TIME SERIES SALIENCY MAPS: EXPLAINING MODELS ACROSS MULTIPLE DOMAINS

Anonymous authors

000

001

003 004

006

008 009

010 011

012

013

014

015

016

017

018

019

021

024

025

026027028

029

031

033

034

036 037

040

041

042 043

044

046

047

048

051

052

Paper under double-blind review

ABSTRACT

Traditional saliency map methods, popularized in computer vision, highlight individual points (pixels) of the input that contribute the most to the model's output. However, in time-series they offer limited insights as semantically meaningful features are often found in other domains. We introduce Cross-domain Integrated Gradients, a generalization of Integrated Gradients. Our method enables feature attributions on any domain that can be formulated as an invertible, differentiable transformation of the time domain. Crucially, our derivation extends the original Integrated Gradients into the complex domain, enabling frequency-based attributions. We provide the necessary theoretical guarantees, namely, path independence and completeness. Our approach reveals interpretable, problem-specific attributions that time-domain methods cannot capture, on three real-world tasks: wearable sensor heart rate extraction, electroencephalography-based seizure detection, and zero-shot time-series forecasting. We release an open-source Tensorflow/PyTorch library to enable plug-and-play cross-domain explainability for time-series models. These results demonstrate the ability of cross-domain integrated gradients to provide semantically meaningful insights into time series models that are impossible with traditional time-domain saliency.

1 Introduction

Saliency maps are visual tools to explain deep learning models. Popularized in computer vision, they highlight input points that contribute the most to the model's output. For images, the original input domain, pixels, aligns naturally with human perception, since neighboring pixels form coherent objects that are understood by human vision. This makes pixel-level saliency intuitive and semantically meaningful. Similarly, in natural language processing, word-level attributions can be informative, as words inherently bear semantic meaning.

In contrast, in time series, this intuition breaks down. In the time domain, groups of temporally adjacent points - the equivalent of the pixel - do not necessarily form intuitive *concepts*. Rather, such *concepts* are found in intricate interactions between points, linking them to higher-level abstractions such as oscillating frequency patterns or statistically independent formations. As a consequence, highlighting individual time points does not provide meaningful insight into the behavior of the model.

Signal processing practice has long faced this challenge, where signal interpretation generally relies on the decomposition of the original signal into structured *components*. Through transformations, the original time domain is mapped to the component domain, capturing the higher-level interaction, and linking the input to semantically meaningful concepts. The choice of decomposition and component domain depends on the nature of the signals and the task. For example, the Fourier transform decomposes the original signal into sinusoid oscillations, while the Independent Component Analysis (ICA) decomposes the signal into statistically independent components. Such transformations map the time signals into structured, semantically rich domains, providing more intuitive interpretations of the signal's contents.

Building on this insight, we argue that visual explanations of time-series models should be expressed in interpretable domains, even when the model processes time points. We empirically demonstrate that the explainability power of available saliency-based methods is limited in the time domain.

This motivates the need for saliency map tools that can visualize feature importance across multiple domains.

To address this, we develop Cross-domain Integrated Gradients, a novel method to visualize feature importance across multiple domains. Based on the principles of IG Sundararajan et al. (2017), we derive the formulas, axioms, and proofs required to apply IG across domains. We validate our method following the exact same steps as IG Sundararajan et al. (2017). We show that cross-domain IG maintains the Completeness property, hence satisfying *Sensitivity* and *Implementation Invariance*. We apply our method to real-world time-series models and applications, demonstrating that descriptive domains can be very powerful in understanding model behavior.

In this work, we introduce the following novel contributions:

- We propose a generalization of the Integrated Gradients that enables cross-domain explainability for any invertible transformation, including non-linear ones.
- We derive a generalization of the Integrated Gradients for real-valued functions with a complex domain, enabling the generation of frequency-domain saliency maps.
- We demonstrate how different domains allow for a better understanding of model behavior on time-series data.
- We release an open-source Python library, compatible with tensorflow and pytorch, for cross-domain time series explainability. The code for reproducing the results of this paper, along with the library, is available in the supplemental material. Upon acceptance, we will also include the corresponding open-source Github links.

2 RELATED WORKS

Saliency map interpretation. Saliency maps as a means of interpreting the behavior of the model have been popularized in computer vision. These methods generate an output mapping each individual input pixel to a significance score. Several methods have been proposed for this mapping. Activation-based methods, such as GradCAM Selvaraju et al. (2017) and later variations Chattopadhay et al. (2018), generate saliency based on deep layer activations. Gradient-based methods such as Integrated Gradients (IG) Sundararajan et al. (2017); Kapishnikov et al. (2021) generate significance scores by using the model's output gradients with respect to its inputs. Similarly, Layer-wise Relevance Propagation (LRP) methods Bach et al. (2015) propose rules to propagate the model output backwards by splitting the overall output among individual input features.

Time domain explainability. Saliency map methods have been applied to time series applications, either by direct application of computer vision-derived methods Jahmunah et al. (2022); Tao et al. (2024) or by developing dedicated time series saliency approaches Queen et al. (2023); Liu et al. (2024). To streamline comparisons between time-domain interpretability, Ismail et al. (2020) proposed an extensive synthetic, multi-channel benchmark. In all cases, these approaches focus on identifying significant regions of the time-domain input which contribute the most to the model's output. Such regions of interest are events that trigger the model's output.

Cross domain interpretability. The current time domain saliency methods have limitations, as highlighted time points do not always explain the underlying mechanisms Theissler et al. (2022). Furthermore, Chung et al. (2024) demonstrate that such methods are not robust to frequency perturbations. These limitations diminish the explanatory power of the generated saliency map. To address this issue, they proposed a perturbation method in the time-frequency domain, attributing the model output to time-frequency components. However, frequency perturbations can strongly affect model performance and, therefore, explainability due to out-of-distribution effects Sundararajan et al. (2017). Similarly, Vielhaben et al. Vielhaben et al. (2024) proposed the *virtual inspection layer* placed after the model input to transform the saliency map of the time domain to the frequency and time frequency domains, proposing dedicated relevance propagation rules for the frequency transform.

Saliency map evaluation. Evaluating saliency maps is not a trivial task. A major challenge lies in disentangling saliency map errors from model errors Kim et al. (2021); Akhavan Rahnama (2023), complicating validation by comparison with ground truth saliency. Sundararajan et al. (2017) propose

solving this by relying on a set of desirable axioms, bypassing necessity for empirical evaluations. Validation based on insert / deletion is another approach Hama et al. (2023); Ismail et al. (2020). These methods empirically evaluate the effect of removing/retaining the most important input features, reinforcing trust in the saliency map method under examination.

Despite progress in time-series saliency, existing methods (i) operate solely in the time domain, (ii) rely on perturbation-based attributions only in the frequency domain, or (iii) require transform-specific hand-crafted relevance-propagation rules valid only in the frequency domain. In contrast, our work provides a principled generalization of Integrated Gradients that supports any invertible, differentiable transform, including complex-valued domains, while preserving axiomatic properties and enabling semantically meaningful attributions across diverse time series applications.

3 PRELIMINARIES

3.1 PROBLEM STATEMENT AND MOTIVATION

We consider a function $f:\mathcal{D}_s\to\mathbb{R}$ representing a deep learning model. The input $x\in\mathcal{D}_s$ is constructed from a continuous-time signal $x(t)\in\mathbb{R}$ after discretizing it at a sampling frequency f_s [Hz] and considering a window of length L seconds: $x=[x_0,...,x_{n-1}], n=f_s\cdot L$. Now consider a transform $T:\mathcal{D}_S\to\mathcal{D}_T$ that maps the original time domain to a semantically rich explanation target domain \mathcal{D}_T . Our task is to construct an informative saliency map that assigns a significance score to each characteristic $z_i=T(x)_i$ in the explanation domain.

Saliency maps developed in computer vision applications, and in particular IG, provide explanations in the same domain as the model's input, that is, $\mathcal{D}_T = \mathcal{D}_S$. Applying these methods to time-series models results in maps expressed in the time domain.

Proposition 1. The time domain is not always informative in explaining f.

We motivate Proposition 1 through a synthetic example. We provide additional real-world examples in Section 5 after formally defining our method.

3.2 TIME DOMAIN EXPLANATION LIMITATIONS

Consider that the input x is sampled from signals $x(t) = cos(2\pi\xi t + \phi)$. In this setup, there are two classes of samples depending on the oscillating frequency ξ :

$$y = \begin{cases} 1, & \xi \sim \mathcal{N}(1.0, 0.5) \\ 2, & \xi \sim \mathcal{N}(4.0, 0.5) \end{cases}$$
 (1)

We design a classifier f to distinguish between these two classes. We opt to manually construct f so that we have full mechanistic understanding of its inner workings. We choose a CNN architecture composed of a single convolutional layer with two channels followed by a ReLU activation and global average pooling $f(x) = AvgPool\left(ReLU(\boldsymbol{w}*\boldsymbol{x})\right)$. The kernel of the first channel is a low-pass filter (cutoff at 2.5Hz), while the second channel kernel is a high-pass filter with the same cutoff (see Figure 1).

Ideally, the model should be fully explained by describing its inner mechanism. In this particular scenario, we have designed f for this purpose, and hence a formal detailed explanation is available.

Mechanistic Interpretation 1. Convolutional channel i allows only frequencies of class i to pass through the output; otherwise, the channel's output is almost zero, not activating. The ReLU and Average Pooling mechanism extract the amplitude of the signal Kechris et al. (2024a). Hence, the channel i of the model output is only active when samples from class i are processed, leading to the correct classification of the input.

That depth in model understanding is not easily available in larger models, which have been learned from samples. Hence, saliency maps are often used as a proxy. We provide IG explanations of the model f for samples from both classes expressed in the time and frequency domains (Figure 1). Although time points are periodically highlighted as *more important*, it is not exactly clear how this input tilts the model towards producing its output.

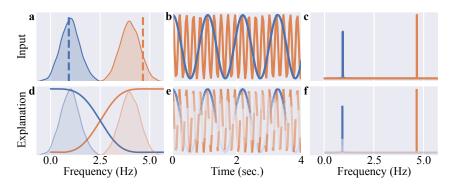


Figure 1: Mechanistic interpretation along with Time and Frequency domain saliency maps. (a) Distributions of the main frequency, ξ , for classes one and two. For producing the saliency maps, we sample one input for each class (vertical dashed lines). (b) The sampled inputs presented in the time and (c) frequency domains. (d) Illustration of the Mechanistic Interpretation1. We plot the frequency response for the first and second channels of the CNN. The sample distributions (a) are also overlayed. (e) Saliency maps expressed in the time and (f) frequency domains.

In contrast, a saliency map expressed in the frequency domain, which we introduce in Section 4, highlights the frequency components that contribute to the final output: for the samples of class one, only the 1 Hz component contributes to the model's output, and accordingly, for class two, the 4 Hz component. Here, this saliency map is much more interpretable. It provides useful information and better aligns with the mechanistic understanding (Mechanistic Interpretation 1) of this model. In Section 4, we show analytically that the frequency-expressed IG, for the data distribution and model of this example, is directly linked to its mechanistic explanation.

3.3 Integrated Gradients

To explain the output of a model f on an input x with a baseline $\hat{x} \in \mathbb{R}^n$, IG generates a saliency map as Sundararajan et al. (2017):

$$IG_i(\boldsymbol{x}) = (x_i - \hat{x}_i) \int_0^1 \frac{\partial f}{\partial x_i} \bigg|_{\boldsymbol{x}' + t \cdot (\boldsymbol{x} - \hat{\boldsymbol{x}})} dt$$
 (2)

with each element $IG_i(x)$ of the map corresponding to the significance of the input feature x_i : saliency is expressed in the same domain as the input. The IG definition relies on two key points from the theory of integrals over differential forms: the line integral definition and Stokes' theorem.

Line integral definition. The IG can be derived from the definition of the integral of the differential form df along the line $\gamma(t) = \hat{x} + t(x - \hat{x})$:

$$\int_{\gamma} df = \int \gamma^* df = \int_0^1 \sum_{i=0}^N \frac{\partial f}{\partial x_i} \gamma_i'(t) dt = \sum_{i=0}^N \int_0^1 \frac{\partial f}{\partial x_i} \gamma_i'(t) dt = \sum_{i=0}^N (x_i - \hat{x}_i) \int_0^1 \frac{\partial f}{\partial x_i} dt \quad (3)$$

where $\gamma^* df$ is the pullback of df by γ : $\gamma^* df = \sum_{i=0}^N \frac{\partial f}{\partial x_i} \gamma_i'(t) dt$ Do Carmo (1998). Each individual element of the IG map $IG_i(x)$ corresponds to each element of the last sum of eq. 3.

Stoke's Theorem. The *Completeness* axiom of the IG Sundararajan et al. (2017): $f(x) - f(\hat{x}) = \sum IG_i$ is a consequence of the Stokes' Theorem for the case of integral of 1-form: $\int_{\gamma} df = \int_{\partial \gamma} f = f(x) - f(\hat{x})$, which guarantees path independence: the value of the integral is only dependent on the first and last points of the path, not the path itself.

3.4 SALIENCY MAPS EVALUATIONS

Saliency map evaluation is challenging (Section 2), therefore we adopt a broad, complementary validation protocol that triangulates evidence from theory, controlled experiments, qualitative sanity checks, and dataset-level stress tests:

- 1. **Axiomatic soundness.** We show that Cross-domain IG maintains the Completeness property, hence satisfying *Sensitivity* and *Implementation Invariance* Sundararajan et al. (2017).
- 2. **Mechanistic alignment.** Based on the example in Section 3.2, we theoretically show that cross-domain IG can align with the model's internal mechanisms when the target domain is appropriate (Section 4.2).
- 3. **Qualitative applications.** We show representative examples, Section 5, demonstrating the full Cross-Domain IG workflow and how it can uncover data/model insights.
- Quantitative sufficiency/necessity. We run insertion-deletion tests on real-world time-series datasets.

4 METHODS

In this section, we define Cross-Domain IG (Section 4.1), and derive it based on the IG principles from Section 3.3. We then analyse it in the complex frequency domain using a simple yet representative convolutional network, highlighting its relation to the network's properties (Section 4.2). This analysis also provides theoretical grounding for the connection between frequency-domain IG and the Mechanistic Interpretation discussed in Section 3.2. Finally, we detail our method's implementation.

4.1 Cross-domain IG derivation

Let $f: \mathcal{D}_s \to \mathbb{R}$ a deep neural network, operating on a domain $\mathcal{D}_s \subseteq \mathbb{R}^n$. Also, denote $x, \hat{x} \in \mathcal{D}_s$ the input and baseline samples, respectively, as defined by the IG method. We introduce an invertible, differentiable transformation $T: \mathcal{D}_S \to \mathcal{D}_T$ and its inverse T^{-1} , also differentiable, with z = T(x) and $z = T^{-1}(z)$ and z =

Definition 4.1 (Cross-domain Integrated Gradients). Given a model $f: \mathcal{D}_s \to \mathbb{R}$, a transform $T: \mathcal{D}_S \to \mathcal{D}_T$ and its inverse T^{-1} , input and baseline samples $x, \hat{x} \in \mathcal{D}_s$ and $\gamma(t)$ the line from z = T(x) to $\hat{z} = T(\hat{x})$ the Cross-Domain IG is defined as:

$$IG_i^{\mathcal{D}_T}(\boldsymbol{z}) = 2 \int_0^1 \operatorname{Re} \left\{ \frac{\partial (f \circ T^{-1})}{\partial z_i} \Big|_{\boldsymbol{\gamma}(t)} \cdot (z_i - \hat{z}_i) \right\} dt \tag{4}$$

Note that the original IG, eq. 2, and $IG^{\mathcal{D}_T}$ explain the exact same functionality since f(x) and $(f \circ T^{-1})(z)$ are equivalent. However, their output saliency maps are expressed in different domains. We now derive Definition 4.1 from first principles of the original IG method, Section 3.3.

Derivation sketch. The original IG is only defined for real inputs. To enable complex-valued transformations, such as the Fourier transform, we extend IG for real-valued functions g with complex inputs z, referred to as Complex IG. Our derivation builds on the two key points in Section 3.1:

- 1. **Line integral definition.** We begin our derivation by defining a function u that is equivalent to g(z). Just like in the case of real inputs, eq. 2, we elaborate on the line integral $\int_{\gamma} du$. The end goal is to end up with a sum of integrals $\sum_{i} \int ...dt$ similar to eq. 3. In the final step, each IG element is defined as the corresponding integral term of the final sum, $\int ...dt$.
- 2. **Stokes' Theorem.** We define u and derive complex IG to ensure path independence and satisfy the *Completeness axiom*, which may fail for functions of several complex variables Lebl (2019). To this end, we first state and prove Lemma 4.1 as an intermediate result. Using Lemma 4.1, we then derive Definition 4.1 using Wirtinger calculus.

Lemma 4.1. Let $g: \mathbb{C}^n \to \mathbb{R}$, z = p + jq, with $p, q \in \mathbb{R}^N$, $\gamma(t) = \hat{z} + t(z - \hat{z})$, $t \in [0, 1]$ the line from the baseline point \hat{z} to the input point z and $n(t) = \text{Re}\{\gamma(t)\}$ and $m(t) = \text{Im}\{\gamma(t)\}$, $n(t), m(t) \in \mathbb{R}^n$. Then the IG of q in z is given by:

$$IG_i^{\mathbb{C}^n}(z) = \int_0^1 \left(\frac{\partial g}{\partial p_i} n_i'(t) + \frac{\partial g}{\partial q_i} m_i'(t) \right) dt$$
 (5)

A detailed proof of Lemma 4.1 can be found in Appendix B. From Lemma 4.1, and considering $g(z) = f\left(T^{-1}(z)\right)$ and the complex differential form Range (1998) $dg = \partial g + \overline{\partial} g$ we can write the complex integrated gradient definition as:

$$IG_i^{\mathbb{C}^n} = 2 \int_0^1 \operatorname{Re} \left\{ \frac{\partial g}{\partial z_i} \gamma_i'(t) \right\} dt \tag{6}$$

The complete derivation can be found in Appendix C. Notice that Cross-domain IG maintains the Completeness property since $\int_{\gamma} du = u(\boldsymbol{a}(1)) - u(\boldsymbol{a}(0)) = g(\boldsymbol{z}) - g(\hat{\boldsymbol{z}}) = f(\boldsymbol{x}) - f(\hat{\boldsymbol{x}})$, where $u : \mathbb{R}^{2n} \to \mathbb{R}$ s.t. $g(\boldsymbol{p} + j\boldsymbol{q}) = u([\boldsymbol{p}, \boldsymbol{q}])$ and $\boldsymbol{a} = [\boldsymbol{n}, \boldsymbol{m}]$.

Remark 1. Although definition 4.1 defines a linear path of integration, in our derivation, eq. 6, the path of integration is a general curve $\gamma(t)$. This enables incorporating into cross-domain IG alternative integration paths/methods to reduce sensitivity to noise Yang et al. (2023); Kapishnikov et al. (2021).

Cross-Domain IG for real-valued inputs. If g processes real-valued inputs, then eq. 6 is equivalent to eq. 2: since g(z) = g(p+j0), $\partial g/\partial q = 0$, $\partial g/\partial z = (1/2)\partial g/\partial p$. Thus, if $\mathcal{D}_T \subseteq \mathbb{R}^n$ the cross-domain IG can equivalently be expressed as $IG_i^{\mathcal{D}_T}(z) = (z_i - \hat{z}_i) \int \frac{\partial (f \circ T^{-1})}{\partial z_i} dt$.

Remark 2. In IG Sundararajan et al. (2017), the baseline \hat{x} is defined as the point without information about the original model inference. The authors argued that most deep networks admit such a neutral input. For cross-domain IG, if \hat{x} exists, and T is invertible, then \hat{z} is trivially defined. Crucially, cross-domain IG enables baselines that were not easily defined, e.g., filtering specific components from \hat{x} to form \hat{z} .

4.2 COMPLEX IG ON A SIMPLE MODEL

Adebayo et al. (2018) analytically study a minimal single-layer convolutional network, demonstrating that IG can collapse into an *edge detector*, producing misleading saliency maps. Although this exposes a failure mode of the IG in the input domain, we show that Complex-IG faithfully reflects the inner mechanisms of a simple convolutional network in the frequency domain. In direct parallel, we derive a closed-form link between the complex IG saliency map of a CNN and the frequency response of its filters. Building on the example in Section 3.2, we work on a simple CNN and prove that Complex-IG highlights each filter's gain at its corresponding input frequency.

Let f be a convolutional neural network composed of a single convolutional layer (1 channel) followed by a ReLU operation and Global Average Pooling: $f(\boldsymbol{x}) = AvgPool \ (ReLU(\boldsymbol{w}*\boldsymbol{x}))$. We begin with the case in which f processes windows sampled from single-component sinusoidal signals $x(t) = a_j \cdot cos(2\pi\xi_j t + \phi), \ a_j > 0$. Then, the output $f(\boldsymbol{x})$ is Kechris et al. (2024a): $f(\boldsymbol{x}) = \frac{a_jb_j}{\pi}$, with b_i the amplification of the filter \boldsymbol{w} at frequency ξ_iHz : $b_i = \|\sum_n w_n e^{-2\pi\xi_i n}\|$. We employ the Complex IG method on f with baseline input $\hat{\boldsymbol{x}} = \boldsymbol{0}, \ f(\boldsymbol{0}) = 0$. This yields $IG_i^{\mathbb{C}^n} = 0, \ \forall i \neq j$ and $\sum_i IG_i^{\mathbb{C}^n} = f(\boldsymbol{x}) - f(\hat{\boldsymbol{x}})$. Thus,

$$IG_j^{\mathbb{C}^n} = f(\boldsymbol{x}) = \frac{a_j b_j}{\pi} \tag{7}$$

This links $IG_j^{\mathbb{C}^n}$ to the output frequency content a_jb_j and, by extension, to the convolutional filter's frequency response. An example for the model of Section 3.2 is presented in Figure 5 (Appendix E).

4.3 IMPLEMENTATION

Autograd (pytorch / tensorflow) allows for automatic differentiation with complex variables using Wirtinger calculus Kreutz-Delgado (2009). Thus, the complex IG can be directly approximated

by autograd, using Definition 4.1 or Lemma 4.1, with the detail that Autograd (in both libraries) calculates the conjugate of the complex partial derivative. For the integral calculation, we use a summation approximation similar to Sundararajan et al. (2017). The algorithms for estimating cross-domain IG for the case of $\mathcal{D}_T \subseteq \mathbb{R}^n$ and the two implementations on $\mathcal{D}_T \subseteq \mathbb{C}^n$ (Lemma 4.1 and Definition 4.1) are presented in Algorithms 1 and 2, 3 in the appendix, respectively.

Remark 3. The numerical approximation of the integral in Definition 4.1 requires multiple differentiations, which increases computational overhead. Although the original IG also suffers from similar overhead, our method requires an additional step due to the inverse transform step (see line 9 in Algorithm 3 in the Appendix).

5 APPLICATIONS

We deploy cross-domain IG in a range of time series applications and models. We selected applications on all three main time-series tasks: regression (section 5.1), classification (section 5.2) and forecasting (Section 5.3). In all three cases, the models are trained to infer on inputs in the time domain. For each application, first, we study the properties of the input signal from a signal processing perspective. We then define an interpretability task: what do we want to learn about our model's behavior through a saliency map? Based on this domain knowledge and interpretability task, we select an appropriate explanation space yielding semantically meaningful saliency maps. We conclude each example with a remark on actionable insights based on cross-domain attributions. Time-Domain IG attributions of these examples can be found in the Appendix G, and additional examples in Appendix H. We also perform feature insertion/deletion evaluation in Appendix F.

5.1 HEART RATE EXTRACTION FROM PHYSIOLOGICAL SIGNALS

We use the KID-PPG Kechris et al. (2024b), a deep convolutional model with attention, to extract heart rate (HR) from photoplethysmography (PPG) signals collected from a wrist-worn wearable device. We use signals from the PPGDalia dataset Reiss et al. (2019). For a time window small enough for the HR frequency, ξ_{hr} , to be considered constant, a clean PPG signal can be modeled as Kechris et al. (2024b): $x(t) = a_1 cos(2\pi \cdot \xi_{hr} \cdot t + \phi) + a_2 cos(2 \cdot \pi(2\xi_{hr}) \cdot t + \phi)$, with $a_1 > a_2$. However, external signals are also usually present in PPG recordings Reiss et al. (2019); Kechris et al. (2024b). These *interferences* are not created by the heart and are preventing the model from making accurate HR inferences.

Remark 4. KID-PPG processes PPG signals that contain both heart-related components and external inference. A trustworthy model should base the inferred hear-rate on heart-related signals only, filtering out all other sources of noise.

Interpretability task. Given a PPG sample and KID-PPG's HR inference, determine whether the model is focusing on heart-related information or external interference.

Problem-specific transformation. Since our understanding of this application is mostly frequency-based, we have selected the frequency domain using the Fourier transform as the explanation target domain. Hence, the frequency-domain IG highlights individual frequencies as important to the final model inference. This allows us to investigate whether the HR inference is produced from components related to the heart or external interference.

An illustration of two PPG inputs and the corresponding frequency-domain IGs are presented in Figure 2. The frequency IG identifies samples in which the model infers heart rate from external interference, thus limiting the reliability of the model's output.

Remark 5. Frequency-domain IG highlights whether KID-PPG inference is **trustworthy** (based on heart oscillations) or **spurious** (based on motion-induced artifacts).

5.2 ELECTROENCEPHALOGRAPHY-BASED EPILEPTIC SEIZURE DETECTION

We use the zhu-transformer Zhu and Wang (2023) which performs seizure detection on scalpelectroencephalography (EEG). We analyze a recording from the Physionet Siena Scalp EEG Database

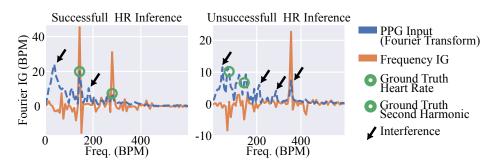


Figure 2: **Frequency-domain IG on heart rate inference model.** The **PPG** signal includes components from the **heart rate** and other components attributed to external interference (\rightarrow) , e.g. motion. **Left:** Sample with a small inference error 0.93 beats-per-minute (BPM). The **IG** highlights the two heart components located at hr and $2 \cdot hr$ (second harmonic), with more weight given to the actual heart rate frequency. **Right: PPG** sample with high inference error (26.78 BPM). **IG** coefficients highlight frequency components which are not related to the heart.

v1.0.0 Detti (2020); Detti et al. (2020); Goldberger et al. (2000). In EEG a single channel captures the electrical activity of multiple *sources*: e.g., epileptic activity, muscle interference, or electrical noise.

Remark 6. A seizure classification model processes the aggregated activity of all sources in the EEG. The model should isolate only the epileptic activity, filtering all others, to reach a trustworthy inference.

Interpretability task. Given an EEG recording and the corresponding zhu-transformer seizure classification, we want to identify the *sources* on which the model based its inference.

Problem-specific transformation. We chose Independent Component Analysis Lee and Lee (1998) (ICA) as our transform of choice. ICA isolates the activity of each individual source to a source-specific channel (Independent Component), assuming statistical independence between the sources. This allows the ICA-domain IG to produce attributions for each individual isolated source, therefore providing insights on our interpretability task (Figure 3).

Remark 7. ICA-IG highlights whether zhu-transformer inference is based on known components of epileptic seizure activity or other components irrelevant to the seizure - further reinforcing trust in the model decision.

5.3 FOUNDATION MODEL TIME SERIES FORECASTING

We use TimesFM Das et al. (2024) time-series foundation model to explain forecasting outputs. We perform zero-shot forecasting, without any fine-tuning, on a time series with exponential trend and seasonal components (Figure 4).

Remark 8. A time-series forecasting model should model equally successfully both the trend and the season to reach a low-error long-horizon forecast.

Interpretability task. Given a time-series input and the TimesFM forecast, determine if the trend or the season is more difficult to model in the long-horizon forecast setting.

Task-specific transform. To isolate the relevant *concepts* we chose Seasonal-Trend decomposition using LOESS (STL) Cleveland et al. (1990) to decompose the input time series into trend and seasonality components.

This attribution domain allows us to study the model's behavior for long-term forecasting horizons where the forecast error increases: the model underestimates the overall trend, while the seasonal component estimation presents a smaller error.

Remark 9. Seasonal-Trend IG reveals that TimesFM underweights the trend, degrading long-horizon forecasts. This offers concrete insights to improve model behaviour.

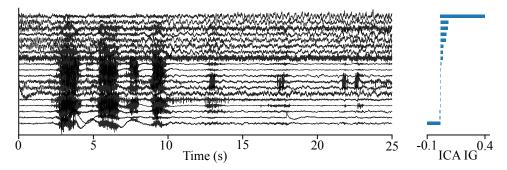


Figure 3: ICA-domain IG on seizure detection model. The ICA components are sorted from the component with the highest IG significance (top) to the lowest (bottom). Left: 19 output channels calculated from ICA on the original EEG channels. The first channel contains the majority of the epileptic activity, which is visible as an evolving pattern of spike-and-wave discharges at $\sim 4.5~\rm Hz$. Some epileptic activity can also be found in the second channel. Significant muscle artifacts are isolated in the 9th-19th channels between 4 and 10 seconds. Right: IG saliency map calculated on the channel components. The map identifies the first channel as the most significant channel in detecting this sample as epileptic. Some significance, although much less, is also given to the next four channels. The channels corresponding to interference components do not get any significance in the output of the classifier. The last channel *tends to tilt* the classifier towards a non-epileptic output.

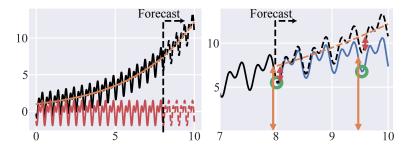


Figure 4: Seasonal-Trend IG on time series foundation model. Left: Input time series decomposed via STL into trend and seasonality. Right: Zero-shot forecasting using TimesFM with Seasonal-Trend IG. For a small horizon, one step ahead prediction (first circle), TimesFM forecasts accurately. Of output, 7.5 units are attributed to trend (\updownarrow), aligning with ground truth (dashed orange) and similarly -1.96 units to seasonality (\updownarrow). For a longer horizon (second circle) the forecast absolute error rises from 0.2 to 2.14. Most of it stems from the model's underestimation of the trend (21% relative error), while the seasonal effect is correctly captured by the model (5.1% relative error).

6 CONCLUSIONS

We introduced a novel generalization of the Integrated Gradients method, which enables saliency map generation in any invertible differentiable transform domain, including complex spaces. As transforms capture high-level interactions between input points, our methods enhance model explainability, especially in time-series data where individual time-point features are often uninformative. We demonstrated versatility of Cross-domain Integrated Gradients, applying it on a diverse set of time-series tasks, model architectures and explanation target domains. Fields where time signals are extensively used, such as healthcare, finance and environmental monitoring, could benefit from domain-specific saliency maps. In particular, with the recent rise of time-series foundation models, our method provides a strong investigation tool for inspecting model behavior. We release an open-source library to enable broader adoption of cross-domain time-series explainability.

7 ETHICAL STATEMENT

Risks may arise if the selected explanation target domain is not appropriate or saliency maps are over-interpreted. It is important to note that the saliency map provides only feature significance scores. Interpreting these scores requires domain expertise. We encourage a holistic interpretation approach to integrating domain knowledge with cross-domain saliency maps. We also caution that this method alone cannot function as a definitive proof of the behavior of the model. Responsible usage of the method should take into consideration model, data and transformation limitations, especially in high-stakes settings, such as in healthcare. We elaborate on the limitations of our method in Appendix K

REFERENCES

- Julius Adebayo, Justin Gilmer, Michael Muelly, Ian Goodfellow, Moritz Hardt, and Been Kim. Sanity checks for saliency maps. *Advances in neural information processing systems*, 31, 2018.
- Amir Hossein Akhavan Rahnama. The blame problem in evaluating local explanations and how to tackle it. In *European Conference on Artificial Intelligence*, pages 66–86. Springer, 2023.
- Sebastian Bach, Alexander Binder, Grégoire Montavon, Frederick Klauschen, Klaus-Robert Müller, and Wojciech Samek. On pixel-wise explanations for non-linear classifier decisions by layer-wise relevance propagation. *PloS one*, 10(7):e0130140, 2015.
- Aditya Chattopadhay, Anirban Sarkar, Prantik Howlader, and Vineeth N Balasubramanian. Grad-cam++: Generalized gradient-based visual explanations for deep convolutional networks. In 2018 IEEE winter conference on applications of computer vision (WACV), pages 839–847. IEEE, 2018.
- Hyunseung Chung, Sumin Jo, Yeonsu Kwon, and Edward Choi. Time is not enough: Time-frequency based explanation for time-series black-box models. In *Proceedings of the 33rd ACM International Conference on Information and Knowledge Management*, pages 394–403, 2024.
- Robert B Cleveland, William S Cleveland, Jean E McRae, Irma Terpenning, et al. Stl: A seasonal-trend decomposition. *J. off. Stat*, 6(1):3–73, 1990.
- Abhimanyu Das, Weihao Kong, Rajat Sen, and Yichen Zhou. A decoder-only foundation model for time-series forecasting. In *Forty-first International Conference on Machine Learning*, 2024.
- Paolo Detti. Siena scalp eeg database v1.0.0. Physionet, 2020.
- Paolo Detti, Giampaolo Vatti, and Garazi Zabalo Manrique de Lara. EEG synchronization analysis for seizure prediction: A study on data of noninvasive recordings. *Processes*, 8:846, 2020.
- Manfredo P Do Carmo. *Differential forms and applications*. Springer Science & Business Media, 1998.
- Ary L. Goldberger, L. A. Amaral, L. Glass, J. M. Hausdorff, P. C. Ivanov, R. G. Mark, J. E. Mietus, G. B. Moody, C. K. Peng, and H. E. Stanley. Physiobank, physiotoolkit, and physionet: components of a new research resource for complex physiologic signals. *Circulation*, 2000.
- Naofumi Hama, Masayoshi Mase, and Art B Owen. Deletion and insertion tests in regression models. *Journal of Machine Learning Research*, 24(290):1–38, 2023.
- Aapo Hyvärinen, Juha Karhunen, and Erkki Oja. *Independent Component Analysis*. Wiley, 1 edition, May 2001. ISBN 9780471405405 9780471221319. doi: 10.1002/0471221317. URL https://onlinelibrary.wiley.com/doi/book/10.1002/0471221317.
- Aya Abdelsalam Ismail, Mohamed Gunady, Hector Corrada Bravo, and Soheil Feizi. Benchmarking deep learning interpretability in time series predictions. *Advances in neural information processing systems*, 33:6441–6452, 2020.
- Vicneswary Jahmunah, Eddie YK Ng, Ru-San Tan, Shu Lih Oh, and U Rajendra Acharya. Explainable detection of myocardial infarction using deep learning models with grad-cam technique on ecg signals. *Computers in Biology and Medicine*, 146:105550, 2022.

- Andrei Kapishnikov, Subhashini Venugopalan, Besim Avci, Ben Wedin, Michael Terry, and Tolga Bolukbasi. Guided integrated gradients: An adaptive path method for removing noise. In Proceedings of the IEEE/CVF conference on computer vision and pattern recognition, pages 5050–5058, 2021.
- Christodoulos Kechris, Jonathan Dan, Jose Miranda, and David Atienza. Dc is all you need: describing relu from a signal processing standpoint. *arXiv preprint arXiv:2407.16556*, 2024a.
 - Christodoulos Kechris, Jonathan Dan, Jose Miranda, and David Atienza. Kid-ppg: Knowledge informed deep learning for extracting heart rate from a smartwatch. *IEEE Transactions on Biomedical Engineering*, 2024b.
 - Joon Sik Kim, Gregory Plumb, and Ameet Talwalkar. Sanity simulations for saliency methods. *arXiv* preprint arXiv:2105.06506, 2021.
 - Marius Klug and Klaus Gramann. Identifying key factors for improving ICA-based decomposition of EEG data in mobile and stationary experiments. *The European Journal of Neuroscience*, 54(12): 8406–8420, December 2021. ISSN 1460-9568. doi: 10.1111/ejn.14992.
 - Ken Kreutz-Delgado. The complex gradient operator and the cr-calculus. *arXiv preprint arXiv:0906.4835*, 2009.
 - Jiri Lebl. Tasty bits of several complex variables. Lulu. com, 2019.
 - Te-Won Lee and Te-Won Lee. *Independent component analysis*. Springer, 1998.
 - Zichuan Liu, Tianchun Wang, Jimeng Shi, Xu Zheng, Zhuomin Chen, Lei Song, Wenqian Dong, Jayantha Obeysekera, Farhad Shirani, and Dongsheng Luo. Timex++: Learning time-series explanations with information bottleneck. *arXiv preprint arXiv:2405.09308*, 2024.
 - Owen Queen, Tom Hartvigsen, Teddy Koker, Huan He, Theodoros Tsiligkaridis, and Marinka Zitnik. Encoding time-series explanations through self-supervised model behavior consistency. *Advances in Neural Information Processing Systems*, 36:32129–32159, 2023.
 - R Michael Range. *Holomorphic functions and integral representations in several complex variables*, volume 108. Springer Science & Business Media, 1998.
 - Attila Reiss, Ina Indlekofer, Philip Schmidt, and Kristof Van Laerhoven. Deep ppg: Large-scale heart rate estimation with convolutional neural networks. *Sensors*, 19(14):3079, 2019.
 - Ramprasaath R Selvaraju, Michael Cogswell, Abhishek Das, Ramakrishna Vedantam, Devi Parikh, and Dhruv Batra. Grad-cam: Visual explanations from deep networks via gradient-based localization. In *Proceedings of the IEEE international conference on computer vision*, pages 618–626, 2017.
 - Mukund Sundararajan, Ankur Taly, and Qiqi Yan. Axiomatic attribution for deep networks. In *International conference on machine learning*, pages 3319–3328. PMLR, 2017.
 - Rui Tao, Lin Wang, Yingnan Xiong, and Yu-Rong Zeng. Im-ecg: An interpretable framework for arrhythmia detection using multi-lead ecg. *Expert Systems with Applications*, 237:121497, 2024.
 - Andreas Theissler, Francesco Spinnato, Udo Schlegel, and Riccardo Guidotti. Explainable ai for time series classification: a review, taxonomy and research directions. *Ieee Access*, 10:100700–100724, 2022.
 - Johanna Vielhaben, Sebastian Lapuschkin, Grégoire Montavon, and Wojciech Samek. Explainable ai for time series via virtual inspection layers. *Pattern Recognition*, 150:110309, 2024.
- Irene Winkler, Stefan Haufe, and Michael Tangermann. Automatic Classification of Artifactual ICA-Components for Artifact Removal in EEG Signals. *Behavioral and Brain Functions*, 7(1):30, August 2011. ISSN 1744-9081. doi: 10.1186/1744-9081-7-30. URL https://doi.org/10.1186/1744-9081-7-30.

Ruo Yang, Binghui Wang, and Mustafa Bilgic. Idgi: A framework to eliminate explanation noise from integrated gradients. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 23725–23734, 2023.

Yuanda Zhu and May D Wang. Automated seizure detection using transformer models on multichannel eegs. In 2023 IEEE EMBS International Conference on Biomedical and Health Informatics (BHI), pages 1–6. IEEE, 2023.

A CROSS-DOMAIN IG ALGORITHMS

649650651652

671672673

674

701

648

```
Algorithm 1 Real Target Domain IG
```

```
653
           Input: f(\cdot), x, \hat{x}, n_{iter}
654
           Output: IG
655
             1: i \leftarrow 1
656
             2: sum \leftarrow 0
657
             3: tape \leftarrow tensorflow.GradientTape()
658
             4: X' \leftarrow T(\hat{x})
659
             5: for i \leq n_{iter} do
660
             6:
                     z \leftarrow T(x)
             7:
                      z \leftarrow \hat{z} + (z - \hat{z}) \cdot i / n_{iter}
661
                      tape.watch(z)
662
                      x_{rec} \leftarrow T^{-1}(z)
             9:
663
                      y \leftarrow f(x_{rec})
            10:
664
                      dy \leftarrow tape.gradient(y, z)
            11:
665
                      sum \leftarrow sum + dy
            12:
666
            13:
                      i \leftarrow i + 1
667
            14: end for
668
            15: sum \leftarrow sum/n_{iter}
669
            16: IG = (z - \hat{z}) \cdot sum
670
```

Algorithm 2 Complex Target Domain IG

```
675
            Input: f(\cdot), x, \hat{x}, n_{iter}
676
            Output: IG
677
              1: i \leftarrow 1
678
              2: sum real \leftarrow 0
              3: sum\_imag \leftarrow 0
679
              4: tape\_real \leftarrow tensorflow.GradientTape()
680
              5: tape\_imag \leftarrow tensorflow.GradientTape()
681
              6: \hat{z} \leftarrow T(\hat{x})
682
              7: for i \leq n_{iter} do
683
                       X \leftarrow T(x)
684
                       z \leftarrow \hat{z} + (z - \hat{z}) \cdot i / n_{iter}
              9:
685
                       re\_z \leftarrow \text{Re}\{z\}
            10:
686
                       im\_z \leftarrow \operatorname{Im}\{z\}
            11:
687
                       tape\_real.watch(re\_z)
            12:
688
            13:
                       tape\_imag.watch(im\_z)
689
                       \hat{z} \leftarrow re\_z + j \cdot im\_zx_{rec} \leftarrow T^{-1}(\hat{z})
            14:
            15:
690
            16:
                       y \leftarrow f(x_{rec})
691
                                                                                                                                  \triangleright Calculate \frac{\partial g}{\partial p_i}
            17:
                       re\_dy \leftarrow tape\_real.gradient(y, re\_z)
692
                                                                                                                                  \triangleright Calculate \frac{\partial g}{\partial g_i}
693
                       im\_dy \leftarrow tape\_imag.gradient(y, im\_z)
            18:
694
            19:
                       sum\_real \leftarrow sum\_real + re\_dy
695
            20:
                       sum\_imag \leftarrow sum\_imag + im\_dy
696
            21:
                       i \leftarrow i+1
697
            22: end for
            23: sum\_real \leftarrow sum\_real/n_{iter}
698
            24: sum\_imag \leftarrow sum\_imag/n_{iter}
699
            25: IG = \text{Re}\{z - \hat{z}\} \cdot sum\_real + \text{Im}\{z - \hat{z}\} \cdot sum\_imag
700
```

Algorithm 3 Complex Target Domain IG with complex differential

```
703
             Input: f(\cdot), x, \hat{x}, n_{iter}
704
             Output: IG
705
              1: i \leftarrow 1
706
              2: sum \leftarrow 0
              3: tape \leftarrow tensorflow.GradientTape()
708
              4: \hat{z} \leftarrow T(\hat{x})
709
              5: for i \leq n_{iter} do
710
                        z \leftarrow T(z)
                        z \leftarrow \hat{z} + (z - \hat{z}) \cdot i / n_{iter}
711
                        tape.watch(z)
712
                        x_{rec} \leftarrow T^{-1}(z)y \leftarrow f(x_{rec})
              9:
713
             10:
714
                        dy \leftarrow tape.gradient(y, X)
             11:
715
             12:
                        sum \leftarrow sum + \overline{dy}
716
             13:
                        i \leftarrow i + 1
717
             14: end for
718
             15: sum \leftarrow sum/n_{iter}
             16: IG = 2 \operatorname{Re} \{ (z - \hat{z}) \cdot sum \}
720
```

B PROOF OF LEMMA 4.1

Lemma. Let $g: \mathbb{C}^n \to \mathbb{R}$, z = p + jq, with $p, q \in \mathbb{R}^N$, $\gamma(t) = \hat{z} + t(z - \hat{z})$, $t \in [0, 1]$ the line from the baseline point \hat{z} to the input point z and $n(t) = \text{Re}\{\gamma(t)\}$ and $m(t) = \text{Im}\{\gamma(t)\}$, $n(t), m(t) \in \mathbb{R}^n$. Then the IG of g in z is given by:

$$IG_i^{\mathbb{C}^n}(z) = \int_0^1 \left(\frac{\partial g}{\partial p_i} n_i'(t) + \frac{\partial g}{\partial q_i} m_i'(t) \right) dt \tag{8}$$

Proof. Let $u: \mathbb{R}^{2n} \to \mathbb{R}$ such that $g(z) = u(w), \forall z = p + jq, w = [p, q]$. For the differential form of u:

$$du := \sum_{i=0}^{2N} \frac{\partial u}{\partial w_i} dw_i \tag{9}$$

Similarly to the g(z)-u(w) equivalence, we consider the equivalence between $\gamma(t)$ and $a(t) = [n(t), m(t)] \in \mathbb{R}^{2n}$. Then the pullback of du by a is:

$$\boldsymbol{a}^* du := \sum_{i=0}^{2N} \frac{\partial u}{\partial w_i} a_i'(t) dt \tag{10}$$

Denoting with a'_i the i-th element of da/dt. The line integral of u along the line defined by a is:

$$\int_{\gamma} du = \int_{\gamma} \mathbf{a}^* du = \int_0^1 \sum_{i=0}^{2N} \frac{\partial u}{\partial w_i} a_i'(t) dt = \sum_{i=0}^{2N} \int_0^1 \frac{\partial u}{\partial w_i} a_i'(t) dt$$
 (11)

Due to the equivalence between w and p, q and u and g the latter sum can be formulated as:

$$\int_{\gamma} du = \sum_{i=0}^{N} \left(\int_{0}^{1} \frac{\partial g}{\partial p_{i}} n_{i}'(t) dt + \int_{0}^{1} \frac{\partial g}{\partial q_{i}} m_{i}'(t) dt \right) = \sum_{i=0}^{N} \int_{0}^{1} \left(\frac{\partial g}{\partial p_{i}} n_{i}'(t) + \frac{\partial g}{\partial q_{i}} m_{i}'(t) \right) dt \quad (12)$$

which concludes the derivation.

C Derivation of Definition 4.1

From Lemma 4.1 we conclude to Definition 4.1 by considering $g(z) = f(T^{-1}(z))$ and the complex differential form Range (1998):

$$dg = \partial g + \overline{\partial}g \tag{13}$$

759 (1998) γ by γ is

with $\partial g = \sum \partial g/\partial z_i dz_i$, $\overline{\partial} g = \sum \partial f/\partial \overline{z_i} \overline{dz_i}$. The complex partial derivatives are defined as Range (1998) $\partial/\partial z_i = 1/2(\partial/\partial p - j\partial/\partial q)$ and $\partial/\partial \overline{z_i} = 1/2(\partial/\partial p + j\partial/\partial q)$. Then the pullback of dg by γ is :

 $\gamma^* dg = \sum \frac{\partial g}{\partial z_i} \gamma_i'(t) dt + \sum \frac{\partial g}{\partial \overline{z_i}} \overline{\gamma_i'(t)} dt$ (14)

Since $g \in \mathbb{R}$, $\partial g/\partial \overline{z} = \overline{(\partial g/\partial z)}$, thus:

$$\gamma^* dg = 2 \operatorname{Re} \left\{ \sum \frac{\partial g}{\partial z_i} \gamma_i'(t) dt \right\}$$
 (15)

Expanding the product into its real and imaginary parts produces the same form as eq. 12:

$$\gamma^* dg = 2 \operatorname{Re} \left\{ \sum \frac{1}{2} \left(\frac{\partial g}{\partial p_i} - j \frac{\partial g}{\partial q_i} \right) (n_i' + j m_i'(t)) dt \right\} = \sum \left(\frac{\partial g}{\partial p_i} n_i'(t) + \frac{\partial g}{\partial q_i} m_i'(t) \right)$$
(16)

Thus, the complex integrated gradient definition can be rewritten as:

$$IG_i^{\mathbb{C}^n} = 2\int_0^1 \operatorname{Re}\left\{\frac{\partial g}{\partial z_i}\gamma_i'(t)\right\}dt$$
 (17)

D RELATION TO VIRTUAL INSPECTION LAYERS

We demonstrate here the equivalence between eq. 6 and the Virtual Inspection Layer Vielhaben et al. (2024) for the case of the Discrete Fourier Transform (DFT) domain saliency maps.

Denote the DFT transform z = Tx with :

$$T_{nk}^{-1} = \frac{1}{\sqrt{N}} e^{2\pi k n/N} \tag{18}$$

Thus from eq.6

$$IG_k^{DFT} = 2 \int_0^1 \text{Re} \left\{ \sum_{n=0}^{N-1} \frac{\partial f}{\partial x_n} T_{nk}^{-1} (z_k - \hat{z_k}) \right\} dt = \sum_{n=0}^{N-1} \text{Re} \left\{ T_{nk}^{-1} (z_k - \hat{z_k}) \right\} 2 \int_0^1 \frac{\partial f}{\partial x_n} dt$$
$$= 2 \sum_{n=0}^{N-1} \text{Re} \left\{ T_{nk}^{-1} (z_k - \hat{z_k}) \right\} \frac{IG_n}{x_n - \hat{x}_n}$$

Denoting $(z_k - \hat{z_k}) = r_k e^{j\phi k}$ then

$$\operatorname{Re}\left\{T_{nk}^{-1}(z_k - \hat{z_k})\right\} = \frac{r_n}{\sqrt{N}}\cos\left(\frac{2\pi kn}{N} + \phi_k\right) \tag{19}$$

And finally,

$$R_k = 2r_k \sum \cos\left(\frac{2\pi kn}{N} + \phi_k\right) \frac{R_n}{x_n - \hat{x}_n} \tag{20}$$

Which is equivalent to the method of Vielhaben et al. (2024)

E RELATIONSHIP BETWEEN FREQUENCY-DOMAIN IG AND FREQUENCY RESPONSE

We probe the two convolutional channels of section 3.2 with sinusoid signals at varying frequencies, ξ_i :

$$x_i(t) = \cos(2\pi\xi_i t + \phi) \tag{21}$$

For each input we perform frequency-domain IG which yields a saliency map described by eq. 7. We aggregate all produced IGs and compare them to each filter's frequency response:

$$b_i = \|\sum_n w_n e^{-2\pi\xi_i n}\|$$
 (22)

The results are presented in Figure 5.

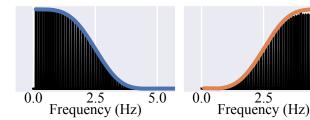


Figure 5: Frequency response (blue - orange) and frequency integrated gradients (black) for the two channels of the model of Section 3.2. We probe the model, performing frequency IG on samples with varying base frequencies.

5.0

F FEATURE-LEVEL INSERTION-DELETION

We perform insertion-deletion evaluation tests on the three examples presented in Section 5. Our evaluation indicates that component-level attributions provide more faithful and concentrated evidence for the models' predictions than time-domain attributions: adding top-rated component features rapidly reconstructs the output, while removing them destroys it.

F.1 HEART RATE EXTRACTION FROM PHYSIOLOGICAL SIGNALS

We follow the following procedure:

- 1. **Select** k% **features**, either in time or in frequency domain. For the frequency and time domain IG we select the k components with the highest IG score. For the random intervention, we randomly sample k% unique frequency bins.
- 2. Insert/delete k components to generate modified samples x_{mod} .
- 3. **Infer** heart rate with x_{mod} input.
- 4. Compare $f(x_{mod})$ with the original heart rate inference before any interventions f(x).

An example of inference after inserting/deleting input features is presented in Figure 6. We plot the heart rate inference throughout the entire 2-hour session of subject 15 from the PPG-Dalia dataset. The results for the entire PPGDalia dataset are summarised in Table 1.

Top k%-features	3.125 %	25%	50%	
Deletion ↑				
Frequency IG Time IG Random	66.39 10.13 8.53	133.56 50.86 37.03	127.13 104.84 68.34	
Insertion ↓				
Frequency IG Time IG Random	37.98 94.58 123.71	20.08 57.27 100.39	9.86 58.61 66.67	

Table 1: Insertion-deletion evaluation dropping the k% most important features. Deletion/Insertion distance (expressed in Beats per Minute- BPM) from the original HR inference averaged across 15 subjects of PPGDalia.

F.2 ELECTROENCEPHALOGRAPHY-BASED EPILEPTIC SEIZURE DETECTION

We used the Physionet Siena Scalp EEG Database v1.0.0 Detti (2020); Detti et al. (2020); Goldberger et al. (2000). For each subjects' sessions we retrieved the first sample that is detected as seizure by the zhu-transformer. For each sample, we generated ICA-domain IG saliency maps and

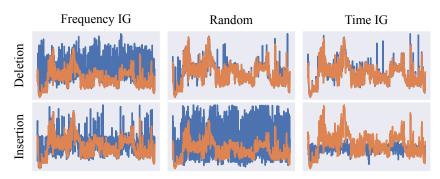


Figure 6: **Example of heart rate inference after deleting features.** We plot the entire session of subject 15 from PPGDalia. For each insertion/deletion we retain/delete 3.125% of the input features. For the Fourier and time IG these are the frequency bins and time-points with the highest assigned IG score. In the random case we randomly drop 3.125% of the frequency bins. We plot the **original HR inference** over the duration of the session and the model's output after **modifying the input** accordingly.

performed insertion/deletion with the most important IC. We kept track of the change in the seizure classification probability, $\Delta p = p(x_{mod}) - p(x)$, as we:

- 1. Delete the most important component and perform inference,
- 2. Maintain the most important component, delete the rest of the components and perform seizure classification.

We compared these results with randomly choosing an IC component and performing the same insertion/deletion evaluation.

	ICA IG	Random
Deletion $\Delta p \uparrow$	0.1776 0.0696	0.0083
Insertion $\Delta p \downarrow$	0.0696	0.4396

Table 2: Insertion-deletion evaluation on the seizure detection model.

G EXAMPLE TIME-DOMAIN ATTRIBUTIONS

Figures 7, 8 and 9 present the time-domain attributions from the examples of Section 5. In all three cases interpreting the time-domain saliency maps is difficult and of limited utility.

Heart rate inference. Time-domain IG highlights individual time-points of the PPG input. However, it is difficult to assess:

- 1. Does an individual time-point contribute to the heart or interference components? In the time-domain both the effect of heart and interference are mixed, and each time-point contains information from both of these components. In contrast, in images when there is component (object) overlap, one component blocks the other and a single pixel carries single-component information.
- 2. Which time-points should be the most important/influential? From domain knowledge we know that oscillations around the ground truth heart rate should be the ones affecting the model's output. However, we do not have any such insights in the time domain and the component overlap further complicates oscillation indification in time.

Consequently these saliency maps do not allow us to answer to the interpretability task of Section 5.1.

Seizure detection. Similarly to the heart rate example, here it is not easy to visually identify the seizure-related oscillations in the time-domain saliency map.

Time series forecasting. The time-domain IG highlights mostly the last input time-points.

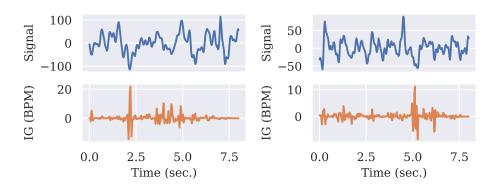


Figure 7: **Time-domain IG for HR inference**. We present the same two inputs as in Figure 2. For each time point in the input we assign a significance value. **Top:** Raw time-domain input which is processed by the model. **Bottom:** IG saliency map expressed in the original time domain.

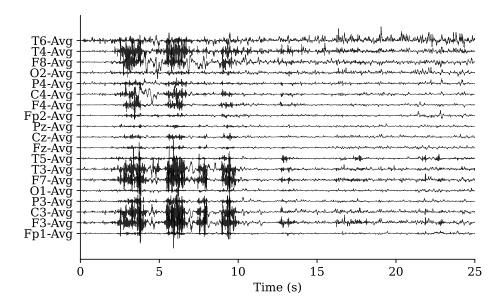


Figure 8: **Time-domain IG for seizure classification**. For each time point on each channel we assign a significance value.

H ADDITIONAL EXAMPLES

We present additional Cross-domain IG examples in Figures 10, 11 and 12.

I EEG AND ICA

The raw EEG input is presented in Figure 13.

The implementation of the zhu-transformer we used can be found here $https://github.com/esl-epf1/zhu_2023$.

The application of ICA in EEG signals is based on the general assumption that the EEG data matrix $X \in \mathbb{R}^{N \times M}$ is a linear mixture of different sources (activities) $S \in \mathbb{R}^{N \times M}$ with a mixing matrix $A \in \mathbb{R}^{N \times N}$ such that X = AS, where N is both the number of sources and EEG channels, and

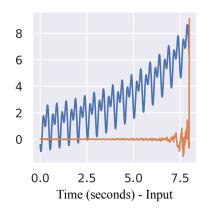


Figure 9: **Time-domain IG for time-series forecasting**. We plot the raw time-domain **input** along with the **IG importance** for each time-point in the input.

M is the number of samples in the dataset. Sources are assumed to be statistically independent and stationary. These assumptions can be leveraged to compute an inverse unmixing matrix $W = A^{-1} (\in \mathbb{R}^{N \times N})$, such that S = WX. Finding W is an ill-posed problem without an analytical solution which can be estimated by means of different ICA algorithms Hyvärinen et al. (2001); Klug and Gramann (2021). ICA is used in EEG to decompose the signal into independent components that separate the signal of interest from various sources of artifacts Winkler et al. (2011). In this work, for ICA we selected the FastICA algorithm implemented in sklearn (max_iter = $3 \cdot 10^4$, tol = $1 \cdot 10^{-8}$).

The independent channels estimated using ICA are presented in Figure 14.

J GENERATED TIME SERIES FOR TIMESFM FORECASTING

We generate a synthetic time series signal, x(t), composed of an exponential trend, $x_{trend}(t)$, and a seasonal component, $x_{seasonal}(t)$:

$$x_{trend}(t) = e^{\frac{t}{\alpha}}$$

$$x_{seasonal}(t) = sin(2\pi \cdot \xi \cdot t + \phi) + sin(2\pi \cdot 2\xi \cdot t + \phi)$$

$$x(t) = x_{trend}(t) + x_{seasonal}(t)$$

For the example in Section 5.3 $\alpha=4,\ \xi=2Hz$. For the samples presented in Appendix H they were randomly sampled from $\alpha\sim U(4.0,7.0)$ and $\xi\sim U(3.0,8.0)[Hz]$. A window of 512 time points, starting at t=0, are given as input to TimesFM which generates forecasts up to 128 time points in the future from t=512. The input time series and STL decomposition are presented in more detail in Figure 15.

K LIMITATIONS

Our method requires an invertible, differentiable transform and a carefully selected baseline point. Consequently, we excluded non-invertible transforms, and further investigation is needed for approximate-invertible cases. Baseline selection also plays a role in the final saliency map. We focused on the zero-signal as the baseline point - future work should include an extensive investigation into the effects of the baseline selection. The current implementation also focuses on a linear integration path, reflecting the original IG. However, other non-linear paths, e.g., Guided IG Kapishnikov et al. (2021), should be explored. Finally, multiple transforms can be combined to provide a multi-faceted saliency map, such as ICA combined with frequency domains - and automatic transform selection could help streamline the process. We leave *ensemble* domains and automatic domain selection as future work.

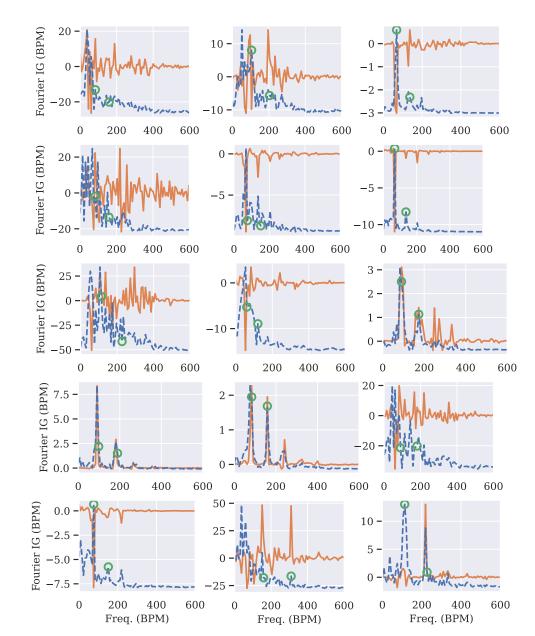


Figure 10: Frequency-domain IG for heart rate inference model.

L EXPERIMENTS COMPUTE RESOURCES

All experiments were run on an NVIDIA Tesla V100 with 32GB memory.

M USE OF LLMS

We used a large language model (LLM) solely for light copy-editing (grammar and wording).

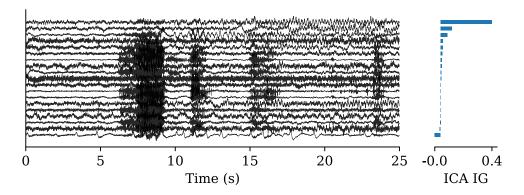


Figure 11: **ICA-domain IG for seizure detection model.** Similarly to the example presented in Section 5.2, the first channel contains the majority of the seizure components. IC channels that contain mostly interference are assigned a very small IG score.

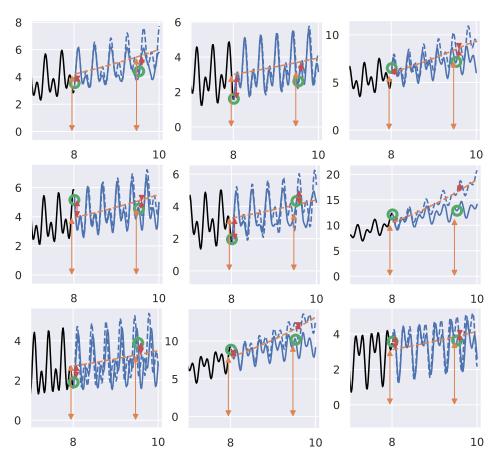


Figure 12: Seasonal-Trend IG for TimesFM forecasts. We generate synthetic samples by sampling them as described in Appendix J.

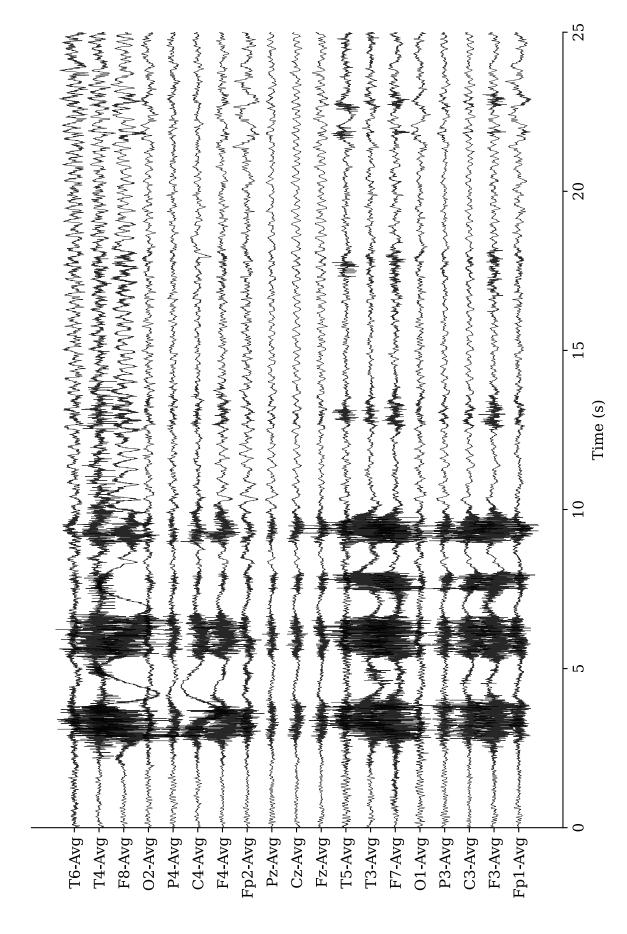


Figure 13: EEG signal in the original channel space.

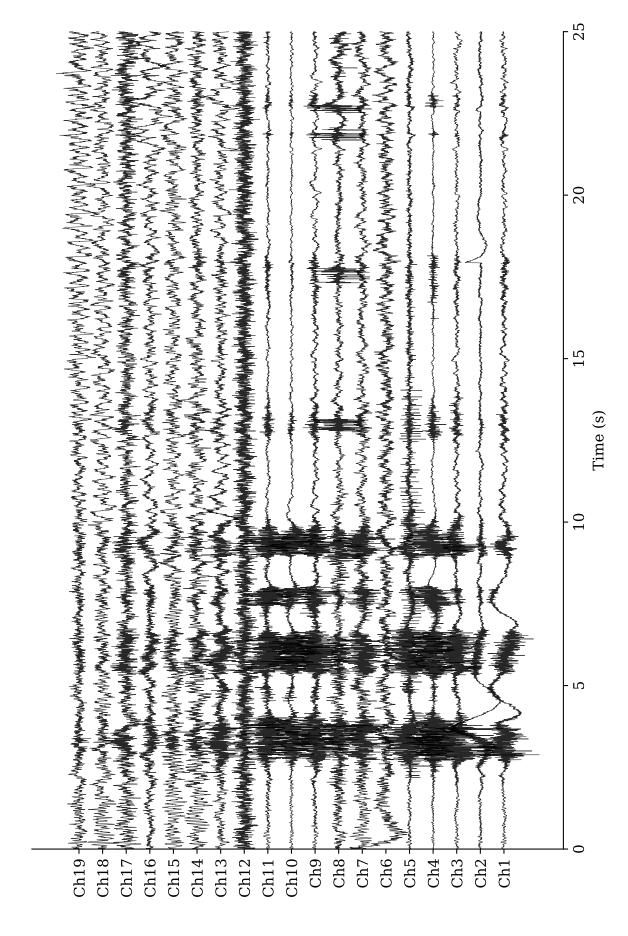


Figure 14: EEG signal in the Independent Component space.

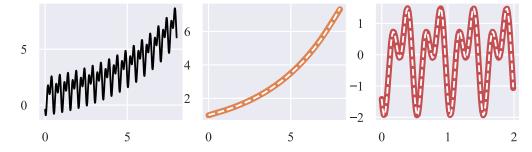


Figure 15: **Input time series for forecasting and successful STL decomposition. Left:** time series with a **trend** and a **seasonal** component. **Center:** The decomposed **trend** component and ground truth trend (white dashed line). **Right:** The decomposed **seasonal** component and ground truth seasonality (white dashed line).