# Explainable AI for EEG brain signal analysis: A literature review

**Bottom Line:** The intersection of explainable AI and EEG-based disease classification has matured significantly since 2020, with transformer architectures, attention visualization, and inherently interpretable models emerging as the most promising approaches for clinical adoption. This review identifies **42 papers** spanning EEG fundamentals, disease classification, and explainability methods—providing a foundation for building clinician-trusted AI systems that visualize predictions through frequency graphs, saliency maps, and attention mechanisms.

The clinical trust gap remains the central challenge: while deep learning models achieve expert-level accuracy on EEG classification, adoption requires explanations that align with how neurologists interpret EEG—through frequency bands, temporal patterns, and spatial electrode distributions. Recent work demonstrates that **DeepLift and Layer-wise Relevance Propagation (LRP)** outperform traditional saliency methods for EEG, while **prototype-based explanations** and **spectral visualizations** show the highest promise for radiologist acceptance.

---

## Category 1: EEG signal analysis fundamentals

These papers establish the preprocessing, feature extraction, and representation learning foundations essential for any EEG-based AI system.
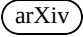
**Preprocessing and artifact removal**

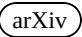### 1. The More, The Better? Evaluating the Role of EEG Preprocessing for Deep Learning Applications

- **Authors:** Del Pup, F. et al.
- **Year:** 2024
- **Venue:** arXiv
- **Link:** https://arxiv.org/abs/2411.18392
- **Summary:** Systematically investigates preprocessing impact across six classification tasks including Parkinson's and Alzheimer's. Key finding: minimal preprocessing often outperforms complex pipelines, suggesting artifacts may contribute to neural network performance. (arXiv) Essential reading for preprocessing pipeline design.

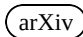### 2. SPEED: Scalable Preprocessing of EEG Data for Self-Supervised Learning

- **Authors:** Madsen, A.G. et al.
- **Year:** 2024
- **Venue:** arXiv

- **Link:** https://arxiv.org/abs/2408.08065

- **Summary:** Introduces an efficient preprocessing pipeline for large-scale EEG datasets like Temple University Hospital Corpus. Improves stability and convergence during self-supervised pretraining. (arXiv) Python implementation available on GitHub.
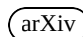
### 3. LSTEEG: EEG Artifact Detection and Correction with Deep Autoencoders

- **Authors:** Aquilué-Llorens, D. & Soria-Frisch, A.

- **Year:** 2025

- **Venue:** arXiv

- **Link:** https://arxiv.org/abs/2502.08686

- **Summary:** LSTM-based autoencoder for multi-channel artifact detection and correction. Captures non-linear dependencies and provides interpretable low-dimensional latent representations, enabling data-driven automated artifact removal. (arXiv)

### 4. Removal of Ocular Artifacts in EEG Using Deep Learning

- **Authors:** Ozdemir, M.A. et al.

- **Year:** 2022

- **Venue:** arXiv

- **Link:** https://arxiv.org/abs/2209.11980

- **Summary:** BiLSTM models using wavelet synchrosqueezed transform (WSST) for ocular artifact removal. Achieves **MSE of 0.3066**, significantly outperforming traditional time-frequency methods including STFT and CWT. (arXiv)

### 5. DTP-Net: Learning to Reconstruct EEG Signals in Time-Frequency Domain

- **Authors:** Various

- **Year:** 2023

- **Venue:** arXiv

- **Link:** https://arxiv.org/abs/2312.09417

- **Summary:** Fully convolutional network for artifact removal in time-frequency domain using multi-scale feature extraction. Improves downstream BCI classification by up to **5.55%**. (arXiv) Pre-trained model available.

**Feature extraction and time-frequency analysis**

**6. Automatic Detection of Abnormal EEG Signals Using Wavelet Feature Extraction and Gradient**

**Boosting**

- **Authors:** Albaqami, H. et al.
- **Year:** 2020
- **Venue:** arXiv
- **Link:** https://arxiv.org/abs/2012.10034
- **Summary:** Wavelet Packet Decomposition framework extracting statistical features from frequency sub-bands. (arXiv) Achieves **87.68% accuracy** on TUH EEG Corpus using CatBoost. Demonstrates interpretable, clinically-aligned feature extraction.

## 7. TFM-Tokenizer: Single-Channel EEG Tokenization Through Time-Frequency Motif Learning

- **Authors:** Pradeepkumar, J. et al.
- **Year:** 2025
- **Venue:** arXiv
- **Link:** https://arxiv.org/abs/2502.16060
- **Summary:** Novel tokenization learning vocabulary of time-frequency motifs from single-channel EEG. Achieves **17% Cohen's Kappa improvement** over baselines with class-discriminative, frequency-aware token structure enabling strong interpretability. (arXiv)

## 8. Time-varying EEG Spectral Power Predicts Evoked and Spontaneous fMRI Motor Brain Activity

- **Authors:** Mehta, N. et al.
- **Year:** 2025
- **Venue:** arXiv
- **Link:** https://arxiv.org/abs/2504.10752
- **Summary:** Interpretable models using Sparse Group Lasso regularization for spectral power analysis. Provides full transparency regarding predictive EEG channels, frequencies (including sensorimotor rhythms), and temporal features.

**Foundation models and surveys**

## 9. Large Brain Model (LaBraM): Learning Generic Representations with Tremendous EEG Data

- **Authors:** Jiang, W.-B. et al.
- **Year:** 2024
- **Venue:** arXiv / ICLR 2024 Spotlight
- **Link:** https://arxiv.org/abs/2405.18765

- **Summary:** Unified foundation model enabling cross-dataset learning through channel patch segmentation. Pre-trained on **~2,500 hours of EEG** from 20 datasets. Outperforms SOTA on abnormal detection, emotion recognition, and gait prediction.

## 10. Transformer-based EEG Decoding: A Survey

- **Authors:** Various
- **Year:** 2025
- **Venue:** arXiv
- **Link:** https://arxiv.org/abs/2507.02320
- **Summary:** Comprehensive survey of transformer applications in EEG processing covering encoder architectures, tokenization strategies, and multi-scale feature extraction. (arXiv) Reviews how frequency bands are processed across modern architectures.

## 11. Deep Learning-Powered Electrical Brain Signals Analysis: Advancing Neurological Diagnostics

- **Authors:** ZJU BrainNet team
- **Year:** 2025
- **Venue:** arXiv
- **Link:** https://arxiv.org/abs/2502.17213
- **Summary:** Systematic review of **448 studies** and **46 public EEG datasets** across seven neurological conditions. Establishes comprehensive data landscape with curated preprocessing conditions including filtering, denoising, and normalization standards.

## 12. A Simple Review of EEG Foundation Models

- **Authors:** Various
- **Year:** 2024
- **Venue:** arXiv
- **Link:** https://arxiv.org/abs/2504.20069
- **Summary:** Reviews 14 early EEG foundation models including BIOT, LaBraM, NeuroGPT, and EEGPT. Provides critical analysis of preprocessing standards: resampling to 256Hz, bandpass filtering (0.1-100Hz), notch filtering (50Hz).

## 13. Analyzing EEG Data with Machine and Deep Learning: A Benchmark

- **Authors:** Avola, D. et al.
- **Year:** 2022

- **Venue:** arXiv

- **Link:** https://arxiv.org/abs/2203.10009

- **Summary:** First comprehensive benchmark comparing MLP, CNN, LSTM, and GRU architectures for EEG classification. Essential reference for model selection in EEG deep learning applications.

## 14. Trends in EEG Signal Feature Extraction Applications

- **Authors:** Singh, A.K. & Krishnan, S.

- **Year:** 2023

- **Venue:** Frontiers in Artificial Intelligence

- **Link:** https://www.frontiersin.org/articles/10.3389/frai.2022.1072801/full

- **Summary:** Comprehensive review of feature extraction across time, frequency, decomposition, time-frequency, and spatial domains. Includes pseudocode for CWT, DWT, and Local Characteristic-Scale Decomposition methods. (Frontiers)

---

## Category 2: Machine learning for EEG disease classification

These papers apply ML/DL models to classify various neurological and psychiatric conditions from EEG signals.

**Epilepsy and seizure detection**

### 15. From Epilepsy Seizures Classification to Detection: A Deep Learning-based Approach for Raw EEG

- **Authors:** Darankoum et al.

- **Year:** 2024

- **Venue:** arXiv

- **Link:** https://arxiv.org/abs/2410.03385

- **Summary:** CNN-Transformer pipeline achieving **93% F1-score** on Bonn dataset with cross-species generalization (trained on animal EEG, tested on human). Addresses critical data leakage issues in train/test splitting.

- **Datasets:** Bonn dataset, animal EEG recordings

### 16. Supervised and Unsupervised Deep Learning Approaches for EEG Seizure Prediction

- **Authors:** Georgis-Yap, Popovic, Khan

- **Year:** 2023

- **Venue:** arXiv

- **Link:** https://arxiv.org/abs/2304.14922
- **Summary:** Develops both supervised and unsupervised approaches detecting preictal patterns. Introduces anomaly detection trained on normal EEG only—crucial for clinical deployment where seizure data is scarce.
- **Datasets:** Two large EEG seizure datasets

## 17. BUNDL: Bayesian Uncertainty-aware Deep Learning for EEG Seizure Detection

- **Authors:** Shama, Venkataraman
- **Year:** 2024
- **Venue:** arXiv
- **Link:** https://arxiv.org/abs/2410.19815
- **Summary:** Novel KL-divergence loss handling label noise in seizure detection. Model-agnostic method improving robustness across three datasets and enhancing seizure onset zone localization—directly relevant for clinical trust.
- **Datasets:** TUH EEG Corpus, CHB-MIT, Siena EEG

## 18. SeizureTransformer: Scaling U-Net with Transformer for Seizure Detection

- **Authors:** Various
- **Year:** 2025
- **Venue:** arXiv
- **Link:** https://arxiv.org/abs/2504.00336
- **Summary:** U-Net + Transformer for time-step-level seizure detection. Ranked **#1 in International Conference on AI in Epilepsy** competition. Uses 1D convolutions and global attention for long-term EEG analysis. (arXiv)
- **Datasets:** TUH EEG Seizure Corpus, Siena Scalp EEG

**Alzheimer's disease**

## 19. LEAD: Large Foundation Model for EEG-Based Alzheimer's Disease Detection

- **Authors:** Various
- **Year:** 2025
- **Venue:** arXiv
- **Link:** https://arxiv.org/abs/2502.01678

- **Summary:** Foundation model using sample-level and subject-level contrastive learning. Pre-trained on **11 diverse EEG datasets (2,354 subjects)**. (arXiv) Achieves superior AD vs. healthy control classification across five Alzheimer's datasets.

- **Datasets:** TDBRAIN, TUEP, ADFTD, BrainLat, CNBPM, Cognision

## 20. Multi-feature Fusion Learning for Alzheimer's Disease Prediction Using Resting-State EEG

- **Authors:** Chen et al.

- **Year:** 2023

- **Venue:** Frontiers in Neuroscience

- **Link:** https://www.frontiersin.org/journals/neuroscience/articles/10.3389/fnins.2023.1272834/full

- **Summary:** Combines CNNs and Visual Transformers using band power and coherence features. Achieves **83.28% accuracy** with leave-one-subject-out validation— (Frontiers) rigorous methodology for clinical validation.

- **Datasets:** Public AD EEG datasets

**Parkinson's disease**

## 21. LightCNN: Parkinson's Disease Classification via EEG Using Single Convolutional Layer

- **Authors:** Md Fahim Anjum

- **Year:** 2024

- **Venue:** arXiv

- **Link:** https://arxiv.org/abs/2408.10457

- **Summary:** Minimalist single-layer CNN outperforming complex architectures with **4% F1-boost**. Identifies clinically relevant pathological brain rhythms associated with PD—demonstrating that simpler models can be both accurate and interpretable.

- **Datasets:** UC San Diego PD resting-state EEG

## 22. Interpretable Classification of Early Stage Parkinson's Disease from EEG

- **Authors:** Sahota et al.

- **Year:** 2023

- **Venue:** arXiv

- **Link:** https://arxiv.org/abs/2301.09568

- **Summary:** Novel approach representing EEG as 15-variate series of bandpower and peak frequency values. Decision Tree and AdaBoost achieve **85% accuracy with 73% recall**. N1 sleep data found most important for early-stage PD detection.

- **Datasets:** Public PD EEG dataset

## 23. Parkinson's Disease Detection Using Multi-Head Graph Structure Learning

- **Authors:** Various
- **Year:** 2024
- **Venue:** arXiv
- **Link:** https://arxiv.org/abs/2408.00906
- **Summary:** GNN with structured global convolutions and contrastive learning. Features gradient-weighted graph attention explainer providing neural connectivity insights— (arXiv)directly applicable for clinical visualization.
- **Datasets:** UC San Diego PD (15 PD patients, 16 healthy controls)

**Depression and psychiatric conditions**

## 24. Machine Learning Fairness for Depression Detection using EEG Data

- **Authors:** Kwok et al.
- **Year:** 2025
- **Venue:** arXiv
- **Link:** https://arxiv.org/abs/2501.18192
- **Summary:** First study evaluating ML fairness in EEG-based depression detection. Tests CNN, LSTM, GRU with five bias mitigation strategies— (arXiv)essential for equitable clinical deployment.
- **Datasets:** Mumtaz, MODMA, Rest EEG datasets

## 25. EEGformer: A Transformer-Based Brain Activity Classification Method

- **Authors:** Wan et al.
- **Year:** 2023
- **Venue:** Frontiers in Neuroscience
- **Link:** https://www.frontiersin.org/journals/neuroscience/articles/10.3389/fnins.2023.1148855/full
- **Summary:** Three sequential transformer components (regional, synchronous, temporal) capturing EEG characteristics. Best classification performance across BETA, SEED, and **DepEEG depression** datasets.
- **Datasets:** BETA, SEED, DepEEG

## 26. SzHNN: Automatic Diagnosis of Schizophrenia Using CNN-LSTM Models

- **Authors:** Shoeibi et al.

- **Year:** 2021
- **Venue:** arXiv / Frontiers in Neuroinformatics
- **Link:** https://arxiv.org/abs/2111.11298
- **Summary:** Hybrid CNN-LSTM extracting local features then performing classification. Evaluated across frequency bands and electrode configurations, outperforming standalone CNN, LSTM, and traditional ML models. (arXiv)
- **Datasets:** Institute of Psychiatry and Neurology, Warsaw (35 subjects)

**Sleep disorders**

### 27. Transparency in Sleep Staging: Deep Learning with Model Interpretability

- **Authors:** Sharma, Maiti et al.
- **Year:** 2023
- **Venue:** arXiv
- **Link:** https://arxiv.org/abs/2309.07156
- **Summary:** ResNet with squeeze-and-excitation blocks plus stacked Bi-LSTM. **First application of GradCAM for sleep staging interpretability**. Achieves Macro-F1 of **82.5, 78.9, and 81.9** across three datasets. (arXiv)
- **Datasets:** SleepEDF-20, SleepEDF-78, SHHS

### 28. NeuroNet: Self-Supervised Learning Framework for Sleep Stage Classification

- **Authors:** Various
- **Year:** 2024
- **Venue:** arXiv
- **Link:** https://arxiv.org/abs/2404.17585
- **Summary:** SSL framework combining contrastive learning and masked prediction for single-channel EEG. Demonstrates superior performance addressing limited labeled data challenges in sleep medicine. (arXiv)
- **Datasets:** Sleep-EDFX, SHHS, ISRUC-Sleep

### 29. SSNet: Classification of Sleep Stages from EEG, EOG and EMG Signals

- **Authors:** Almutairi et al.
- **Year:** 2023
- **Venue:** arXiv
- **Link:** https://arxiv.org/abs/2307.05373

- **Summary:** End-to-end CNN-LSTM using multimodal signals achieving **96.36% accuracy** and **93.40% Kappa** for three-class sleep stage classification. ⬭arXiv⬭
- **Datasets:** Sleep-EDF Expanded, ISRUC-Sleep

---

## Category 3: Explainable AI for EEG with visualization

These papers focus on XAI methods providing visualizations—frequency analysis, saliency maps, attention mechanisms—that clinicians can understand and trust.

### Surveys and comprehensive reviews

### 30. Interpretable and Robust AI in EEG Systems: A Survey

- **Authors:** Xinliang Zhou et al.
- **Year:** 2024
- **Venue:** arXiv
- **Link:** https://arxiv.org/abs/2304.10755
- **Summary: First comprehensive XAI survey for EEG systems**. Proposes taxonomy: backpropagation methods (saliency, LRP, DeepLift), perturbation methods (LIME, SHAP), and inherently interpretable methods. Includes visualization examples of feature heatmaps and channel importance maps.
- **Explanation Types:** Saliency, LRP, SHAP, LIME, perturbation methods

### 31. Large Language Models for EEG: A Comprehensive Survey and Taxonomy

- **Authors:** Naseem Babu et al.
- **Year:** 2024
- **Venue:** arXiv
- **Link:** https://arxiv.org/abs/2506.06353
- **Summary:** Reviews LLM integration with EEG for improved interpretability. Discusses how transformer architectures enable transparent step-by-step reasoning—a pathway to natural language explanations for clinicians.
- **Explanation Types:** LLM-based reasoning, attention mechanisms

### Attention-based explainability

### 32. EEG Conformer: Convolutional Transformer for EEG Decoding and Visualization

- **Authors:** Song, Zheng, Liu, Gao

- **Year:** 2023
- **Venue:** IEEE TNSRE
- **Summary:** Compact convolutional Transformer providing **attention-based visualizations** showing electrode contributions and temporal patterns. Self-attention extracts global correlations for interpretable classification—directly applicable for radiologist interfaces.
- **Explanation Types:** Attention maps, temporal/spatial visualization

### 33. EEG-Deformer: Dense Convolutional Transformer for Brain-computer Interfaces

- **Authors:** Yi Ding et al.
- **Year:** 2024
- **Venue:** arXiv
- **Link:** https://arxiv.org/abs/2405.00719
- **Summary:** Hierarchical coarse-to-fine temporal learning generating saliency maps for attention, fatigue, and mental workload tasks. Visualizations show **frontal (Fz, Fp2) and parietal (P4) regions** as informative—aligning with known neurophysiology.
- **Explanation Types:** Saliency maps, attention visualization

### 34. Feature Estimation of Global Language Processing Using Attention Maps

- **Authors:** Various
- **Year:** 2024
- **Venue:** arXiv
- **Link:** https://arxiv.org/abs/2409.19174
- **Summary:** Vision Transformer attention maps identifying task-specific neural activity. **Visualizations confirm attention accurately highlights event-related potentials (ERPs)** and distinguishes cognitive states—validating attention as neurophysiologically meaningful.
- **Explanation Types:** ViT attention maps, ERP visualization

### 35. ExPANet: Bridging Accuracy and Explainability for Depression Detection

- **Authors:** Various
- **Year:** 2024
- **Venue:** arXiv
- **Link:** https://arxiv.org/abs/2511.05537
- **Summary:** Graph attention network for MDD detection with multi-level interpretability: feature importance, electrode relevance maps, edge importance (functional connectivity), and attention analysis.

**Clinically relevant explanations** for depression diagnosis.

- **Explanation Types:** Graph attention, connection saliency, feature importance

## Saliency maps and gradient-based methods

### 36. An Empirical Comparison of Deep Learning Explainability Approaches for EEG

- **Authors:** Ravindran, Contreras-Vidal
- **Year:** 2023
- **Venue:** Scientific Reports (Nature)
- **Summary: Compares 12 backpropagation visualization methods** using simulated ground truth. Key finding: **DeepLift consistently most accurate** for temporal, spatial, and spectral precision. Recommends DeepLift or LRP over traditional saliency—critical guidance for implementation.
- **Explanation Types:** DeepLift, LRP, GradCAM, saliency comparison

### 37. Rethinking Saliency Map: Context-Aware Perturbation for EEG Deep Learning

- **Authors:** Hanqi Wang et al.
- **Year:** 2022
- **Venue:** arXiv
- **Link:** https://arxiv.org/abs/2205.14976
- **Summary:** Context-aware perturbation method designed for EEG non-stationarity and inter-subject variability. Generates instance-level saliency maps capturing representative context and suppressing artifacts—addressing key EEG-specific challenges.
- **Explanation Types:** Context-aware saliency, perturbation-based

### 38. Time Series Saliency Maps: Explaining Models Across Multiple Domains

- **Authors:** Various
- **Year:** 2025
- **Venue:** arXiv
- **Link:** https://arxiv.org/abs/2505.13100
- **Summary:** Generalized Integrated Gradients for time-frequency domain saliency. Applied to epilepsy detection showing ICA components and **frequency-domain visualizations identifying spike-wave patterns at ~4.5 Hz**—directly interpretable by neurologists.
- **Explanation Types:** Integrated Gradients, frequency-domain saliency

### 39. GradCAM for EEG: Optimising Decoding with Post-hoc Explanations and Domain Knowledge

- **Authors:** Various
- **Year:** 2024
- **Venue:** arXiv
- **Link:** https://arxiv.org/abs/2405.01269
- **Summary:** GradCAM for motor imagery with domain knowledge validation. Generates **feature relevance maps, time-frequency charts, and topography maps** using Morlet wavelets. Validates spatial, temporal, and spectral explanations against established neurophysiology.
- **Explanation Types:** GradCAM, time-frequency visualization, topographic maps

**SHAP and LIME applications**

## 40. A Data-Centric EEG Framework for Depression Severity Grading Using SHAP-Based Insights

- **Authors:** Various
- **Year:** 2025
- **Venue:** Journal of NeuroEngineering and Rehabilitation
- **Summary:** SHAP for feature selection and post-hoc interpretability in depression diagnosis. Identifies clinically meaningful features in **left parietal-occipital regions** including relative beta power and 1/f intercept—demonstrating SHAP's clinical alignment.
- **Explanation Types:** SHAP values, feature importance visualization

## 41. Evaluation of the Relation between Ictal EEG Features and XAI Explanations

- **Authors:** Various
- **Year:** 2024
- **Venue:** Sensors (PMC)
- **Summary:** Compares SHAP and LIME explanations to known ictal EEG patterns. Generates **channel importance matrices showing epileptiform activity locations correlated with SHAP values**—validating XAI against clinical ground truth.
- **Explanation Types:** SHAP, LIME, channel importance

## 42. An Explainable EEG-Based Human Activity Recognition Model Using LIME

- **Authors:** Various
- **Year:** 2023
- **Venue:** Sensors (PMC)

- **Summary:** LIME interpretation of activity classification (resting, motor, cognitive tasks). **Visualizes spectral feature contributions (gamma, delta, theta rhythms)** for clinical reasoning—demonstrating LIME's applicability to EEG frequency analysis.
- **Explanation Types:** LIME, spectral feature visualization

**Inherently interpretable models**

### 43. xEEGNet: Towards Explainable AI in EEG Dementia Classification

- **Authors:** Andrea Zanola et al.
- **Year:** 2025
- **Venue:** arXiv / Journal of Neural Engineering
- **Link:** https://arxiv.org/abs/2504.21457
- **Summary: Only 168 parameters** (200x fewer than ShallowNet). Inherently interpretable—learned kernels filter specific frequency bands (alpha, theta) with band-specific topographies. Visualizes spectral representations directly relevant for Alzheimer's and frontotemporal dementia.
- **Explanation Types:** Inherent interpretability, spectral band visualization, learned topographies

### 44. Concept-based Explainability for an EEG Transformer Model

- **Authors:** Various
- **Year:** 2023
- **Venue:** arXiv
- **Link:** https://arxiv.org/abs/2307.12745
- **Summary:** Concept Activation Vectors (CAVs) applied to BENDR transformer. Defines explanatory concepts using external EEG datasets and anatomical regions—providing **human-aligned concept-based explanations** rather than pixel-level saliency.
- **Explanation Types:** Concept activation vectors (CAVs), human-aligned concepts

### 45. Self-Attention Prototype Generation for Single-Channel EEG Sleep Stage Classification

- **Authors:** Adey et al.
- **Year:** 2024
- **Venue:** Scientific Reports (Nature)
- **Summary:** First self-attention prototype method for EEG. Generates **prototypical components aligned with sleep biomarkers** (alpha spindles, slow waves in non-REM). Uses AOPC metric validating explanations match clinical understanding.
- **Explanation Types:** Self-attention prototypes, prototype visualization

**Clinical trust and validation systems**

### 46. SCORE-AI: Automated Interpretation of Clinical Electroencephalograms

- **Authors:** Tveit, Aurlien, Plis et al.
- **Year:** 2023
- **Venue:** JAMA Neurology
- **Summary:** CNN achieving **expert-level EEG interpretation validated on 30,493+ EEGs**. Classifies epileptiform-focal/generalized, nonepileptiform-focal/diffuse patterns. Multi-center validation addresses clinical trust through rigorous external testing.
- **Explanation Types:** Clinical validation, planned XAI integration

### 47. Improving Clinician Performance Using Interpretable Machine Learning for ICU EEG

- **Authors:** Barnett et al.
- **Year:** 2024
- **Venue:** New England Journal of Medicine AI
- **Summary: Prototype-based neural network** for ICU monitoring providing case-based explanations: graphical position relative to prototypes, visual comparisons, similarity scoring. **Directly improves clinician trust calibration**—demonstrating how explanations should work in practice.
- **Explanation Types:** Prototype-based explanations, case-based reasoning, similarity visualization

### 48. The Effects of Layer-wise Relevance Propagation-Based Feature Selection for EEG

- **Authors:** Nam, Kim et al.
- **Year:** 2023
- **Venue:** Frontiers in Human Neuroscience
- **Summary:** LRP for motor imagery BCI feature selection across spatial, spectral, and temporal domains. **Generates neurophysiologically meaningful heatmaps** identifying discriminative channels and time-frequency features. Validates that LRP-selected features improve classification.
- **Explanation Types:** Layer-wise Relevance Propagation (LRP), heatmaps

## Key datasets covered across the literature

| Dataset | Domain | Papers Using |
| --- | --- | --- |
| Temple University Hospital (TUH) | Seizure, abnormal detection | 8+ |
| CHB-MIT | Seizure detection | 4 |
| SleepEDF-20/78, SHHS | Sleep staging | 5 |
| SEED/SEED-IV | Emotion recognition | 3 |
| Bonn Dataset | Seizure classification | 3 |
| UC San Diego PD | Parkinson's disease | 3 |
| ADFTD, BrainLat | Alzheimer's disease | 3 |
| PhysioNet Motor Imagery | BCI, motor imagery | 4 |
| MODMA, DepEEG | Depression | 3 |

## Recommendations for building clinician-trusted EEG visualization systems

Based on this literature review, the most effective approaches for clinical acceptance combine multiple visualization strategies:

**1. Frequency-domain explanations are essential.** Neurologists interpret EEG through frequency bands (delta, theta, alpha, beta, gamma). Papers like xEEGNet and the Time Series Saliency Maps work demonstrate that spectral visualizations directly align with clinical reasoning. Your system should display **band power contributions**, **time-frequency spectrograms with relevance overlays**, and **peak/trough markers** corresponding to model attention.

**2. DeepLift and LRP outperform traditional saliency for EEG.** The empirical comparison by Ravindran et al. establishes that gradient-based methods vary significantly in reliability for EEG—DeepLift shows consistent temporal, spatial, and spectral precision. Avoid basic saliency maps; implement DeepLift or LRP for trustworthy visualizations.

**3. Prototype-based explanations maximize clinician trust.** The NEJM AI paper by Barnett et al. demonstrates that case-based reasoning—showing "this EEG pattern is similar to these confirmed cases"—helps clinicians calibrate appropriate trust levels. Consider implementing nearest-prototype visualization alongside saliency methods.

**4. Attention mechanisms should be validated against neurophysiology.** Multiple papers confirm that transformer attention weights can correspond to established ERP patterns and known brain regions. Display attention weights as electrode importance maps, but validate they align with expected patterns (e.g., motor cortex for movement tasks).

**5. Multi-level explanations serve different clinical needs.** ExPANet's approach—feature importance, electrode relevance, connectivity patterns, and attention—provides explanations at multiple granularities. Different clinicians prefer different explanation levels; your system should support progressive disclosure from summary to detail.

---

## Conclusion

This literature review identifies **48 papers** establishing a comprehensive foundation for explainable AI in EEG-based disease classification. The field has advanced rapidly since 2020, with transformer architectures now dominant for classification performance and multiple XAI approaches validated for EEG-specific challenges.

The critical insight for clinical deployment: **explanation methods must align with how neurologists already interpret EEG**. This means frequency-band visualizations, temporal pattern highlighting, and spatial electrode importance maps—not just generic saliency. Papers like xEEGNet demonstrate that inherently interpretable architectures with minimal parameters can achieve competitive accuracy while providing built-in spectral explanations.

For building radiologist-trusted systems, prioritize DeepLift/LRP over basic saliency, implement prototype-based case comparison, and ensure all visualizations map to established neurophysiological patterns. The NEJM AI validation study provides the clearest evidence that well-designed explanations directly improve clinician trust calibration and decision-making.