



TOWARDS TRUSTWORTHY DEEP LEARNING MODELS FOR VHF PROPAGATION: GENERALIZATION, UNCERTAINTY AND EXPLAINABILITY

Ahmad-Mahbubul Alam¹, Bernard Uguen¹, Thierry Marsault²

¹IETR, University of Rennes 1, Rennes, France, {ahmad-mahbubul.alam, bernard.uguen}@univ-rennes.fr

²DGA-MI, Rennes, France, thierry.marsault@intradef.gouv.fr

Keywords: deep correlation alignment (Deep CORAL), conformalized quantile regression (CQR), saliency maps, Grad-CAM, attention mechanism

Mots clés: alignement profond de corrélation, régression quantile conformisée, cartes de saillance, Grad-CAM, mécanisme d'attention

Résumé/Abstract

Les modèles d'apprentissage profond (DL) pour la propagation sans fil sont confrontés à trois défis critiques : la robustesse aux changements de domaine, la quantification de l'incertitude et l'interprétabilité. Cet article aborde ces trois problèmes critiques, progressant vers des modèles DL plus généralisables, fiables et explicables pour les communications sans fil. Nous nous appuyons sur une architecture hybride combinant des métadonnées scalaires avec un profil de paramètre de diffraction, utilisant des blocs résiduels et d'attention pour capturer les effets dépendants du terrain à longue portée dans la bande très haute fréquence (VHF). Nous démontrons que l'adaptation de domaine utilisant l'alignement profond de corrélation (Deep CORAL) améliore substantiellement la généralisation sous les changements de domaine. Pour garantir des prédictions dignes de confiance, nous employons la régression quantile conformisée (CQR), qui fournit des garanties d'incertitude sans distribution et à échantillon fini et corrige les prédictions trop confiantes, produisant des intervalles de prédiction bien calibrés à 90% qui améliorent la fiabilité pour les scénarios de services d'urgence. De plus, en utilisant des techniques d'explicabilité telles que les cartes de saillance, les visualisations d'attention et Grad-CAM, nous montrons que le modèle capture des caractéristiques physiquement significatives. Cette approche intégrée de modélisation adaptative au domaine, consciente de l'incertitude et interprétable ouvre la voie à un déploiement fiable du DL dans les systèmes de communication sans fil du monde réel.

Deep learning (DL) models for wireless propagation face three critical challenges: robustness to domain shifts, uncertainty quantification, and interpretability. This paper addresses these three critical issues, advancing towards more generalizable, reliable, and explainable DL models for wireless communication. We build upon a hybrid architecture combining scalar metadata with a diffraction parameter profile, using residual and attention blocks to capture long-range terrain-dependent effects in the very high frequency (VHF) band. We demonstrate that domain adaptation using deep correlation alignment (Deep CORAL) substantially improves generalization under domain shifts. To ensure trustworthy predictions, we employ conformalized quantile regression (CQR), which provides finite-sample, distribution-free uncertainty guarantees and corrects for overconfident predictions, yielding well-calibrated 90% prediction intervals that enhance reliability for emergency service scenarios. Furthermore, using explainability techniques including saliency maps, attention visualizations, and Grad-CAM, we show that the model captures physically meaningful features. This integrated approach of domain-adaptive, uncertainty-aware, and interpretable modeling paves the way for trustworthy DL deployment in real-world wireless communication systems.

1 Introduction

Accurate wireless signal propagation prediction is essential for network planning and optimization. While physics-based models are interpretable, they are often computationally expensive. Deep learning (DL) provides a powerful alternative, but its practical deployment is limited by three key challenges: poor robustness to domain shifts, lack of reliable uncertainty quantification, and limited interpretability [1]. Building on TriResNet [2], we replace the multi-channel environmental tensor with a single diffraction parameter profile, reducing the input dimensionality by a factor of three while improving accuracy, and address domain shifts, uncertainty quantification, and interpretability.

Most DL models assume that training and testing data are independent and identically distributed (IID), an assumption that is often violated in wireless propagation scenarios. Although domain adaptation methods such as deep correlation alignment (Deep CORAL) have been widely applied in vision, NLP, and speech [3], their application to propagation modeling remains limited. We address cross-domain generalization under simultaneous shifts in frequency bands and transmit antenna heights, demonstrating that this technique substantially improves generalization across heterogeneous domains.

Standard deep learning models produce point estimates without quantifying uncertainty, which limits their use in safety-critical network planning. Probabilistic methods, such as Bayesian neural networks and Monte Carlo Dropout [4–6],

are computationally expensive and rely on strong modeling assumptions. We observe that standard quantile regression produces overconfident and poorly calibrated prediction intervals. To address this, we adopt conformalized quantile regression (CQR), a distribution-free approach that generates efficient prediction intervals with valid coverage. Unlike prior work [1], which mainly focuses on spatial uncertainty, for example LoS/NLoS boundaries, we evaluate uncertainty across key propagation factors, including distance, frequency, antenna height, and obstructions. This effectively transforms point forecasts into risk-aware predictions for robust network planning.

The black-box nature of DL models raises concerns about interpretability and trustworthiness. While explainable artificial intelligence (XAI) has been explored in wireless network control [7–9], its use in propagation prediction remains limited. To address this, we apply complementary interpretability techniques to both metadata and spatial profiles. Saliency analysis of the metadata aligns with propagation theory where distance, frequency, and obstruction increase path loss, while antenna heights and LoS conditions reduce it. Meanwhile, attribution methods including attention, and Grad-CAM consistently highlight physically meaningful regions within the diffraction parameter profiles which reflect terrain-induced diffraction effects. Combined with calibrated uncertainty, this enables transparent and reliable predictions.

2 Problem Definition and Data Preparation

2.1 Problem Formulation

We aim to develop a neural network g_{nn} parameterized by θ_{nn} that predicts reference signal received power (RSRP) at unseen locations:

$$\hat{y}_{dB} = g_{nn}([\mathbf{v}, \mathbf{m}]; \theta_{nn}), \quad (1)$$

where inputs comprise the diffraction parameter profile \mathbf{v} , which encodes knife-edge diffraction values along the transmitter-receiver path, and the metadata vector \mathbf{m} , which contains scalar parameters such as antenna heights, frequency, and distance. Beyond this core regression task, our framework simultaneously ensures robustness to domain shifts, provides calibrated uncertainty intervals, and delivers physically interpretable insights.

2.2 Data Collection and Preprocessing

The dataset comprises raw measurements collected by DGA-MI across diverse French terrains (rural plains, forests, hills, semi-urban) for military applications. Signal strength was averaged over 100-meter blocks to suppress small-scale fading, focusing the model on large-scale path loss and shadowing. Terrain profiles account for Earth curvature using the standard 4/3 Earth radius model [10], and LoS conditions are determined by comparing terrain elevations to the direct transmitter-receiver path. Effective antenna heights are derived from linear fit intercepts in each antenna’s vicinity, consistent with the terrain regression approach specified in ITU-R P.1546 [11]. All profiles are linearly interpolated to match the longest path (84.18 km, 1403 points), preserving continuity and yielding better performance than zero-padding.

2.3 Input Features

The input features consist of metadata \mathbf{m} and the diffraction parameter profile \mathbf{v} . The metadata vector \mathbf{m} includes operating frequencies (30–80 MHz), Tx-Rx distances (up to 50 km), transmitter heights (mostly above 30 m, maximum 80 m), and other scalar parameters listed in Table 1.

Table 1: Metadata Vector \mathbf{m} .

Transmitting antenna height [h_c] (m)	Effective transmitter antenna height [h_{effc}] (m)
Effective receiver antenna height [h_{effr}] (m)	Tx-Rx separation distance [distance] (m)
Operating frequency [frequency] (MHz)	Binary indicator (1 = LoS, 0 = NLoS) [LoS]
Knife-edge diffraction parameter for main obstacle [eng] [10]	Distance from transmitter to main diffracting obstacle [d_{main}] (m)

The diffraction parameter profile \mathbf{v} encodes knife-edge diffraction values along the Tx-Rx path at 60 m resolution, allowing the model to learn terrain-induced attenuation without explicit terrain or clutter data. At each point, v is calculated as

$$v = h_{obs} \cdot \sqrt{\frac{2(d_1 + d_2)}{\lambda d_1 d_2}}, \quad (2)$$

where d_1 and d_2 are distances from the obstacle to the transmitter and receiver, respectively, h_{obs} is the height of the terrain plus clutter above the line-of-sight at distance d_1 from the transmitter, and λ is the wavelength. Tx/Rx endpoints are excluded to avoid extreme v values and numerical instability. The physical context of these endpoints is preserved through the antenna height features contained in the metadata vector \mathbf{m} . Both metadata and the v profile are standardized (zero mean, unit variance) using training set statistics.

3 Model Architecture and Training Configuration

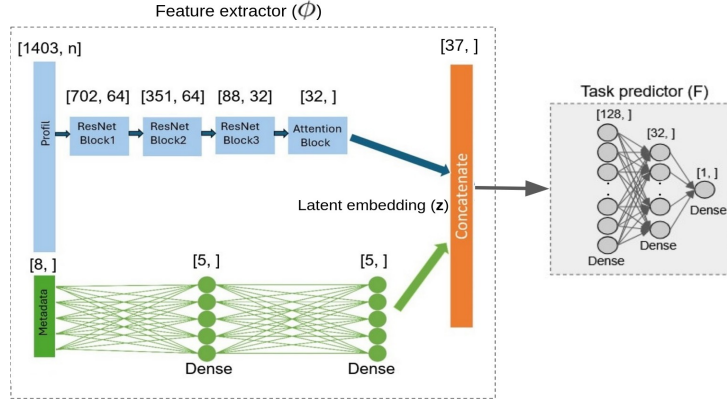


Figure 1: Hybrid architecture combining metadata and diffraction profile.

The model builds on the hybrid architecture TriResNet [2], with separate branches for metadata and spatial profiles. The architecture processes two input types: metadata vector \mathbf{m} through a two-layer FNN branch, and diffraction parameter profile $\mathbf{v} \in \mathbb{R}^{\{1403 \times 1\}}$ through a CNN branch with three 1D ResNet blocks. Residual connections mitigate vanishing gradients and increase the receptive field to capture long-range terrain effects, while an additive attention mechanism after the third ResNet block highlights critical regions of the profile [2]. Outputs from both branches are concatenated to form the latent embedding $\mathbf{z} = \phi(\mathbf{x})$, where $\mathbf{x} = [\mathbf{m}, \mathbf{v}]$ and ϕ is the feature extractor. The embedding \mathbf{z} then feeds into the task predictor F , which outputs $\hat{y} = F(\mathbf{z})$ for the regression task.

All models were trained with the Adam optimizer using an initial learning rate $\eta_0 = 1 \times 10^{-3}$, decaying exponentially every 5 epochs to a minimum of 5×10^{-6} . Training was performed for 145 epochs with batch size 32, selecting the best model based on validation loss.

4 Building Trustworthy Deep Learning Models

4.1 Domain Adaptation and Generalizability

To evaluate generalizability under realistic operational conditions, we introduce a frequency-based domain shift between training and deployment, as models trained on one frequency band often fail to generalize to another due to differing propagation characteristics. The dataset is partitioned by frequency: the source domain (training) includes samples with frequency ≥ 50 MHz, while the target domain (testing) includes samples with frequency < 50 MHz. These domains also differ in transmitter antenna height, with target samples corresponding to lower heights, making the adaptation task both realistic and challenging. Within the source domain, 80% of samples are used for training and 20% for validation. The target domain is split into 60% for unsupervised domain alignment during training and 40% as a held-out test set.

To address this domain shift, we apply Deep CORAL [3], which explicitly matches second-order statistics of latent features between domains by training on mixed mini-batches of labeled source and unlabeled target data to minimize:

$$L = \frac{1}{n_s} \sum_{i=1}^{n_s} (y_i - F(\phi(\mathbf{x}_i^S)))^2 + \lambda_{coral} (|\mu_S - \mu_T|_2^2 + |\mathbf{C}_S - \mathbf{C}_T|_F^2), \quad (3)$$

where (\mathbf{x}_i^S, y_i) denotes the i -th labeled source sample, n_s is the number of source samples in the batch, μ_S , μ_T and \mathbf{C}_S , \mathbf{C}_T are the batch-wise means and covariance matrices of source and target embeddings $\phi(\mathbf{x})$, $\|\cdot\|_F$ denotes the Frobenius norm, λ_{coral}

balances the source MSE and the CORAL loss to encourage domain-invariant embeddings while maintaining predictive accuracy on the source domain, and F is the task predictor.

4.2 Conformal Prediction for Uncertainty Quantification

We adopt CQR [12], a rigorous framework that constructs prediction intervals with distribution-free, finite-sample coverage guarantees. The dataset is partitioned into a training set (60%), a calibration set (20%), and a testing set (20%). The target coverage is set to $1 - \alpha$, where α is the allowed miscoverage rate, i.e., $\alpha = 0.1$ for 90% coverage. The implementation proceeds as follows. The model is trained on the training set to predict lower and upper quantiles by minimizing the pinball loss:

$$L_\tau(y, \hat{y}) = \max(\tau(y - \hat{y}), (\tau - 1)(y - \hat{y})), \quad (4)$$

where y is the true RSRP, \hat{y} is the predicted quantile, and $\tau \in (0, 1)$ is the quantile level. For the i -th calibration sample, the conformity score is defined as:

$$S_i = \max(\hat{y}_l(\mathbf{x}_i) - y_i, y_i - \hat{y}_h(\mathbf{x}_i)), \quad (5)$$

where $\hat{y}_l(\mathbf{x}_i)$ and $\hat{y}_h(\mathbf{x}_i)$ denote the predicted lower and upper quantiles, respectively. A value $S_i > 0$ indicates undercoverage, i.e., the true value lies outside the predicted interval, while $S_i \leq 0$ indicates correct coverage. To guarantee the target coverage $1 - \alpha$, we compute q as the $\lfloor (n + 1)(1 - \alpha) \rfloor$ th order statistic of the conformity scores $\{S_i\}_{i=1}^n$. The calibrated prediction interval is given by:

$$C(\mathbf{x}) = [\hat{y}_l(\mathbf{x}) - q, \hat{y}_h(\mathbf{x}) + q]. \quad (6)$$

4.3 Explainability and Interpretability Analysis

To verify that the model captures physically meaningful propagation effects, we employ complementary explainability techniques: saliency maps [13] for scalar metadata, and Grad-CAM [14] alongside attention mechanisms [2] for spatial profile analysis.

For scalar metadata, feature importance is estimated using saliency maps, which compute the gradient of the predicted attenuation with respect to each input feature. Averaging across samples gives global importance, with the sign indicating whether a feature increases or decreases the prediction. For spatial inputs, we leverage two complementary methods. Grad-CAM [14] computes gradients of the output with respect to the final convolutional layer, highlighting which spatial regions most influence the prediction. The attention block [2] after the third ResNet block learns normalized weights that indicate the importance of each position along the propagation path. Their agreement provides strong evidence of physically meaningful learning.

5 Results and Analysis

5.1 Effect of Domain Adaptation

Under an IID split, the model achieved strong performance, with an RMSE of 3.15 dB and R^2 of 0.94 on the test set, demonstrating its ability to accurately capture propagation trends from metadata and diffraction profiles. However, when a domain shift was introduced (source: ≥ 50 MHz with high Tx antennas; target: < 50 MHz with low Tx antennas), performance degraded significantly. Training on the source domain and evaluating on the target yielded an RMSE of 30.63 dB and R^2 of 0.40, highlighting the sensitivity of deep learning models to distribution shifts in wireless propagation.

To mitigate this issue, we evaluated Deep CORAL with different trade-off parameters $\lambda \in \{1, 10, 100, 1000\}$. All configurations improved performance over the non-adapted baseline, with $\lambda = 10$ yielding the best results. As shown in Table 2, Deep CORAL reduced the RMSE from 30.63 dB to 8.5 dB and achieved an R^2 of 0.78 on the target domain.

	No Adaptation	Deep CORAL ($\lambda=10$)
RMSE (dB)	30.63	8.5
R^2	0.40	0.78

Table 2: Performance of Deep CORAL Domain Adaptation on the Target Domain.

To analyze the effect of domain adaptation, we visualized latent embeddings \mathbf{z} using t-SNE before and after training, as shown in Fig. 2. Before adaptation, source and target embeddings formed distinct clusters. After training with $\lambda = 10$, Deep CORAL achieved significant domain alignment between source and target embeddings. These visualizations confirm that domain alignment correlates with target-domain performance, demonstrating that domain adaptation is essential for generalization.

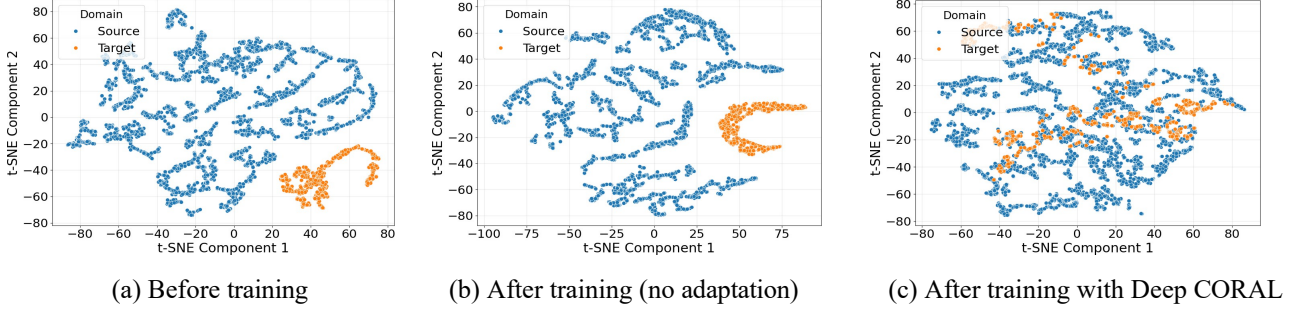


Figure 2: t-SNE visualization of latent embeddings \mathbf{z} with and without domain adaptation.

5.2 Uncertainty Quantification with CQR

We first analyze the raw prediction intervals produced by the baseline QR model without conformal calibration. As shown in Fig. 3, the model captures the heteroscedastic nature of wireless propagation, with prediction interval widths varying systematically across different propagation conditions.

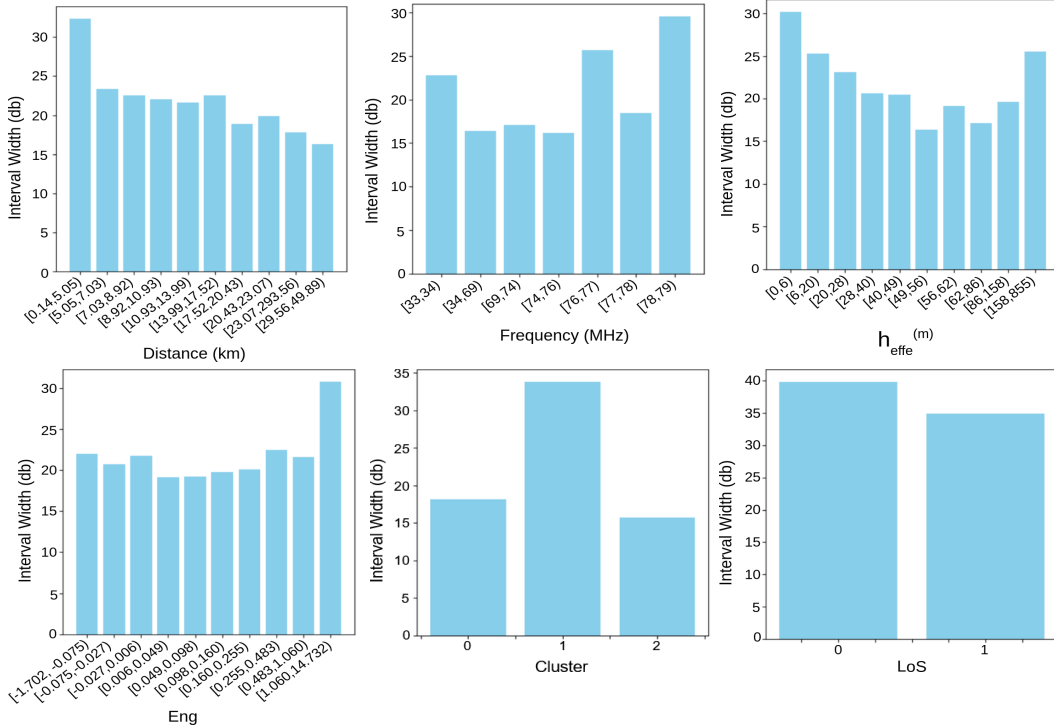


Figure 3: Feature-wise analysis of uncalibrated prediction intervals.

Uncertainty increases in challenging propagation regimes, with short distances below 5 km showing higher variability due to strong multipath. Frequency and transmitter height are correlated in the dataset, with low frequencies paired with low antennas and high frequencies with high antennas. The largest uncertainty occurs at these extreme combinations, reflecting terrain interactions at low heights and obstruction effects at high frequencies. Environmental factors also influence uncertainty. The engagement parameter, which quantifies obstacle intrusion into the Fresnel zone, shows increasing interval widths with larger values, indicating stronger diffraction effects. Additionally, based on k-means clustering of terrain roughness and clutter density, we identified three environmental clusters. Cluster 1, corresponding to mountainous areas with high terrain roughness and clutter density, exhibits the largest uncertainty. This behavior is

consistent with non-line-of-sight conditions ($LoS = 0$), where irregular terrain and obstacles introduce complex diffraction and shadowing effects that are more difficult to predict.

Despite capturing scenario-dependent variability, the uncalibrated QR model exhibits clear overconfidence, falling short of the desired 90% coverage as shown in Figure 4. To restore statistical validity, we apply CQR.

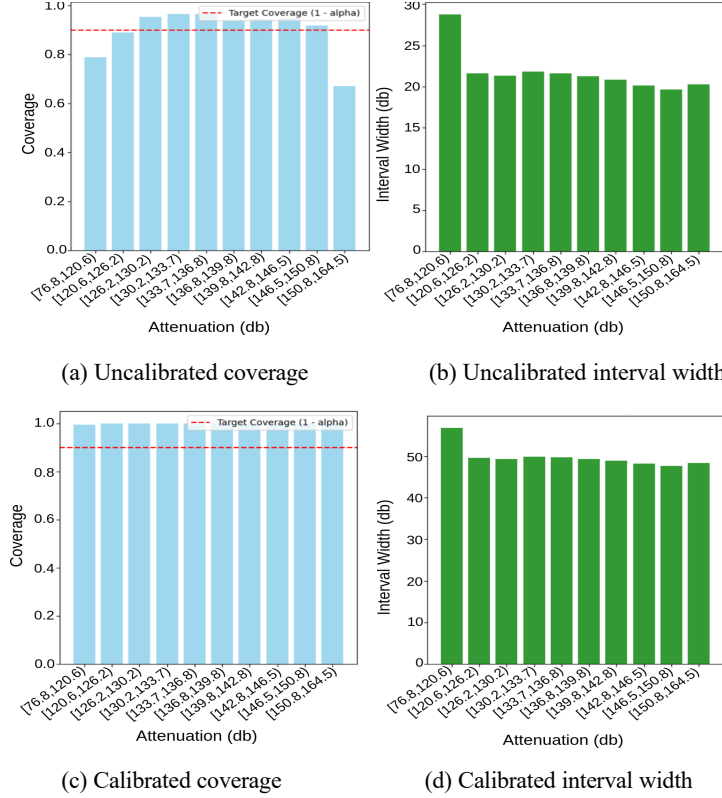


Figure 4: Uncalibrated (top row) and calibrated (bottom row) coverage and interval width across attenuation levels.

After calibration, as shown in Fig. 4, empirical coverage reaches the target 90% across the entire attenuation range. The average interval width increases from 23 dB to 50 dB, reflecting the necessary trade-off between sharpness and reliability. The baseline model captures where uncertainty rises, and CQR provides statistically valid confidence intervals, converting an overconfident regressor into a reliable uncertainty-aware predictor for wireless network planning. In Figs. 3 and 4, each bar contains an equal number of sample points to ensure a fair comparison across the range of values.

5.3 Model Interpretability

We analyze scalar metadata contributions using saliency maps, which measure the sensitivity of predicted attenuation to each input feature. As shown in Fig. 5, distance, frequency, and the obstruction engagement parameter (eng) exhibit strong positive saliency, indicating higher values increase attenuation. In contrast, antenna-related features (h_{e_s} , h_{eff_e} , h_{eff_r}) and the line-of-sight indicator ($LoS = 1$) show negative saliency, meaning higher antennas and clear line-of-sight reduce attenuation. These results confirm that the model learns physically consistent relationships from the metadata.

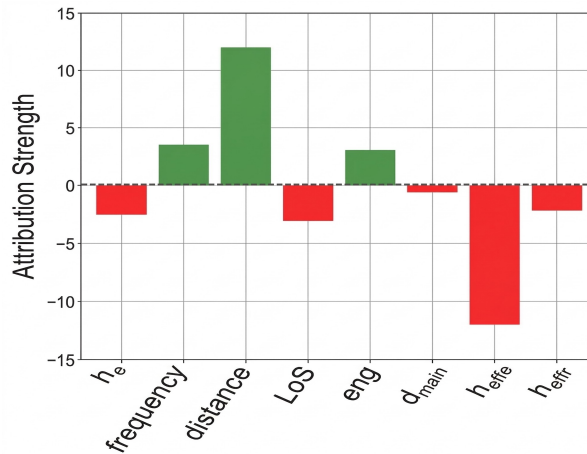


Figure 5: Metadata saliency map showing feature importance for attenuation prediction.

Fig. 6 shows terrain elevation, diffraction parameter (v), Grad-CAM, and attention profiles for a representative link, plotted against normalized distance from transmitter (0) to receiver (1). Grad-CAM, computed from the last convolutional layer before the attention mechanism, highlights terrain-driven features such as dominant obstructions, elevation peaks, deep valleys, and regions with large v values.

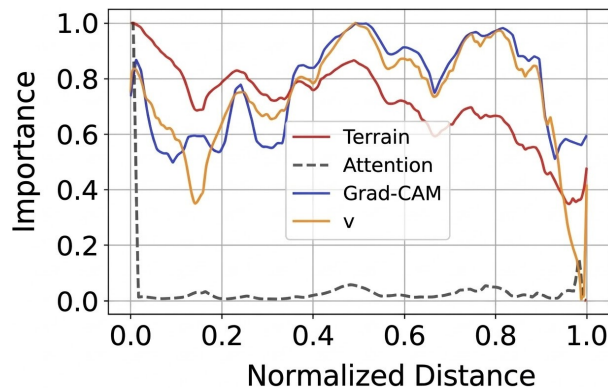


Figure 6: Terrain elevation, diffraction parameter (v), Grad-CAM, and attention profiles for a representative link.

On the other hand, attention profiles strongly emphasize the endpoints corresponding to the transmitter and receiver locations, consistent with the metadata saliency results, which identify antenna-related features, particularly the effective transmitter antenna height (h_{eff_t}), as dominant predictors. Although the endpoint focus visually dominates the profile, attention mechanism also assigns elevated weights to critical terrain features, such as deep valleys and sharp elevation changes, which are similarly highlighted by Grad-CAM activations. This confirms that the model attends to physically relevant features across the entire path.

Overall, the consistent alignment between metadata saliency, attention, and Grad-CAM indicates that the hybrid model relies on physically meaningful propagation factors, including distance, frequency, line-of-sight conditions, diffraction regions, and transmitter/receiver antenna heights.

6 Conclusion

This paper presented a unified framework for VHF propagation prediction addressing domain generalization, uncertainty estimation, and interpretability. The proposed hybrid architecture combines scalar metadata with diffraction parameter profiles, leveraging residual blocks and attention mechanisms to capture the long-range terrain-dependent effects in VHF propagation.

Experimental results showed that frequency domain shift significantly degrades model performance, while domain adaptation mitigates this effect by learning domain-invariant representations, reducing RMSE from 30.63 dB to 8.5 dB.

Although the experiments focused on frequency and antenna height variation, the framework can be extended to other real-world shifts such as terrain or atmospheric conditions.

To ensure reliable predictions, conformalized quantile regression was used to calibrate prediction intervals, correcting the overconfidence of standard quantile regression and achieving well-calibrated intervals with guaranteed 90% coverage. This enables risk-aware propagation predictions for practical network planning and emergency communication scenarios. Explainability analyses using saliency maps, attention, and Grad-CAM showed that the model relies on physically meaningful propagation factors, including distance, frequency, line-of-sight conditions, diffraction regions, and antenna heights. Integrating domain adaptation, uncertainty estimation, and explainability enables trustworthy deep learning models for wireless propagation that are accurate, robust, and interpretable.

Future work will explore generalization across additional propagation scenarios, develop improved conformal calibration strategies to reduce unnecessary interval widening when the desired coverage is already achieved, and further investigate model explainability.

References

- [1] R. He, N. D. Cicco, B. Ai, M. Yang, Y. Miao, and M. Boban, "COST CA20120 INTERACT framework of artificial intelligence-based channel modeling," *IEEE Wireless Communications*, vol. 32, no. 4, pp. 200–207, 2025.
- [2] A.-M. Alam, B. Uguen, and T. Marsault, "Hybrid Deep Learning Model for VHF-Band Propagation Prediction Using ResNet Architecture and Attention Mechanism," in *Proc. 20th Eur. Conf. Antennas Propag. (EuCAP)*, Dublin, Ireland, Apr. 2026, accepted for publication.
- [3] B. Sun and K. Saenko, "Deep CORAL: Correlation alignment for deep domain adaptation," in *Proc. European Conf. Computer Vision (ECCV)*, Springer, 2016, pp. 443–450.
- [4] Z. Fu, N. Fujita, and Y. Terada, "Enhancing reliability in model-based DL reconstruction: A systematic study of MC dropout for uncertainty quantification," in *Proc. ISMRM Annu. Meeting*, 2023.
- [5] S. Brahma, T. Schaffter, C. Kolbitsch, and A. Kofler, "Data-efficient uncertainty quantification for radial cardiac cine MR image reconstruction," in *Proc. ISMRM Annu. Meeting*, 2022.
- [6] Y. Gal and Z. Ghahramani, "Dropout as a Bayesian approximation: Representing model uncertainty in deep learning," in *Proc. International Conf. Machine Learning (ICML)*, 2016, pp. 1050–1059.
- [7] N. Khan, S. Coleri, A. Abdallah, A. Celik, and A. M. Eltawil, "Explainable and robust artificial intelligence for trustworthy resource management in 6G networks," *IEEE Communications Magazine*, vol. 62, no. 4, pp. 50–56, 2023.
- [8] T. Senevirathna, V. H. La, S. Marcha, B. Siniarski, M. Liyanage, and S. Wang, "A survey on XAI for 5G and beyond security: Technical aspects, challenges and research directions," *IEEE Communications Surveys & Tutorials*, vol. 27, no. 2, pp. 941–973, 2024.
- [9] S. K. Jagatheesaperumal, Q.-V. Pham, R. Ruby, Z. Yang, C. Xu, and Z. Zhang, "Explainable AI over the Internet of Things (IoT): Overview, state-of-the-art and future directions," *IEEE Open J. Commun. Soc.*, vol. 3, pp. 2106–2136, 2022.
- [10] T. S. Rappaport, *Wireless Communications: Principles and Practice*, 2nd ed., Upper Saddle River, NJ, USA: Prentice Hall, 2002.
- [11] A. G. Longley and P. L. Rice, "Prediction of Tropospheric Radio Transmission Loss Over Irregular Terrain: A Computer Method—1968," U.S. Department of Commerce, ESSA Technical Report ERL 79-ITS 67, 1968.
- [12] Y. Romano, E. Patterson, and E. Candes, "Conformalized quantile regression," in *Proc. Advances in Neural Information Processing Systems (NeurIPS)*, vol. 32, 2019.
- [13] K. Simonyan, A. Vedaldi, and A. Zisserman, "Deep inside convolutional networks: Visualising image classification models and saliency maps," in *Proc. Int. Conf. Learn. Represent. (ICLR)*, 2014.
- [14] R. R. Selvaraju, M. Cogswell, A. Das, R. Vedantam, D. Parikh, and D. Batra, "Grad-CAM: Visual explanations from deep networks via gradient-based localization," in *Proc. IEEE International Conference on Computer Vision (ICCV)*, 2017, pp. 618–626.