, p. 5769–5779.
- [11] M. Arjovsky and L. Bottou, "Towards principled methods for training generative adversarial networks," in *Proceedings of the International Conference on Learning Representations*, Apr. 2017.
- [12] D. P. Kingma and M. Welling, "Auto-encoding variational bayes," *CoRR*, vol. abs/1312.6114, 2013. [Online]. Available: <https://api.semanticscholar.org/CorpusID:216078090>