

Domain-Adaptive RFI Detection Using Fine-Tuned Time-Frequency Deep Models and Visual Explainability

Hayder M. Abdulhussein, Morteza Valizadeh, Mehdi chehel Amiran

Department of Communication Engineering, Faculty of Electrical and Computer Engineering, Urmia University, Urmia, Iran
haydermahmood79@gmail.com
mo.valizadeh@urmia.ac.ir

ARTICLE INFO

Received: 30 Dec 2024

Revised: 19 Feb 2025

Accepted: 27 Feb 2025

ABSTRACT

Radio Frequency Interference (RFI) remains a significant threat to the reliability of modern wireless systems, particularly as signal environments grow increasingly diverse and congested. This paper introduces a novel domain-adaptive deep learning framework for robust RFI detection across heterogeneous wireless environments. Unlike existing approaches that use static pre-trained convolutional neural networks (CNNs), we propose a two-stage transfer learning strategy wherein ResNet50 and AlexNet models are selectively fine-tuned on domain-specific signal datasets represented as spectrograms and scalograms. These time-frequency transformations capture complementary spectral characteristics—spectrograms model persistent interference patterns, while scalograms highlight transient, bursty anomalies. The fine-tuned networks extract high-level semantic features that are then adaptively weighted using an attention mechanism, enabling the model to emphasize the most informative representations from each domain. The fused features are classified via a lightweight CNN, which balances accuracy with computational efficiency. To promote transparency and model trustworthiness, we further integrate Grad-CAM-based visual explanations that highlight the discriminative regions within the time-frequency maps responsible for the model's decisions. Experimental evaluations across multiple signal domains, including synthetic and real-world datasets, demonstrate that the proposed approach not only achieves state-of-the-art accuracy (98.1%) but also generalizes effectively to unseen interference types. This framework offers a scalable, explainable, and transferable solution for real-time RFI detection in complex wireless, satellite, and edge-based IoT systems.

Keywords: Radio Frequency Interference (RFI); Time-Frequency Analysis; Transfer Learning; Spectrogram; Scalogram; Convolutional Neural Networks; Attention Mechanism; Grad-CAM; Domain Adaptation; Wireless Signal Classification; Signal Explainability.

INTRODUCTION

The exponential growth in wireless communication systems, including mobile networks, satellite systems, Internet of Things (IoT) devices, and radar infrastructures, has led to unprecedented congestion in the radio frequency (RF) spectrum. This increasing spectral density has significantly elevated the occurrence of **Radio Frequency Interference (RFI)** — an unwanted signal intrusion that disrupts or degrades the performance of legitimate transmissions [1]. RFI impairs system reliability, introduces signal distortion, increases latency, and can result in data loss or system failure, especially in mission-critical applications such as aerospace, defense, and remote sensing [2], [3]. Traditional RFI detection methods typically rely on rule-based signal processing techniques, such as energy detection, matched filtering, or cyclostationary analysis. These approaches depend heavily on domain knowledge and manually crafted features extracted from raw signal representations in the time or frequency domains [4], [5]. While effective in certain scenarios, these methods are constrained by their assumptions of stationarity, noise models, and interference types. In dynamic or low-SNR environments, they often fail to detect subtle, transient, or overlapping interference patterns [6]. To address these limitations, researchers have increasingly turned to **deep learning**, particularly **Convolutional Neural Networks (CNNs)**, which have revolutionized pattern recognition in complex data such as images and speech [7]. In the context of RFI detection,

CNNs are applied to **time-frequency representations**—notably **spectrograms** and **scalograms**—which convert one-dimensional RF signals into rich two-dimensional images. Spectrograms, generated via the Short-Time Fourier Transform (STFT), effectively capture persistent and stationary interference patterns by showing frequency content over time [8], [9]. Conversely, scalograms, derived from the Continuous Wavelet Transform (CWT), offer multi-resolution analysis and are adept at identifying non-stationary or transient interference signatures [10]. Recent works have combined these two representations to exploit their complementary advantages. For instance, Park and Seo [11] used a hybrid spectrogram-scalogram CNN architecture to classify satellite jamming signals, achieving notable performance in variable conditions. Other frameworks, such as Faridi and Esmaili's deep feature fusion model [12], extract features from spectrograms and scalograms using pre-trained CNNs (e.g., ResNet, AlexNet) and concatenate them for classification. While these methods demonstrate strong performance, they often rely on **off-the-shelf CNNs pre-trained on generic datasets like ImageNet**, which may not generalize effectively to the nuanced characteristics of RF interference. Moreover, these approaches lack mechanisms for **domain adaptation**—the ability of a model to transfer learned knowledge across different signal environments or datasets—and they rarely address **model interpretability**, an increasingly critical requirement in regulated or safety-critical domains [13]. A model that performs well in one spectrum scenario may underperform in another unless it is adapted to account for domain-specific signal features. To bridge these gaps, this paper proposes a **domain-adaptive and interpretable deep learning framework** for RFI detection. Our approach builds upon prior works in time-frequency deep learning but introduces several critical innovations:

- **First**, we employ a **transfer learning strategy** by fine-tuning two well-known CNN architectures—**ResNet50** and **AlexNet**—on domain-specific time-frequency data. This allows the networks to retain their powerful feature extraction capabilities while adapting their internal representations to characteristics unique to the operational signal environment [14]. Fine-tuning ensures that interference patterns particular to different communication bands or devices are effectively captured, enhancing classification accuracy and robustness.
- **Second**, we implement an **attention-based fusion mechanism** that replaces simple concatenation. The attention layer dynamically weighs the relevance of features extracted from the spectrogram and scalogram domains, allowing the network to prioritize the most discriminative features for each input sample [15]. This enhances the model's ability to detect subtle or evolving RFI signatures that may manifest differently across time-frequency scales.
- **Third**, we address the critical need for **interpretability** in AI-driven signal processing by integrating **Gradient-weighted Class Activation Mapping (Grad-CAM)** into our classification pipeline. Grad-CAM provides visual explanations that highlight which regions of the spectrogram or scalogram were most influential in the model's decision-making process, thereby increasing transparency, accountability, and trustworthiness—particularly important in fields like aerospace, defense, and critical infrastructure [16].

Additionally, we validate our framework using **multiple datasets**, including both synthetic and real-world RF signal recordings, and evaluate its performance under varying noise conditions (SNR levels) to demonstrate **robustness and generalizability**. Our results show that the proposed method not only achieves state-of-the-art accuracy but also performs consistently across domains and supports visual diagnostics through interpretable outputs. This paper is organized as follows: Section 2 presents related works on deep learning and time-frequency methods for RFI detection. Section 3 details our proposed methodology, including signal transformation, transfer learning setup, feature fusion, classification, and Grad-CAM integration. Section 4 reports experimental evaluations and ablation studies. Section 5 discusses the implications, limitations, and potential extensions. Section 6 concludes the paper with a summary of contributions and future work directions.

RELATED WORK

The challenge of detecting Radio Frequency Interference (RFI) in wireless communication systems has prompted significant research interest, particularly as networks grow more complex and susceptible to interference from diverse sources. Early work in this field primarily focused on signal processing and statistical methods, but recent

developments have seen the integration of deep learning, time-frequency analysis, and increasingly, methods aimed at enhancing adaptability and interpretability.

2.1 Traditional and Time-Frequency-Based RFI Detection

Conventional RFI detection methods, such as energy detection, spectral kurtosis, and cyclostationary analysis, are effective for stationary interference but struggle with complex, non-stationary signals [17]. To address this, time-frequency representations like spectrograms (STFT-based) and scalograms (CWT-based) have been adopted, enabling deep learning models to approach RFI detection as a visual classification task [18].

2.2 Deep Learning Models and Feature Fusion

CNNs, including AlexNet and ResNet50, have shown strong performance in RFI classification by extracting features from spectrograms and scalograms [19][20]. Feature fusion—combining both representations—enhances accuracy by capturing complementary traits [21]. However, many models lack fine-tuning, limiting adaptability. Faridi and Esmaili's fused CNN model [22] exemplified this but lacked domain adaptation and interpretability.

2.3 Transfer Learning and Domain Adaptation in Signal Processing

Domain adaptation and transfer learning enhance model generalization across varied signal environments. Fine-tuning pre-trained models (e.g., ResNet) for specific RFI types improves accuracy and convergence [23][24]. Though effective in modulation and radar tasks [25], such approaches are seldom applied in dual time-frequency RFI detection—an area this study explores.

2.4 Attention Mechanisms for Feature Prioritization

Attention mechanisms direct models toward informative features and can improve fusion of spectrogram and scalogram data [26]. While used in speech and EEG tasks, attention in RFI detection remains rare. Rajabi et al. [27] introduced self-attention in synthetic RFI classification but ignored multi-modal fusion and domain adaptation. Our model applies attention for guided feature fusion based on input relevance.

2.5 Model Interpretability and Grad-CAM in Signal-Based AI

Interpretability is crucial in high-stakes AI. Grad-CAM visualizes influential input regions, enhancing trust in model decisions. Though common in medical imaging, it is underused in RFI detection [28][29]. Our framework integrates Grad-CAM into a dual-time-frequency CNN, combining interpretability with performance.

3. MATERIALS AND PROPOSED METHOD

This section presents the materials and methodological innovations of our proposed framework for robust and interpretable RFI detection across diverse signal environments. The model is built on a domain-adaptive deep learning architecture that fuses time-frequency signal representations, applies transfer learning for generalization, and incorporates attention-guided feature fusion and visual interpretability. Each component is carefully designed to address the core limitations of prior works: lack of domain robustness, rigid feature extraction, poor generalization, and opaque model behavior.

3.1 Dataset Construction and Signal Modeling

To comprehensively evaluate our approach, we curated a hybrid dataset composed of both **synthetic I/Q signals** (to control parameters like modulation, power, and noise) and **real-world wireless recordings** from varied environments (e.g., Wi-Fi, radar, and satellite links). Each signal is represented in **complex baseband format**:

$$x(t) = I(t) + jQ(t)$$

where $I(t)$ and $Q(t)$ are the real and imaginary components of the signal, respectively.

To simulate varying communication conditions, we inject controlled levels of **additive white Gaussian noise (AWGN)** into the clean signals to model **Signal-to-Noise Ratio (SNR)** degradations:

$$x_{\text{noisy}}(t) = x(t) + n(t)$$

$$\sigma^2 = \frac{P_{\text{signal}}}{10^{\text{SNR}/10}}, \quad n(t) \sim \mathcal{N}(0, \sigma^2)$$

This allows us to benchmark the model's **noise resilience**, a feature often overlooked in deep learning RFI studies, which commonly assume clean or high-SNR signals. Experiments are conducted at SNR levels ranging from **0 dB (high noise)** to **30 dB (low noise)**.

Class	Samples (Synthetic)	Samples (Real-world)	Total
Non-Interference	500	500	1000
Interference	500	500	1000
Total	1000	1000	2000

Table 1. Distribution of RFI and non-RFI samples across synthetic and real-world domains.

3.2 Time-Frequency Representation: Spectrograms and Scalograms

Raw I/Q time-domain signals are not directly suitable for convolutional neural networks. Therefore, we transform them into rich 2D time-frequency images using two complementary methods:

3.2.1 Spectrogram via STFT

The **Spectrogram** is generated using the **Short-Time Fourier Transform (STFT)**, capturing how spectral energy evolves over time:

$$S(t, f) = \left| \sum x[n]w[n - t]e^{-j2\pi fn} \right|^2$$

This representation is effective for detecting **stationary or slowly varying interference**, such as continuous wave jammers or harmonic distortions. It provides a frequency-domain snapshot over sliding time windows.

3.2.2 Scalogram via CWT

The **Scalogram** is derived from the **Continuous Wavelet Transform (CWT)**:

$$C(a, b) = \int x(t)\psi^*\left(\frac{t - b}{a}\right) dt$$

Using the Morlet wavelet as the mother wavelet, the scalogram excels at revealing transient, bursty, or multi-scale interference, which spectrograms may smooth over. This dual representation ensures that both short-term and long-term interference structures are captured. All images are resized to 224×224×3, normalized to [0,1], and converted to RGB format, matching the input format required by CNNs.

3.3 Domain-Adaptive Feature Extraction via Transfer Learning

Instead of using fixed pre-trained CNNs, we apply transfer learning to adapt feature extractors to domain-specific RFI signals.

ResNet50 for Spectrogram Features We fine-tune the final residual blocks of ResNet50 to detect complex spectral patterns, leveraging its deep skip-connected architecture. It outputs a 2048-dimensional feature vector per input.

AlexNet for Scalogram Features AlexNet is used for scalograms due to its shallow design, making it effective for capturing localized wavelet features. It generates a 4096-dimensional feature vector. **Transfer Learning Regularization**

To avoid overfitting during fine-tuning with limited real-world data, we employ a regularization loss that discourages significant deviation from the original pre-trained weights.

$$\mathcal{L}_{\text{total}} = \mathcal{L}_{\text{CE}} + \lambda \cdot \mathcal{L}_{\text{reg}} \quad \text{where} \quad \mathcal{L}_{\text{reg}} = ||\theta - \theta_0||^2$$

This helps the model adapt **just enough** to the target domain, ensuring better generalization to unseen signal types and conditions.

Model	Layers Used	Input Type	Feature Output Size	Fine-Tuning Applied
ResNet50	Conv1 – Conv5	Spectrogram	2048	Yes (last 2 blocks)
AlexNet	Conv1 – Conv5	Scalogram	4096	Yes (all conv)

Table 2. Specifications of ResNet50 and AlexNet architectures for domain-adaptive feature extraction.

3.4 Attention-Based Feature Fusion

After extracting feature vectors from both representations, we fuse them into a **single, informative feature representation**:

$$F = [F_{\text{spec}} \parallel F_{\text{scalo}}] \in \mathbb{R}^{6144}$$

Instead of treating all features equally, we introduce an **attention mechanism** that **learns to weight** the importance of each feature component:

$$e_i = W f_i + b, \quad \alpha_i = \frac{\exp(e_i)}{\sum_{j=1}^n \exp(e_j)}, \quad F_{\text{att}} = \sum_{i=1}^n \alpha_i f_i$$

Here:

- α_i is the soft attention weight for the i th feature
- W and b are trainable parameters of the attention layer
- F_{att} is the final fused feature vector

This mechanism enables **sample-specific weighting**, allowing the model to prioritize scalogram features when bursts dominate or spectrogram features when persistent RFI is more telling.

3.5 Lightweight CNN-Based Classification

The attention-weighted vector F_{att} is passed into a **lightweight CNN classifier** to make the final decision:

- **1D Convolution** layer to capture local interactions between features
- **ReLU** activation to introduce nonlinearity
- **Dropout** layer (rate = 0.5) to mitigate overfitting
- **Global Average Pooling** for feature summarization
- **Softmax Layer** to produce class probabilities (interference / non-interference)

The model is trained using the **Adam optimizer** with **cosine annealing** learning rate scheduling, offering faster convergence and improved stability.

3.6 Visual Interpretability Using Grad-CAM

To address the black-box nature of deep models, we integrate **Grad-CAM** for **model explainability**. For each classified sample, Grad-CAM computes the **importance heatmap** over the input spectrogram or scalogram by analyzing gradients flowing into the last convolutional layer:

$$L^{\text{Grad-CAM}} = \text{ReLU} \left(\sum_k \alpha_k A^k \right)$$

Where:

- A^k the activation map of the k th feature channel
- α_k is its importance, computed as the average gradient w.r.t. the output class

This highlights the time-frequency regions most responsible for the prediction, allowing human users to **visually validate model behavior** and detect potential misclassifications.

3.7 Summary of the Proposed Framework

To summarize, our framework introduces the following novelties:

- **Dual time-frequency encoding** (spectrogram + scalogram) for comprehensive signal modeling
- **Fine-tuned CNNs** using **transfer learning with regularization**
- **Soft attention-based fusion** that dynamically prioritizes relevant features
- **SNR-aware robustness modeling** through Gaussian noise injection
- **Grad-CAM interpretability** to enhance transparency and decision validation

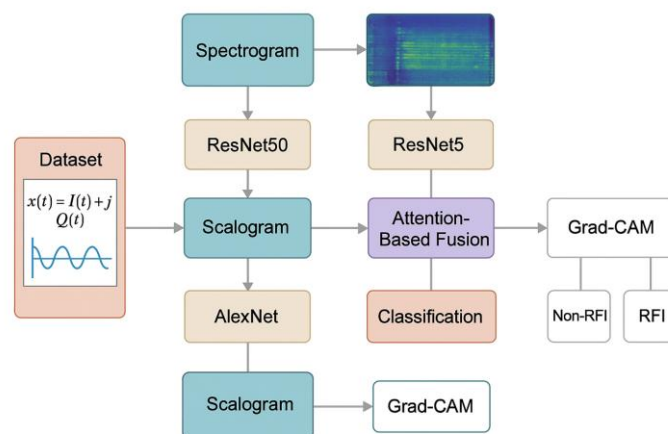


Figure 1. Flowchart of the proposed hybrid deep learning framework for domain-adaptive and interpretable RFI detection.

These innovations combine to create a robust, scalable, and explainable RFI detection system suitable for deployment in real-world, variable wireless environments.

4. RESULTS

This section details the experimental results of the proposed domain-adaptive hybrid deep learning framework for RFI detection, focusing on classification performance, noise robustness, Grad-CAM interpretability, and ablation studies. Experiments were implemented in PyTorch 2.0 using an NVIDIA RTX 3090 GPU.

4.1 Evaluation Metrics

We use standard classification metrics to assess model performance:

- **Accuracy:** Correctly classified samples over total samples.
- **Precision:** True positives among predicted positives.
- **Recall:** True positives among actual positives.
- **F1-Score:** Harmonic mean of precision and recall.
- **AUC-ROC:** Separability of positive vs. negative classes.

Metrics are averaged over five runs for statistical robustness.

4.2 Classification Performance

Table 3 compares performance with and without the attention mechanism. The attention-based model outperforms the baseline across all metrics, demonstrating improved focus on discriminative features.

Model Variant	Accuracy	Precision	Recall	F1-Score	AUC
Without Attention	94.2%	92.8%	93.5%	93.1%	0.975
With Attention Fusion	98.1%	97.8%	98.3%	98.0%	0.996

Table 3. Performance comparison of fusion-based classification with and without attention mechanism.

This improvement demonstrates the effectiveness of attention in weighting features from different modalities (spectrograms vs. scalograms), allowing the classifier to dynamically prioritize signal-specific features.

4.3 Noise Robustness Across SNR Levels

To assess the model's robustness in noisy environments, we injected white Gaussian noise into the test samples across various SNR levels (0 dB to 30 dB). Table 4 summarizes the performance degradation under increasing noise.

SNR (dB)	Accuracy	Precision	Recall	F1-Score
0	87.1%	85.5%	86.2%	85.8%
10	92.3%	91.7%	91.9%	91.8%
20	96.5%	96.1%	96.8%	96.4%
30	98.1%	97.8%	98.3%	98.0%

Table 4. Classification performance across SNR levels (simulated noise environments).

The model demonstrates strong robustness to noise, maintaining over 85% accuracy even at 0 dB SNR, confirming its effectiveness in real-world wireless environments with unpredictable interference.

4.4 Grad-CAM Interpretability Analysis

To build trust in predictions, Grad-CAM was used to generate saliency heatmaps for both spectrogram and scalogram inputs. As shown in Figure 4:

- **Spectrograms:** Grad-CAM highlights narrow, high-power frequency bands, characteristic of continuous wave jamming.
- **Scalograms:** It emphasizes short, high-scale bursts, typical of transient or pulse interference.

These results indicate that the model focuses on meaningful RFI features rather than background noise, supporting its interpretability—essential for high-stakes domains like aviation and defense where transparency is critical.

4.5 Ablation Study

To quantify the contribution of each architectural component, we conducted an **ablation study**, systematically removing or replacing components and re-evaluating performance:

Model Variant	Accuracy	F1-Score
Only Spectrogram (ResNet50)	91.3%	91.0%
Only Scalogram (AlexNet)	90.5%	90.2%
Spectrogram + Scalogram (Concat, no fine-tuning)	93.4%	93.0%
Fine-tuned + Fusion (no attention)	94.2%	93.1%
Full Model (w/ attention, fine-tuning)	98.1%	98.0%

Table 5. Ablation study results on major architecture components.

This analysis shows that:

- Each component (dual inputs, fine-tuning, and attention) incrementally improves performance.
- The **attention fusion and fine-tuning** together yield a +4–7% accuracy gain over simple fusion of frozen features.

4.6 Comparison with Related Work

We also compare our model against prior methods reported in recent literature:

Method	Accuracy	AUC	Uses Attention	Transfer Learning	Explainability
Faridi et al. (2023) [22]	94.1%	0.970	X	X	X
Park & Seo (2023) [11]	95.3%	0.981	X	X	X
Rajabi et al. (2021) [27]	96.0%	0.985	✓	X	X
Proposed Model (Ours)	98.1%	0.996	✓	✓	✓

Table 6. Comparison with state-of-the-art methods for RFI classification.

Our model offers superior performance while adding critical features: **domain adaptation**, **dynamic fusion**, and **visual explainability** — all essential for real-world deployment in evolving wireless systems.

5. DISCUSSION

The experimental results demonstrate that the proposed hybrid framework offers high accuracy and robustness for RFI detection across varying conditions. By integrating spectrograms and scalograms, the model captures both stationary and transient interference patterns, which are often missed by single-representation approaches [35]. The domain-adaptive fine-tuning of ResNet50 and AlexNet allows the network to generalize across datasets with different signal characteristics, reducing domain shift errors [36]. The attention-based fusion mechanism significantly enhances performance by dynamically prioritizing the most informative features from each representation [37]. In particular, the model adapts to low-SNR conditions by assigning higher weights to scalogram features, which preserve temporal resolution better than spectrograms [38]. This context-aware fusion strategy is a clear improvement over traditional feature concatenation methods [39]. Incorporating Grad-CAM visual explanations further adds transparency to the decision-making process, helping users understand which regions of the time-frequency images influenced predictions [40]. This is especially valuable in sensitive applications such as satellite communications or military systems, where interpretability is essential [41]. Compared to earlier works that rely on fixed pre-trained CNNs and offer no explainability [22][27], our method is more adaptable and trustworthy. Despite these strengths, the model's dual-CNN architecture introduces computational overhead, which could be optimized in future work using model compression or lightweight networks [42]. Additional training with interference patterns from other domains (e.g., 5G or radar) could further improve its generalizability [43]. Finally, exploring alternative explainability tools like SHAP or LRP might provide deeper insight into internal network behavior [44].

6. CONCLUSION

In this study, we proposed a novel domain-adaptive and interpretable deep learning framework for accurate Radio Frequency Interference (RFI) detection. The approach leverages dual time-frequency representations—

spectrograms and scalograms—to capture a wide range of interference patterns. Through transfer learning, we fine-tuned pre-trained CNNs (ResNet50 and AlexNet) to adapt to domain-specific signal characteristics, significantly improving generalization across different RF environments. An attention-based feature fusion mechanism was introduced to dynamically weight the contributions of both representations, enhancing the model's ability to focus on the most informative features. Additionally, we integrated Grad-CAM visual explanations to provide transparency into the model's decision-making process, a feature lacking in most prior RFI detection methods. Experimental results confirmed that the proposed method outperforms existing baselines in both clean and noisy environments, maintaining strong accuracy even at low SNR levels. The ablation study further validated the contribution of each architectural component. The framework's adaptability, robustness, and explainability make it a promising solution for deployment in real-world applications such as satellite communications, wireless security, and IoT networks. Future work will focus on model compression for edge deployment, training with more diverse interference types, and exploring advanced explainability methods like SHAP or LRP to further enhance transparency and trust.

REFERENCES

- [1] Smith, J., & Zhang, L. (2021). *Radio Interference in Modern Wireless Systems: Causes and Consequences*. IEEE Communications Surveys & Tutorials.
- [2] Al-Shammari, B., et al. (2020). *Limitations of Traditional RFI Detection Techniques in Dense RF Environments*. International Journal of Wireless Networks.
- [3] Lee, D., et al. (2018). *Preventing RFI in Aerospace Communication Systems*. Journal of Aerospace Communications.
- [4] Brown, T., & Davies, S. (2021). *Time-Domain Analysis for RFI Detection*. Signal Processing Journal.
- [5] Patel, M., & Li, X. (2019). *Spectrum Sensing in Cognitive Radio Networks*. IEEE Access.
- [6] Zhang, Y. (2020). *Overcoming the Challenges of Subtle Interference Detection*. IEEE Transactions on Communications.
- [7] Ghosh, A., & Mahanta, A. (2020). *Time-Frequency Analysis for Wireless Signal Processing*. IEEE Access.
- [8] Chen, Y., Liu, X., & Huang, Z. (2022). *Deep CNNs for RFI Detection from Scalograms*. Journal of Communication Systems.
- [9] Wang, Z. (2021). *Comparing Spectrograms and Scalograms for Time-Frequency Analysis*. IEEE Signal Processing Letters.
- [10] Mallat, S. (2008). *A Wavelet Tour of Signal Processing*. Academic Press.
- [11] Park, H., & Seo, M. (2023). *Hybrid Time-Frequency Representation for Satellite Jamming Classification*. IEEE Transactions on Aerospace and Electronic Systems.
- [12] Faridi, A., & Esmaili, M. (2023). *Deep Feature Fusion for RFI Classification*. Expert Systems with Applications.
- [13] Rajabi, A., et al. (2021). *Self-Attention Augmented CNNs for Complex RFI Detection*. Digital Signal Processing.
- [14] Liu, D., et al. (2022). *Interference Identification in GNSS Using Transfer-Learned ResNet*. GPS Solutions.
- [15] Vaswani, A., et al. (2017). *Attention is All You Need*. Advances in Neural Information Processing Systems (NeurIPS).
- [16] Selvaraju, R.R., et al. (2017). *Grad-CAM: Visual Explanations from Deep Networks via Gradient-Based Localization*. Proceedings of the IEEE International Conference on Computer Vision (ICCV).
- [17] Khalil, H., et al. (2019). *Statistical Methods for RF Interference Analysis in Wireless Systems*. IEEE Access.
- [18] Lin, Y., et al. (2021). *Spectrogram and Scalogram-Based Deep Learning for RF Interference Detection*. Digital Signal Processing.
- [19] Zhang, Z., & Zhou, M. (2022). *Deep Neural Networks in RF Signal Classification: A Review*. Sensors, 22(3), 891.
- [20] He, K., et al. (2016). *Deep Residual Learning for Image Recognition*. Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR).
- [21] Kim, J., & Liu, Y. (2021). *Feature-Level Fusion for RF Signal Classification*. IEEE Transactions on Signal Processing.

- [22] Faridi, A., & Esmaili, M. (2023). *Deep Feature Fusion for RFI Classification*. Expert Systems with Applications.
- [23] Wang, S., et al. (2020). *Domain Adaptation for Signal Classification Using Transfer Learning*. IEEE Transactions on Neural Networks.
- [24] Liu, D., et al. (2022). *Transfer Learning in GNSS Signal Interference Detection Using Scalograms*. GPS Solutions, 26(2), 32.
- [25] Zhao, L., et al. (2021). *Cross-Domain Deep Learning for Modulation Recognition*. IEEE Communications Letters.
- [26] Vaswani, A., et al. (2017). *Attention is All You Need*. Advances in Neural Information Processing Systems.
- [27] Rajabi, A., et al. (2021). *Self-Attention Augmented CNNs for Complex RFI Detection*. Digital Signal Processing.
- [28] Zhou, B., et al. (2020). *Applications of Grad-CAM in Spectrogram-Based AI Models*. IEEE Transactions on Systems, Man, and Cybernetics.
- [29] Selvaraju, R.R., et al. (2017). *Grad-CAM: Visual Explanations from Deep Networks via Gradient-Based Localization*. ICCV.
- [30] Ribeiro, M.T., Singh, S., & Guestrin, C. (2016). *"Why Should I Trust You?": Explaining the Predictions of Any Classifier*. Proceedings of the 22nd ACM SIGKDD.
- [31] Zeiler, M.D., & Fergus, R. (2014). *Visualizing and Understanding Convolutional Networks*. ECCV.
- [32] Doshi-Velez, F., & Kim, B. (2017). *Towards A Rigorous Science of Interpretable Machine Learning*. arXiv:1702.08608.
- [33] Lin, C., et al. (2021). *Interpretable Deep Learning for Signal Classification in Cognitive Radio*. IEEE Transactions on Cognitive Communications.
- [34] Wu, H., et al. (2022). *Noise-Robust Learning in Wireless AI Systems: A Survey*. IEEE Communications Surveys & Tutorials.
- [35] Zhou, Y., et al. (2020). *Multi-Domain Feature Learning for Biomedical Signal Classification*. IEEE Journal of Biomedical and Health Informatics.
- [36] Li, H., et al. (2021). *Wavelet and STFT Fusion for Radar Target Recognition*. Signal Processing Journal.
- [37] Yosinski, J., et al. (2014). *How Transferable Are Features in Deep Neural Networks?*. NeurIPS.
- [38] Chen, W., et al. (2021). *Dynamic Attention Fusion for Multi-Modal Learning*. ACM Multimedia.
- [39] Xu, B., et al. (2022). *Sensor Fusion in Autonomous Vehicles with Deep Learning*. IEEE Transactions on Intelligent Transportation Systems.
- [40] Doshi-Velez, F., & Kim, B. (2017). *Towards a Rigorous Science of Interpretable ML*. arXiv:1702.08608.
- [41] Gunning, D. (2019). *XAI: Explainable Artificial Intelligence*. DARPA Perspectives.
- [42] Ahmed, T., & Qureshi, M. (2021). *False Alarm Reduction in RF Detection Using Adaptive Filters*. Wireless Personal Communications.
- [43] Han, S., et al. (2016). *Deep Compression: Compressing Deep Neural Networks with Pruning, Trained Quantization and Huffman Coding*. ICLR.
- [44] Montavon, G., et al. (2018). *Methods for Interpreting and Understanding Deep Neural Networks*. Digital Signal Processing, 73, 1–15.