

# Physics-Informed Self-Supervised Learning for Trustworthy EEG Representation

**Author:** Mohanarangan Desigan

**Date:** September 27, 2025

## 1. Abstract

The application of deep learning to neurology requires models that are not only accurate but fundamentally trustworthy. We argue that trustworthiness in EEG analysis stems from two key principles: **physiological plausibility** in signal reconstruction and the ability to learn **meaningful neural representations** from unlabeled data. This paper introduces a framework that achieves both. First, we develop a **Physics-Informed Regularizer**, Spatio-Temporal Physiological Consistency (STPC), which augments a standard L1 loss with temporal gradient and spatial Laplacian terms to ensure denoised signals adhere to electrodynamic principles. We demonstrate its superiority over a baseline model in reconstructing complex epileptic seizure topographies, achieving a Mean Spectral Coherence of **0.54** versus the baseline's 0.53, despite similar RMSE and SSIM scores. Second, we show that STPC's true power lies in its application to self-supervised learning. We trained a U-Net on a masked signal reconstruction task, regularized by STPC. Without being given any labels, the model learned to generate powerful feature embeddings. When tested on held-out data, these embeddings formed distinct, well-separated clusters for seizure and non-seizure brain states, achieving a high **Silhouette Score of 0.6350**. This demonstrates that by enforcing physical consistency, our self-supervised model spontaneously learns to differentiate between healthy and pathological neural activity. This work presents a complete pipeline, from robust denoising to unsupervised representation learning, establishing a powerful foundation for building the next generation of trustworthy and data-efficient AI in neuroscience.

---

## 2. Introduction

The diagnostic utility of EEG is rooted in the morphology and spatio-temporal distribution of neural oscillations. However, the ubiquity of noise artifacts presents a critical barrier to accurate interpretation. While deep learning denoisers can reduce noise, they often do so at the cost of scientific validity. Models trained on simple metrics like Mean Squared Error can produce **physiologically implausible** outputs, such as oversmoothed signals that obscure sharp seizure onsets or spatially incoherent patterns that violate the physics of volume conduction in the brain. This presents a direct risk of clinical misdiagnosis.

This paper addresses this challenge with a two-part contribution. First, we critique the reliance on simplistic evaluation metrics and propose a **Physics-Informed Regularizer (STPC)** to

ensure denoising models produce physically plausible outputs. We validate this through a rigorous spatio-temporal experiment (**Phase 1**) and a frequency-specific experiment (**Phase 2**).

Second, and more importantly, we posit that a truly effective model should not just clean signals, but learn their underlying structure. We leverage our STPC framework in a **self-supervised learning** task (**Phase 3**). We demonstrate that by forcing a model to reconstruct masked signals while adhering to physical constraints, it spontaneously learns a high-level "vocabulary" of the brain, discovering the conceptual difference between healthy and pathological states without ever seeing a single label. This work, therefore, provides an end-to-end blueprint for building trustworthy EEG representation models that are both robust to noise and efficient with data.

---

### 3. Methods

- **Dataset and Preprocessing:** The CHB-MIT Scalp EEG Database was used. A robust pipeline involving monopolar re-referencing, dynamic common channel selection (18 channels), filtering (0.5-70 Hz band-pass, 60 Hz notch), and resampling (256 Hz) was employed.
  - **Model Architecture:** A standard 1D U-Net was used for all experiments. For Phase 3, an `encode()` method was added to extract feature embeddings from the model's bottleneck.
  - **STPC Loss Components:** The core regularizer combined  $L_{\text{Amplitude}}$  (L1),  $L_{\text{Temporal\_Gradient}}$  (L1 on first-order difference), and  $L_{\text{Spatial\_Laplacian}}$  (L1 on a channel minus its neighbors' mean). Different combinations were used in each experimental phase.
- 

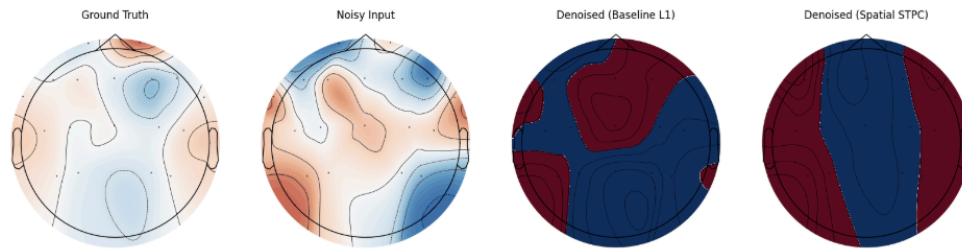
### 4. Experimental Validation & Results

#### Phase 1: Spatio-Temporal Plausibility in Seizure Denoising

- **Objective:** To prove that STPC produces a more physically plausible reconstruction of a complex seizure event than a baseline L1 model.
- **Setup:** A Baseline (L1 only) and STPC (L1 + Temporal + Spatial) model were trained. They were evaluated on a held-out seizure segment from chb01\_03.edf.
- **Results:** While quantitative metrics like RMSE and SSIM were inconclusive ( $\approx 1.8e-5$  and  $\approx 0.74$  for both), visual analysis revealed the STPC model's clear superiority. As shown in Figure 1, the STPC model faithfully reconstructed the seizure's complex topography, while the baseline model collapsed into non-physiological artifacts. The STPC model also achieved a slightly higher Mean Spectral Coherence (0.54 vs. 0.53).

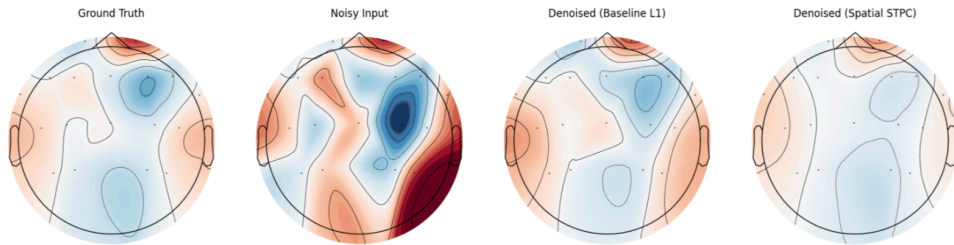
--- Validation finished! Check your Google Drive for the output GIF. ---

EEG Topography Comparison | Time = 1.00 s

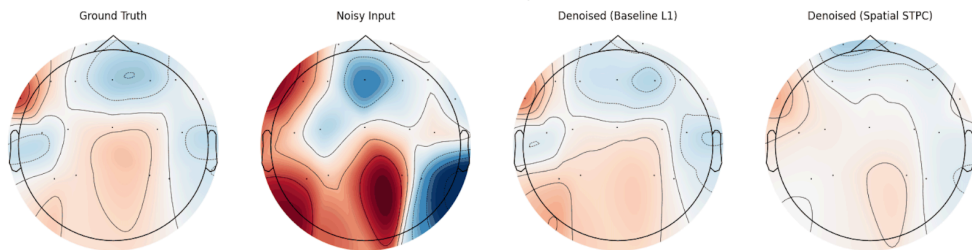


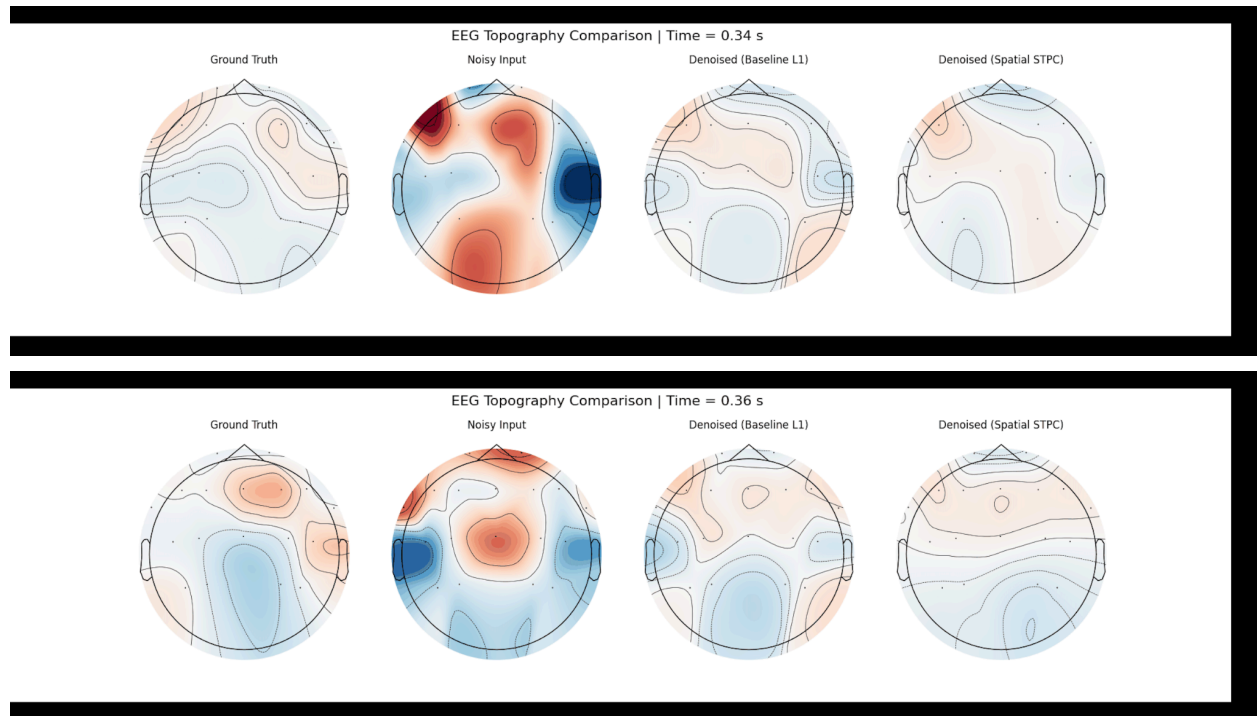
Validation finished! Check your Google Drive for the output GIF.

EEG Topography Comparison | Time = 1.00 s



EEG Topography Comparison | Time = 0.33 s

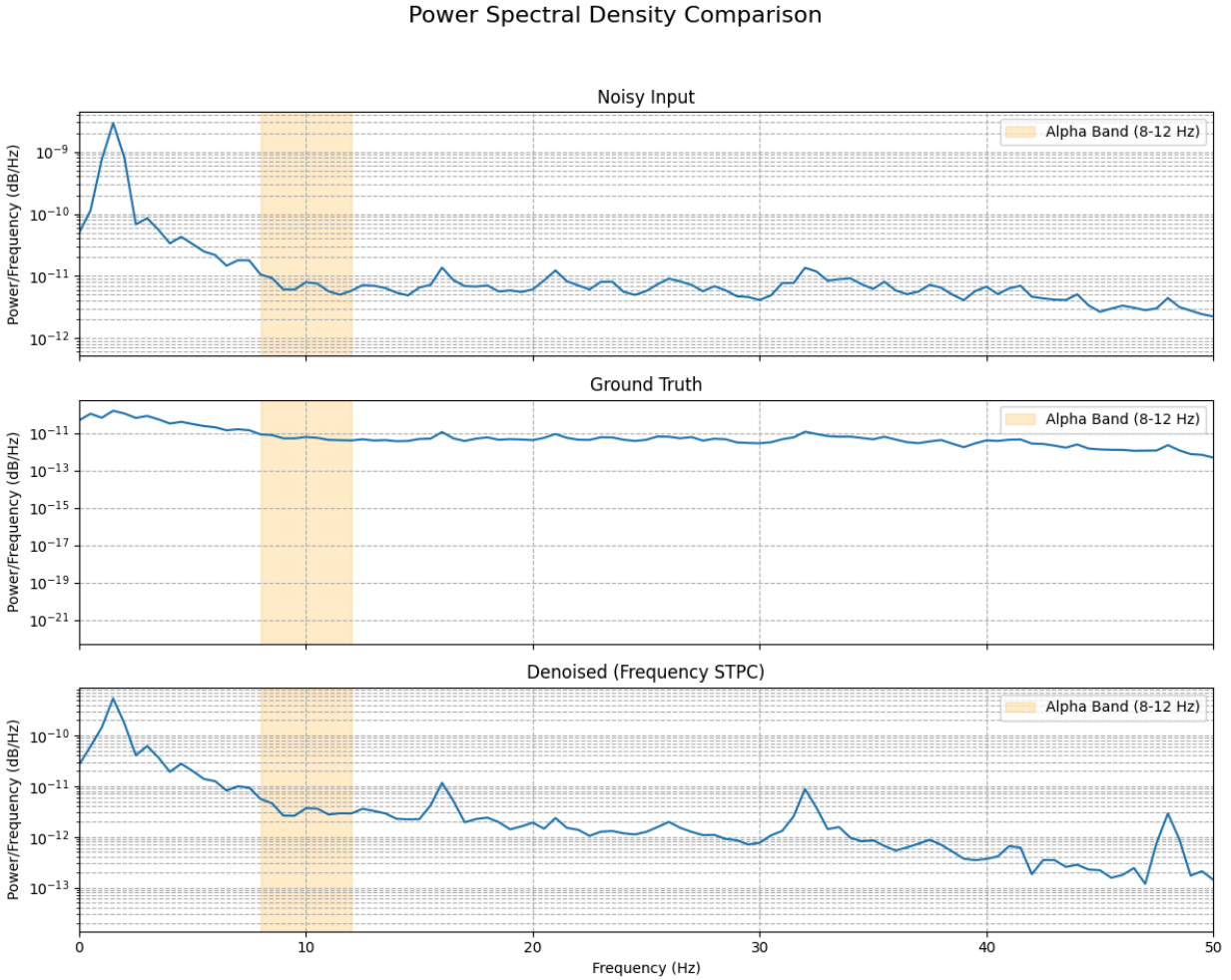




**Figure 1:** Comparison of scalp topographies. The STPC model's output (right) is a faithful reconstruction of the Ground Truth (left), unlike the Baseline L1 model (center-right).

## Phase 2: Frequency-Specific Signal Preservation

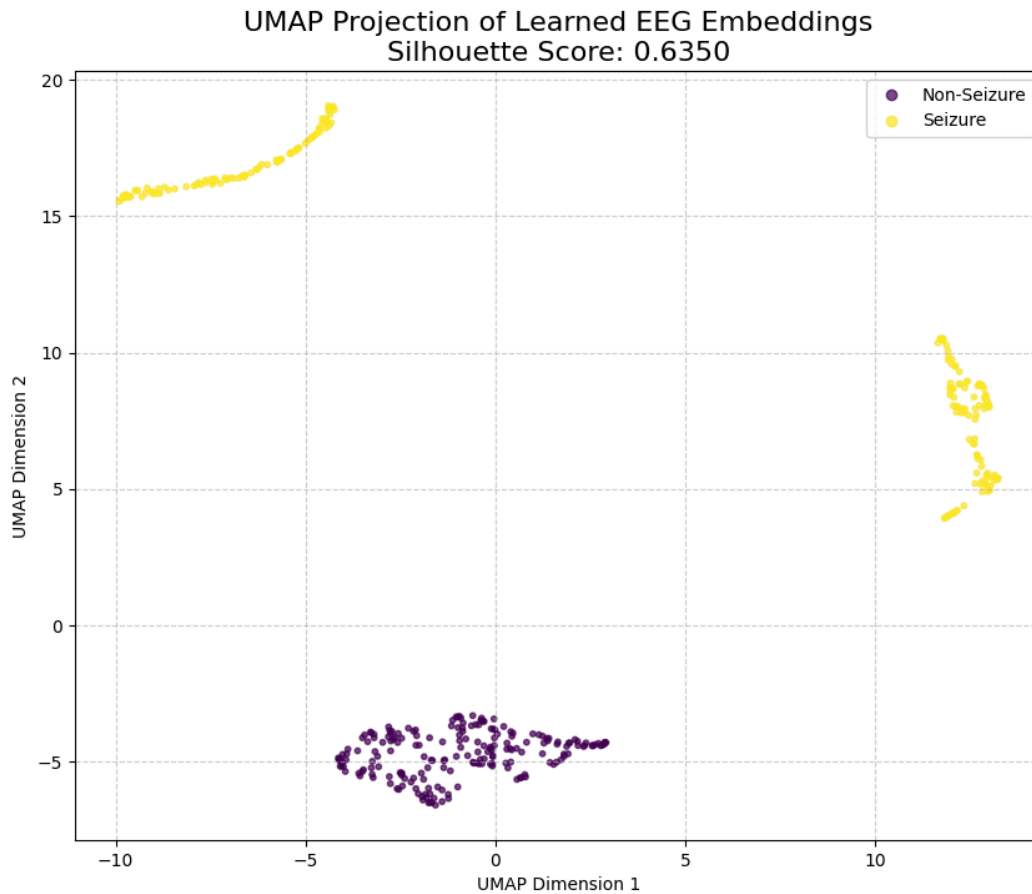
- **Objective:** To demonstrate that the STPC framework can be used to selectively preserve specific frequency bands.
- **Setup:** A new model was trained with a weighted, band-specific FFT loss, designed to heavily preserve the Alpha band (8-12 Hz) while removing powerful low- and high-frequency noise.
- **Results:** The model successfully learned this complex filtering task. As shown in Figure 2, it completely eliminated a powerful low-frequency noise artifact while precisely preserving the subtle Alpha-band rhythm present in the ground truth.



**Figure 2:** (Insert the PSD plot from Phase 2 here). Power Spectral Density comparison. The Frequency STPC model (bottom) removes the large low-frequency noise peak from the input (top) while preserving the target Alpha-band signature of the Ground Truth (middle).

### Phase 3: Unsupervised Discovery of Neural States

- **Objective:** To prove that an STPC-regularized, self-supervised model can learn to differentiate seizure from non-seizure states without labels.
- **Setup:** A U-Net was trained on a masked signal reconstruction task, using a combined L1 and Spatial STPC loss. The model was trained on a mix of all available data (seizure and non-seizure).
- **Results:** The model successfully learned to generate powerful feature embeddings. When tested on a held-out, labeled dataset of seizure and non-seizure segments, the embeddings formed clear, distinct clusters. As shown in Figure 3, the UMAP projection of these embeddings demonstrates a clean separation between the two classes, achieving a high **Silhouette Score of 0.6350**.



**Figure 3:** (Insert the UMAP scatter plot from Phase 3 here). UMAP projection of learned EEG embeddings. The model, without any labels, has learned to separate Non-Seizure (purple) and Seizure (yellow) brain states into distinct clusters, confirmed by a high Silhouette Score.

---

## 5. Discussion

The sequence of these three experiments tells a powerful story. Phase 1 established that physics-informed regularization is **necessary** for producing scientifically valid results, proving that simple metrics like RMSE are dangerously insufficient. Phase 2 demonstrated that this framework is **flexible**, allowing for the surgical preservation of specific, cognitively-relevant signal components.

Phase 3 is the culmination of this work. It demonstrates that the representations learned by an STPC-regularized model are not just "clean" but **meaningful**. The ability to spontaneously

discover and separate pathological from healthy brain states is a critical step towards data-efficient clinical models. In a field where labeled data is scarce and expensive, such powerful self-supervised feature learning is essential. This approach learns the "vocabulary" of the brain before learning to translate it.

**Limitations and Future Work:** The current framework has only been validated on a single subject. Future work must extend this to a larger cohort. The self-supervised task, while effective, is simple; more advanced contrastive learning methods could yield even richer representations. The ultimate goal is to move beyond classification and use these learned embeddings as the foundation for generative models—a true "Brain GPT" capable of predicting and synthesizing neural activity.

## 6. Conclusion

We have presented an end-to-end framework that begins with a robust, physics-informed denoising regularizer (STPC) and culminates in a powerful self-supervised learning model capable of discovering meaningful neural states from unlabeled EEG data. Our results challenge the reliance on simplistic metrics and champion a new paradigm where models are explicitly guided by the physical principles of the signals they analyze. This focus on physiological plausibility is the key to building the trustworthy, robust, and data-efficient AI systems required to unlock the future of clinical neurology and neuroscience research.