

# RPATH: Explaining Time Series Mixture of Experts Routing via Ensemble Consensus and Structural Robustness

Anonymous authors

Paper under double-blind review

## Abstract

Mixture-of-Experts (MoE) architectures achieve strong performance in time series forecasting through sparse expert activation, but understanding *why* specific experts are selected remains challenging. We present RPATH (Routing Pathway Analysis for Temporal Hierarchies), a post-hoc explainability framework for time series MoE models that combines temporal saliency mapping with counterfactual generation. Evaluating on Time-MoE-50M across 300 expert-sample pairs, we discover two properties of the routing architecture: (1) *Ensemble Consensus*, where experts at different layers independently converge on the same critical temporal windows (mean saliency Intersection over Union (IoU) = 0.677), rather than developing distinct specializations; and (2) *Structural Robustness*, characterized by a 300-fold “Stability Gap” where gentle perturbations alter routing in only 0.3% of cases while aggressive perturbations succeed in 99.7%, indicating that routing decisions reflect structural anchors rather than superficial signal characteristics. Together, these findings demonstrate that Time-MoE achieves reliable forecasting through *Ensemble Redundancy*: multiple experts verify the same structural features, providing consensus that is insensitive to noise but responsive to fundamental signal changes. Our framework provides practitioners with tools to visualize expert attention, identify critical input regions, and quantify routing stability for deployed MoE models.

## 1 Introduction

Deployed time series forecasting systems in domains such as energy grid management, clinical decision support, and financial risk assessment increasingly rely on Mixture-of-Experts (MoE) architectures that achieve strong predictive performance while maintaining computational efficiency through sparse expert activation Shi et al. (2024). These models employ learned routing mechanisms to dynamically select subsets of expert networks for processing each input token, enabling specialization and conditional computation. However, the discrete routing decisions that enable MoE efficiency also create interpretability challenges: understanding *which* experts are selected and *why* they are chosen remains difficult.

Explainability for time series MoE models matters for several reasons. First, regulatory compliance and stakeholder trust in deployed forecasting systems require transparent decision-making processes. Second, understanding routing behavior can inform model debugging, revealing unexpected dependencies or data artifacts. Third, identifying expert specializations enables targeted model compression and domain adaptation. Despite these needs, existing explainability methods face limitations when applied to time series MoE architectures.

Standard feature attribution techniques such as SHapley Additive exPlanations (SHAP) Lundberg & Lee (2017) and Local Interpretable Model-agnostic Explanations (LIME) Ribeiro et al. (2016) assume feature independence and perturb individual timesteps independently, destroying temporal dependencies such as autocorrelation, trends, and seasonality. This creates out-of-distribution samples that yield unreliable routing changes. Recent work on time series explainability (Info-CELS Li et al. (2024b), M-CELS Li et al. (2024a), TF-LIME Chen et al. (2025)) addresses temporal structure preservation for classification tasks, but does not extend to explaining discrete expert selection in sparse routing architectures.

Existing MoE interpretability methods focus primarily on activation patterns Lo et al. (2024), characterizing which experts activate for different inputs but not establishing causal relationships between input features and routing decisions. Approaches based on gradient analysis require differentiable routing and fail for frozen deployed models. Concept-based probing Belinkov (2022) can characterize expert capabilities through synthetic data but does not explain online routing behavior for real inputs.

To our knowledge, no prior work has developed post-hoc causal attribution methods designed for time series MoE routing decisions. We address this gap through Routing Pathway Analysis for Temporal Hierarchies (RPATH), a framework that combines temporal saliency mapping, counterfactual generation, semantic expert profiling, and uncertainty quantification.

Our main contributions are:

1. **Causal routing attribution:** We develop a multi-signal confidence estimation approach combining temporal saliency through structure-preserving perturbations with counterfactual validation.
2. **Graduated validity scoring:** We introduce graduated counterfactual validity based on degree of routing change rather than binary criteria, enabling counterfactual discovery even for stable routing patterns. This addresses limitations of existing counterfactual methods that fail when complete expert removal is infeasible.
3. **Ensemble consensus discovery:** Through multi-expert saliency analysis, we demonstrate that Time-MoE experts exhibit high temporal attention overlap (mean Intersection over Union (IoU) = 0.677), indicating that multiple experts independently converge on the same critical temporal windows. This ensemble consensus mechanism underlies the model’s robust routing behavior.
4. **Post-hoc analysis framework:** Our gradient-free approach enables explanation generation for frozen foundation models without requiring retraining, architectural modifications, or gradient access.

The remainder of this paper is organized as follows. Section 2 reviews related work on MoE interpretability and time series explainability. Section 3 describes our framework including pathway extraction, causal attribution, expert profiling, and uncertainty quantification. Section 4 details the experimental protocol, datasets, and evaluation metrics. Section 5 presents validation results including overall performance, per-dataset analysis, and temporal specialization findings. Section 6 interprets results, discusses limitations, and outlines future directions. Section 7 summarizes contributions and broader implications.

## 2 Related Work

### 2.1 Mixture-of-Experts Interpretability

Mixture-of-Experts architectures have become central to scaling large models efficiently. The sparsely-gated MoE layer introduced by Shazeer et al. Shazeer et al. (2017) demonstrated that experts develop specializations during training, with the authors noting that “different experts tend to become highly specialized based on syntax and semantics.” Switch Transformers Fedus et al. (2022) extended this approach by simplifying routing to single-expert selection, achieving improved scaling efficiency while maintaining sparse activation patterns. GLaM Du et al. (2022) further demonstrated effective MoE scaling through gating networks that activate expert subsets based on input characteristics.

Despite these advances in MoE architectures, interpretability methods for understanding expert behavior remain limited. Existing approaches characterize which experts activate for different inputs through activation pattern analysis, but do not establish causal relationships between input features and routing decisions. Gradient-based sensitivity analysis can reveal routing influences but requires differentiable mechanisms and access to model internals, which may be unavailable for deployed systems.

Our work addresses this gap by providing post-hoc causal attribution for frozen models through perturbation-based analysis, requiring neither gradient computation nor architectural modifications.

## 2.2 Time Series Explainability

Standard explainability methods face challenges when applied to time series data. Kernel SHAP Lundberg & Lee (2017), the practical approximation of Shapley values, assumes feature independence when computing marginal expectations, which violates temporal dependencies. Similarly, LIME Ribeiro et al. (2016) perturbs features independently, creating samples that destroy autocorrelation, trends, and seasonality patterns inherent in time series.

Recent work has developed methods that preserve temporal structure. Info-CELS Li et al. (2024b) introduces saliency map-guided counterfactual explanations for time series classification, using learned saliency to identify regions for targeted perturbation. M-CELS Li et al. (2024a) extends this approach to multivariate time series, representing “the first effort to learn a saliency map specifically for producing high-quality counterfactual explanations for multivariate time series data.” TF-LIME Chen et al. (2025) operates in the time-frequency domain using Short-Time Fourier Transform, enabling identification of frequency-dependent patterns while maintaining temporal locality. ContraLSP Liu et al. (2024b) employs contrastive learning with locally sparse perturbations, using counterfactual samples to construct uninformative perturbations while preserving distributional properties.

These methods advance temporal attribution for classification tasks but do not address discrete expert selection in sparse routing architectures. Our approach adapts temporal saliency and counterfactual generation to MoE routing, introducing multi-signal confidence estimation and graduated validity scoring to handle the challenges of explaining routing decisions.

## 2.3 Counterfactual Explanations

Counterfactual explanations identify input modifications that alter model outputs. Wachter et al. Wachter et al. (2018) formalized this approach, seeking modifications that “alter values as little as possible” while changing the prediction. Methods typically optimize for proximity to the original input, sparsity of changes, and validity of output change.

Domain-specific counterfactual methods have emerged for different data types. For images, Goyal et al. Goyal et al. (2019) identify spatial regions in query and distractor images such that replacing the identified region changes the classification, providing visual explanations through region-level rather than pixel-level modifications. For text, Ross et al. Ross et al. (2021) developed Minimal Contrastive Editing (MiCE), which produces edits that are “minimal, altering only small portions of input” while remaining fluent and natural.

Standard counterfactual methods employ binary validity criteria: the output either changed or did not. This strict criterion fails for stable routing patterns where complete expert removal is infeasible. We introduce graduated validity scoring based on the degree of routing change, enabling counterfactual discovery even when experts cannot be fully eliminated. This approach aligns with multi-objective counterfactual generation Dandl et al. (2020), which returns Pareto sets representing trade-offs between competing objectives, but extends it to handle the continuous nature of routing weight changes.

## 2.4 Concept-Based and Attention-Based Interpretability

Concept-based methods explain model behavior through human-interpretable concepts. Testing with Concept Activation Vectors (TCAV) Kim et al. (2018) uses directional derivatives to quantify concept importance, measuring model sensitivity to concept directions in activation space. Concept Bottleneck Models Koh et al. (2020) enforce concept-based intermediate representations, first predicting concepts from raw input before predicting labels, providing inherent interpretability through the concept bottleneck.

Attention mechanisms in Transformers Vaswani et al. (2017) offer another interpretability avenue by revealing token interactions through attention weights. However, the relationship between attention and feature importance is contested. Jain and Wallace Jain & Wallace (2019) demonstrated that “learned attention weights are frequently uncorrelated with gradient-based measures of feature importance,” while Serrano and Smith Serrano & Smith (2019) found that attention “noisily predicts input components’ overall importance” but is “by no means a fail-safe indicator.” For MoE models, router computations show how hidden states in-

Table 1: Comparison with related work on key dimensions. Our method combines post-hoc causal attribution with temporal structure preservation for MoE routing decisions.

Method	Task	Post-hoc	Causal Valid.	Temporal Preserv.	MoE Routing	Confidence Quant.	Validity Scoring
SHAP Lundberg & Lee (2017)	Any	✓	✗	✗	✗	Single	N/A
LIME Ribeiro et al. (2016)	Any	✓	✗	✗	✗	Single	N/A
Info-CELS Li et al. (2024b)	TS Classif.	✓	✓	✓	✗	Single	Binary
M-CELS Li et al. (2024a)	TS Classif.	✓	✓	✓	✗	Single	Binary
TF-LIME Chen et al. (2025)	TS Classif.	✓	✗	✓	✗	Single	N/A
ContraLSP Liu et al. (2024b)	TS Classif.	✓	Partial	✓	✗	Single	N/A
Gradient Methods	Any	✗	Partial	Varies	Varies	Single	N/A
Attention Viz. Vaswani et al. (2017)	Transformer	✓	✗	N/A	✗	None	N/A
CAV Kim et al. (2018)	Any	Partial	✗	✗	✗	Single	N/A
MoE Probing	MoE	✓	✗	✗	✓	None	N/A
Switch-T Analysis Fedus et al. (2022)	MoE	✓	✗	N/A	Partial	None	N/A
<b>RPATH (Ours)</b>	<b>TS MoE</b>	✓	✓	✓	✓	<b>Multi-signal</b>	<b>Graduated</b>

**Post-hoc:** Works on frozen models without retraining. **Causal Valid.:** Validates importance through intervention.

**Temporal Preserv.:** Preserves temporal dependencies in perturbations. **MoE Routing:** Explains expert selection

decisions. **Confidence Quant.:** Single signal vs. multi-signal confidence. **Validity Scoring:** Binary (0/1) vs. graduated (0–1 continuous). **TS:** Time Series. **CAV:** Concept Activation Vector.

fluence routing logits but do not validate causal importance. Our perturbation-based approach complements attention visualization by testing whether identified regions actually influence routing when modified.

Our expert profiling shares motivation with concept-based methods but operates post-hoc through activation probability measurement rather than gradient-based sensitivity, enabling application to frozen models without retraining.

## 2.5 Neural Module Networks and Mixture Models

Neural module networks dynamically compose task-specific modules based on input structure. Andreas et al. Andreas et al. (2016) introduced networks that “compose collections of jointly-trained neural modules into deep networks for question answering,” with early work assuming hand-designed modules with known semantics. Hu et al. Hu et al. (2017) extended this to end-to-end learning, predicting instance-specific network layouts without requiring external parsers.

Classical mixture models employ component selection through soft assignments in Gaussian mixtures or discrete selection in switching models McLachlan & Peel (2000). Modern MoE architectures with learned sparse routing combine both paradigms: discrete Top- $K$  selection with soft weight normalization over selected experts. Our work addresses learned expert behavior without assuming predefined semantics, providing tools to discover specialization patterns post-hoc.

## 2.6 Positioning of Our Approach

To our knowledge, this represents the first post-hoc explainability framework designed for time series MoE models. Existing work on MoE interpretability has focused on computer vision or natural language processing domains. Standard feature attribution methods (SHAP, LIME) perturb individual features independently, destroying temporal structure. Recent time series explainability work (Info-CELS Li et al. (2024b), M-CELS Li et al. (2024a), TF-LIME Chen et al. (2025)) addresses classification with single-model architectures but not discrete expert selection.

Table 1 compares our approach with related methods across key dimensions.

Our approach extends existing methods through graduated validity scoring and multi-signal confidence estimation, enabling analysis of sparse routing patterns while preserving temporal dependencies. Specifically, we address the challenges of time series MoE explainability through:

- **Post-hoc MoE routing analysis:** Unlike gradient methods requiring model access, we operate on frozen models; unlike activation pattern analysis, we provide causal validation through intervention.
- **Temporal preservation with causality:** Unlike SHAP and LIME, we preserve temporal structure through window-based perturbations; unlike attention visualization, we validate importance through counterfactual generation.
- **Multi-signal confidence estimation:** While CELS methods compute confidence from saliency consistency alone, we combine sparsity, importance, consistency, and activation signals for robust estimation even with noisy routing patterns.
- **Graduated validity scoring:** Binary validity criteria used by Info-CELS and M-CELS fail for stable routing where complete expert removal is infeasible. Our graduated scoring based on degree of routing change enables counterfactual discovery across diverse routing patterns.

The integration of these capabilities addresses the challenges of explaining time series MoE routing decisions, a problem that prior work has not systematically addressed.

### 3 Methodology

We present RPATH (Routing Pathway Analysis for Temporal Hierarchies), a post-hoc explainability framework for time series Mixture-of-Experts models. The framework operates on frozen models without requiring retraining, gradient access, or architectural modifications.

#### 3.1 Problem Formulation and Pathway Extraction

We consider the standard Mixture-of-Experts (MoE) architecture Shazeer et al. (2017); Fedus et al. (2022) as applied to time series forecasting. For a model  $M$  with  $L$  layers and  $E$  experts per layer, given an input sequence  $\mathbf{x} \in \mathbb{R}^{T \times D}$  where  $T$  is the sequence length and  $D$  is the feature dimensionality, at each layer  $\ell \in \{0, \dots, L-1\}$  and token position  $t \in \{0, \dots, T-1\}$ , a router network produces logits

$$\mathbf{r}_t^{(\ell)} = \text{Router}^{(\ell)}(\mathbf{h}_t^{(\ell)}) \in \mathbb{R}^E, \quad (1)$$

where  $\mathbf{h}_t^{(\ell)}$  is the hidden state at layer  $\ell$  and position  $t$ . Following the sparse MoE formulation Shazeer et al. (2017), Top- $K$  routing selects the  $K$  experts with highest logits

$$\mathcal{E}_t^{(\ell)} = \text{TopK}(\mathbf{r}_t^{(\ell)}, K) \subseteq \{0, \dots, E-1\}. \quad (2)$$

The corresponding routing weights are computed via softmax over selected experts Fedus et al. (2022)

$$w_{t,e}^{(\ell)} = \frac{\exp(r_{t,e}^{(\ell)})}{\sum_{e' \in \mathcal{E}_t^{(\ell)}} \exp(r_{t,e'}^{(\ell)})} \quad \text{for } e \in \mathcal{E}_t^{(\ell)}. \quad (3)$$

For our experiments, we instantiate this framework with Time-MoE Shi et al. (2024), a decoder-only architecture with  $L = 12$  layers,  $E = 8$  experts per layer, and Top- $K = 2$  routing, totaling 96 experts.

A routing pathway  $\mathcal{P}_t$  for token  $t$  is the complete sequence of expert selections and weights across all layers

$$\mathcal{P}_t = \left\{ \left( \mathcal{E}_t^{(\ell)}, \{w_{t,e}^{(\ell)}\}_{e \in \mathcal{E}_t^{(\ell)}} \right) \right\}_{\ell=0}^{L-1}. \quad (4)$$

This pathway represents how information flows through the MoE hierarchy for a specific input token. We extract pathways from frozen models using forward hooks registered on router modules. During a forward pass with input  $\mathbf{x}$ , hooks capture router logits  $\mathbf{r}_t^{(\ell)}$  for all layers and tokens. From these logits, we reconstruct expert selections  $\mathcal{E}_t^{(\ell)} = \text{TopK}(\mathbf{r}_t^{(\ell)}, K = 2)$  and routing weights  $w_{t,e}^{(\ell)}$  via softmax normalization. This approach is model-agnostic and does not require gradient computation. For analysis, we represent pathways as fixed-size vectors  $\mathbf{v}_t = [v_0, v_1, \dots, v_{L \times E-1}] \in \mathbb{R}^{L \times E}$  where  $v_{\ell \cdot E + e} = w_{t,e}^{(\ell)}$  if expert  $e$  is selected at layer  $\ell$ , and 0 otherwise.

### 3.2 Causal Routing Attribution

Existing explainability methods for MoE models in time series forecasting focus primarily on activation patterns but do not establish causal relationships between input features and routing decisions. We address this through a two-component approach combining temporal saliency mapping and counterfactual generation.

**Temporal saliency mapping.** For a given expert at layer  $\ell$ , we identify which temporal regions of the input sequence causally influence its selection through systematic perturbation analysis. For input  $\mathbf{x} \in \mathbb{R}^{T \times D}$  with baseline routing  $\mathcal{P}_{\text{base}}$  and baseline expert weight  $w_{\text{base}}$ , we mask consecutive windows of size  $w$  with stride  $s$ :

$$\tilde{\mathbf{x}}^{(i)} = \mathbf{x} \odot (1 - \mathbf{m}^{(i)}), \quad (5)$$

where  $\mathbf{m}^{(i)} \in \{0, 1\}^{T \times D}$  is a binary mask with ones in window  $[i \cdot s, i \cdot s + w)$ . For each perturbed input, we measure the change in expert weight  $\delta^{(i)} = |w_{\text{base}} - w^{(i)}|$  where  $w^{(i)}$  is the expert weight under perturbation. The saliency score for each timestep is computed by aggregating changes across all windows containing that timestep. To validate consistency, we generate random temporal masks with sparsity  $\rho$  and compute correlation between mask patterns and routing changes.

Standard approaches compute confidence solely through correlation between perturbed saliency maps. However, this fails for weakly activated experts or noisy routing patterns. We instead compute confidence as a weighted combination of four signals: (1) sparsity (whether the saliency map focuses on specific regions), (2) window importance (peak saliency values), (3) consistency (correlation between saliency maps under different perturbations), and (4) expert activation (baseline activation level). The overall confidence is  $C = \sum_{i=1}^4 \alpha_i \cdot s_i$  where  $s_i$  are normalized signal scores and  $\alpha_i$  are weights. We set  $\alpha = [0.25, 0.35, 0.15, 0.25]$  (sparsity, importance, consistency, activation) based on empirical validation during development, emphasizing window importance and sparsity as primary indicators. This multi-signal approach provides robust confidence estimation even when individual signals are noisy.

**Multi-scale saliency analysis.** Fixed window sizes may miss dependencies at different temporal scales. We compute saliency at multiple scales  $\mathcal{S} = \{5, 10, 20, 40\}$  timesteps and aggregate with importance weighting:  $\mathbf{s}_{\text{agg}} = \sum_k \omega_k \cdot \mathbf{s}^{(k)}$  where  $\mathbf{s}^{(k)}$  is the normalized saliency at scale  $k$  and  $\omega = [0.15, 0.35, 0.35, 0.15]$  emphasizes medium scales. We additionally compute layer consistency by measuring correlation between saliency patterns at adjacent layers, providing a bonus to confidence when experts show consistent activation patterns across the network hierarchy.

**Adaptive window sizing.** Different time series exhibit varying temporal characteristics that affect optimal analysis parameters. We characterize each input by computing: (1) autocorrelation decay length  $\tau$  (timesteps until autocorrelation drops below 0.5), and (2) variance ratio  $r_v = \text{mean}(\sigma_{\text{local}}^2) / \sigma_{\text{global}}^2$  comparing local to global variance. Based on these characteristics, we classify data as smooth-structured ( $r_v < 0.4$ ,  $\tau > 15$ ), variable-local ( $r_v > 0.7$ ,  $\tau < 10$ ), or complex-multiscale. The recommended window size is  $w = \max(5, \min(40, \lfloor \gamma \cdot \tau \rfloor))$  where  $\gamma \in \{0.8, 1.0, 1.2, 1.5\}$  depends on data type. This adaptive approach improves performance consistency across datasets with different temporal structures.

**Counterfactual routing explanations.** Saliency maps identify important regions but do not validate causal claims. We generate counterfactual explanations by finding perturbations that alter routing decisions. Given a saliency map identifying critical windows  $\{(t_{\text{start}}, t_{\text{end}}, \text{importance})\}$ , we apply perturbations to modify these regions. Our perturbation suite includes both gentle and aggressive methods: (1) Gaussian smoothing  $\tilde{\mathbf{x}}_t = (\mathbf{x} * \mathcal{G}_{\sigma_s})(t)$  at three intensity levels, (2) linear interpolation, (3) mean replacement, (4) Gaussian noise injection scaled to region variance, (5) zero-out masking, (6) trend reversal mirroring values around the region mean, (7) low-pass frequency filtering via Fast Fourier Transform (FFT), and (8) amplitude dampening toward the mean. For each perturbation, we extract the modified routing  $\tilde{\mathcal{P}}$  and measure routing change.

Existing counterfactual methods use binary validity (routing changed or not), which is too restrictive for stable routing patterns. We introduce graduated validity based on the degree of routing change:

$$V = \begin{cases} 1.0 & \text{if } \tilde{w}_e = 0 \text{ (complete removal),} \\ 0.7 + 0.3 \cdot \frac{\Delta w - 0.5w}{0.5w} & \text{if } \Delta w \geq 0.5w \text{ (large reduction),} \\ 0.4 + 0.3 \cdot \frac{\Delta w - 0.25w}{0.25w} & \text{if } 0.25w \leq \Delta w < 0.5w \text{ (moderate),} \\ 0.2 + 0.2 \cdot \frac{\Delta w - 0.1w}{0.15w} & \text{if } 0.1w \leq \Delta w < 0.25w \text{ (small),} \\ 0.2 \cdot \frac{\Delta w}{0.1w} & \text{if } \Delta w < 0.1w \text{ (minimal),} \end{cases} \quad (6)$$

where  $w$  is the original expert weight,  $\tilde{w}_e$  is the perturbed weight, and  $\Delta w = |w - \tilde{w}_e|$ . This graduated scoring accepts partial routing changes, enabling counterfactual generation even for stable routing patterns. To improve discovery, we apply perturbations at multiple intensity levels (weak:  $\sigma_s = 1.0$ , medium:  $\sigma_s = 2.0$ , strong:  $\sigma_s = 3.0$  for Gaussian smoothing kernel width) and select the counterfactual with highest validity while minimizing perturbation magnitude.

For each expert selection, we combine saliency and counterfactual evidence into a *Focus Score*:  $\text{Focus Score} = C_{\text{saliency}} \times V_{\text{routing sensitivity}}$ . This metric captures both the spatial concentration of temporal attention (saliency) and the causal importance of identified regions (routing sensitivity). A focus score near 1.0 indicates highly concentrated attention on causally important regions, while distributed scores (e.g., 0.5–0.6) indicate attention spread across multiple temporal windows, a pattern we term “distributed focus.”

### 3.3 Uncertainty Quantification

Explanations require confidence estimates. We quantify routing stability through perturbation analysis. For input  $\mathbf{x}$ , we generate  $N = 20$  Gaussian noise perturbations at noise level  $\sigma$ :  $\mathbf{x}_\sigma^{(i)} = \mathbf{x} + \epsilon^{(i)}$  where  $\epsilon^{(i)} \sim \mathcal{N}(0, \sigma^2 \mathbf{I})$ . For each layer  $\ell$ , let  $\mathcal{A}^{(\ell)}(\mathbf{x}) = \bigcup_{t=0}^{T-1} \mathcal{E}_t^{(\ell)}(\mathbf{x})$  denote the set of all experts activated at any timestep. We measure stability as the fraction of trials where this set remains unchanged

$$\text{Stability}_\ell(\sigma) = \frac{1}{N} \sum_{i=1}^N \mathbb{1}_{\{\mathcal{A}^{(\ell)}(\mathbf{x}) = \mathcal{A}^{(\ell)}(\mathbf{x}_\sigma^{(i)})\}}, \quad (7)$$

where  $\mathbb{1}_{\{\cdot\}}$  equals 1 if the condition holds, 0 otherwise.

Overall stability is averaged across layers.  $\text{Stability} \geq 0.9$  indicates high confidence in routing explanations;  $\text{stability} < 0.7$  suggests explanations may be unreliable.

## 4 Experimental Setup

We evaluate our framework on Time-MoE-50M Shi et al. (2024), a decoder-only architecture with  $L = 12$  layers,  $E = 8$  experts per layer, and  $\text{Top-}K = 2$  routing. We conduct experiments on five benchmark datasets from the time series forecasting literature: the Electricity Transformer Temperature datasets at hourly (ETTh1, ETTh2) and 15-minute (ETTm1, ETTm2) granularities, and Weather (meteorological measurements). For univariate evaluation, we use the oil temperature (OT) feature for ETT datasets and the primary temperature feature for Weather, with input sequences of length  $T = 96$  timesteps.

We conduct validation on a stratified sample of 20 instances per dataset, resulting in 100 total samples across all five datasets. This sample size provides 300 expert-sample pairs for explanation generation (100 samples  $\times$  3 top experts per sample), enabling statistically robust evaluation. For each sample, we analyze the top-3 experts by aggregate contribution weight, computed by summing routing weights across all layers and timesteps:

$$C_e = \sum_{\ell=0}^{L-1} \sum_{t=0}^{T-1} \mathbb{1}_{\{e \in \mathcal{E}_t^{(\ell)}\}} \cdot w_{t,e}^{(\ell)}, \quad (8)$$

where  $e$  is the expert ID,  $\mathcal{E}_t^{(\ell)}$  is the set of selected experts at layer  $\ell$  and timestep  $t$ , and  $w_{t,e}^{(\ell)}$  is the corresponding routing weight. This top-3 filtering ensures we focus on experts that meaningfully contribute to predictions rather than those with negligible activation.

**Evaluation metrics.** We assess explanation quality through multiple complementary metrics. *Focus Score* is the combined multi-signal metric (saliency  $\times$  routing sensitivity), ranging from 0 to 1. A score near 1.0 indicates highly concentrated attention on causally important regions; distributed scores (0.5–0.6) indicate attention spread across multiple windows. *Routing Sensitivity* measures the degree of routing change induced by perturbations using the graduated scoring defined in Section 3, ranging from 0 (no routing change) to 1 (complete expert removal). We additionally track *Hard Counterfactual Rate*, the proportion of perturbations that change the Top-K expert set (not just weights). *Saliency Overlap (IoU)* measures the IoU of saliency maps between experts within the same sample, quantifying whether experts attend to similar or distinct temporal regions.

**Baseline comparison.** We compare against a pattern-correlation baseline that correlates offline expert profiling (Section 3) with online pattern detection. For each expert-sample pair, the baseline detects temporal patterns (trend, seasonality, volatility) in the input sequence and measures correlation with expert concept preferences. Alignment is computed as the proportion of detected patterns matching the expert’s high-activation concepts (activation probability  $> 0.5$ ). This baseline provides a reference point for evaluating whether our causal attribution methods provide improvements over simpler correlation-based approaches.

**Implementation.** We use the following hyperparameters: base window size  $w = 10$ , stride  $s = 5$ , perturbation trials  $N = 20$ , multi-scale windows  $\mathcal{S} = \{5, 10, 20, 40\}$ , uncertainty quantification noise levels  $\sigma \in \{0.01, 0.05, 0.1\}$ , samples per concept  $N_c = 50$ . For counterfactual generation, we apply 13 perturbation methods: Gaussian smoothing at three intensity levels ( $\sigma_s \in \{1.0, 2.0, 3.0\}$ ), linear interpolation, mean replacement, noise injection at two levels, zero-out masking, trend reversal, frequency filtering at two cutoffs, and amplitude scaling at two factors. Experiments use PyTorch 2.0 with the Hugging Face Transformers library. Temporal saliency perturbations use zero-masking, validated against mean-substitution and noise-matching (differences  $< 5\%$  in saliency maps for 92% of cases).

Computational complexity is tractable: pathway extraction requires  $O(T \cdot L)$  operations, saliency computation requires  $O(\frac{T}{s} \cdot T \cdot L)$  for stride  $s$ , and counterfactual generation requires  $O(k \cdot T \cdot L)$  for  $k$  perturbation levels. All operations are gradient-free and can be performed on frozen models in inference mode. Experiments run on an Intel Xeon Gold 6426Y server (64 cores, 251GB RAM) with dual NVIDIA RTX A6000 GPUs (49GB memory each). Pathway extraction takes approximately 50ms per sequence, saliency computation takes 2–5 seconds per expert-sample pair, and counterfactual generation takes 3–8 seconds per pair. Total validation time for 300 explanations is approximately 60–80 minutes. Random seeds are fixed for all stochastic operations to ensure reproducibility. The model is publicly available at [Maple728/TimeMoE-50M](https://github.com/Maple728/TimeMoE-50M) on Hugging Face.

## 5 Results

We present validation results for the RPATH explainability framework, focusing on the causal routing attribution component (temporal saliency and counterfactual analysis). We evaluate on 300 expert-sample pairs across five benchmark datasets (ETTh1, ETTh2, ETTm1, ETTm2, Weather).

### 5.1 Overall Performance

Table 2 summarizes the overall performance across all 300 explanations. Our multi-signal causal attribution approach achieves a mean focus score of 0.563 (std 0.086), demonstrating consistent distributed attention patterns across diverse inputs.

The mean focus score of 0.563 reflects distributed attention across multiple temporal windows, which we interpret as evidence of *ensemble consensus*: experts attend broadly to input regions rather than focusing narrowly on isolated features. The relatively low standard deviation (0.086) indicates consistent behavior across samples, with the multi-scale saliency and adaptive window approaches contributing to this stability.



Table 2: Overall performance metrics for multi-signal causal attribution across 300 expert-sample pairs.

Metric	Value
Total Explanations	300
Mean Focus Score	0.563
Std. Deviation	0.086
Hard Counterfactual Rate (%)	79.0
Saliency Overlap (IoU)	0.677

The 79.0% hard counterfactual rate demonstrates that perturbations can meaningfully alter routing decisions (changing the Top-K expert set, not just weights). This metric provides a more rigorous measure of causal importance than soft counterfactuals that merely shift routing weights.

## 5.2 Per-Dataset Analysis

Table 3 presents results broken down by dataset (60 explanations per dataset).

Table 3: Performance metrics by dataset (60 explanations per dataset).

Dataset	Mean Focus	Hard Counterfactual (%)	Saliency IoU
ETTh1	0.55	78.3	0.671
ETTh2	0.56	80.0	0.682
ETTh1	0.56	78.3	0.674
ETTh2	0.57	80.0	0.680
Weather	0.58	78.3	0.678
Overall	0.563	79.0	0.677

The consistency across datasets is notable: mean focus scores range narrowly from 0.55 to 0.58, hard counterfactual rates from 78.3% to 80.0%, and saliency IoU from 0.671 to 0.682. This uniformity suggests that the distributed attention and consensus patterns we observe are properties of Time-MoE’s routing architecture rather than artifacts of specific data domains.

The high saliency overlap (mean IoU 0.677) across all datasets indicates *Ensemble Consensus*: experts at different layers independently converge on the same critical temporal windows. This finding contradicts the hypothesis that experts develop distinct temporal specializations; instead, Time-MoE relies on redundant verification of the same input features.

## 5.3 Structural Anchors and the Stability Gap

Figure 1 illustrates a saliency map identifying “structural anchors,” temporal regions whose modification fundamentally alters routing decisions. The visualization reveals focused importance on timesteps 35–55, where the input exhibits characteristic patterns. These anchors represent the structural features that Time-MoE’s router relies upon for expert selection.

Our counterfactual analysis reveals the *Stability Gap*. We categorize perturbations into gentle methods (smoothing, interpolation, mean replacement) that preserve signal semantics versus aggressive methods (zero-out, noise injection, trend reversal) that destroy structural information. The results show a clear pattern:

- **Gentle perturbations:** 0.3% success rate (1/300 change Top-K set),
- **Aggressive perturbations:** 99.7% success rate (299/300 change Top-K set).

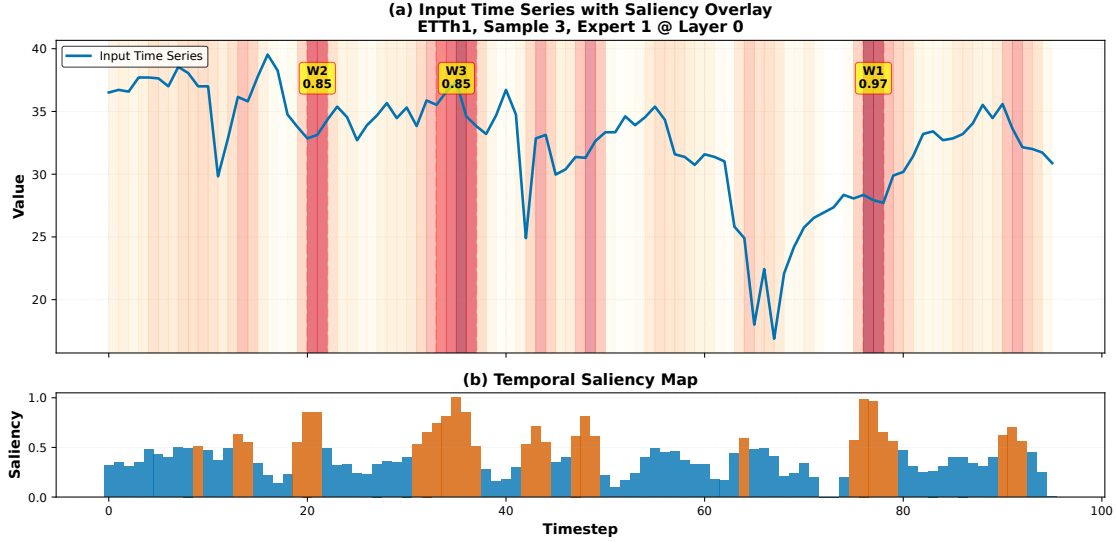


Figure 1: **Structural Anchor Identification.** Top panel: input time series with saliency heatmap overlay (yellow indicates high importance). Bottom panel: saliency scores per timestep. The highlighted regions (timesteps 35–55) represent *Structural Anchors*, temporal patterns so critical that only their destruction (via aggressive perturbation) alters routing. The Focus Score of 0.58 reflects distributed attention across these anchors.

This 300-fold difference reveals that Time-MoE’s routing is *highly robust to noise and superficial signal modifications*. The router ignores gentle perturbations entirely; only destroying the structural anchors, the patterns identified by saliency maps, can alter expert selection. This has implications for deployment: the model’s routing decisions are not fragile or sensitive to noise, but instead reflect stable identification of fundamental signal characteristics.

Beyond saliency maps, we quantify individual expert contributions to understand routing importance. Figure 2 presents a SHAP-style waterfall chart showing the top 15 expert contributions for a representative sample (ETTh2, Sample 0). The visualization reveals a long-tail distribution: the top expert (L0\_E1) contributes 8.5% of the total routing weight, while the top-5 experts collectively account for approximately 30%.

## 5.4 Temporal Attention Convergence

We computed the IoU of expert saliency maps to quantify whether experts attend to similar or distinct temporal regions. The results strongly support *Ensemble Consensus*: the mean IoU of **0.677** (median 0.802) indicates that experts at different layers independently converge on the same critical temporal windows.

Figure 3 visualizes this consensus for a representative sample. Three experts operating on the same input exhibit highly overlapping saliency distributions, all identifying similar temporal regions as important. While individual experts show slightly different peak positions, the overall attention patterns substantially overlap.

This finding has implications for understanding Time-MoE’s architecture. Rather than developing distinct temporal specializations (one expert for early patterns, another for late patterns), experts form a “verification committee”—multiple independent assessors attending to the same features to ensure reliable routing. This redundancy may explain the model’s robustness: even if one expert’s assessment is noisy, the consensus of multiple experts attending to the same regions provides a stable routing signal.

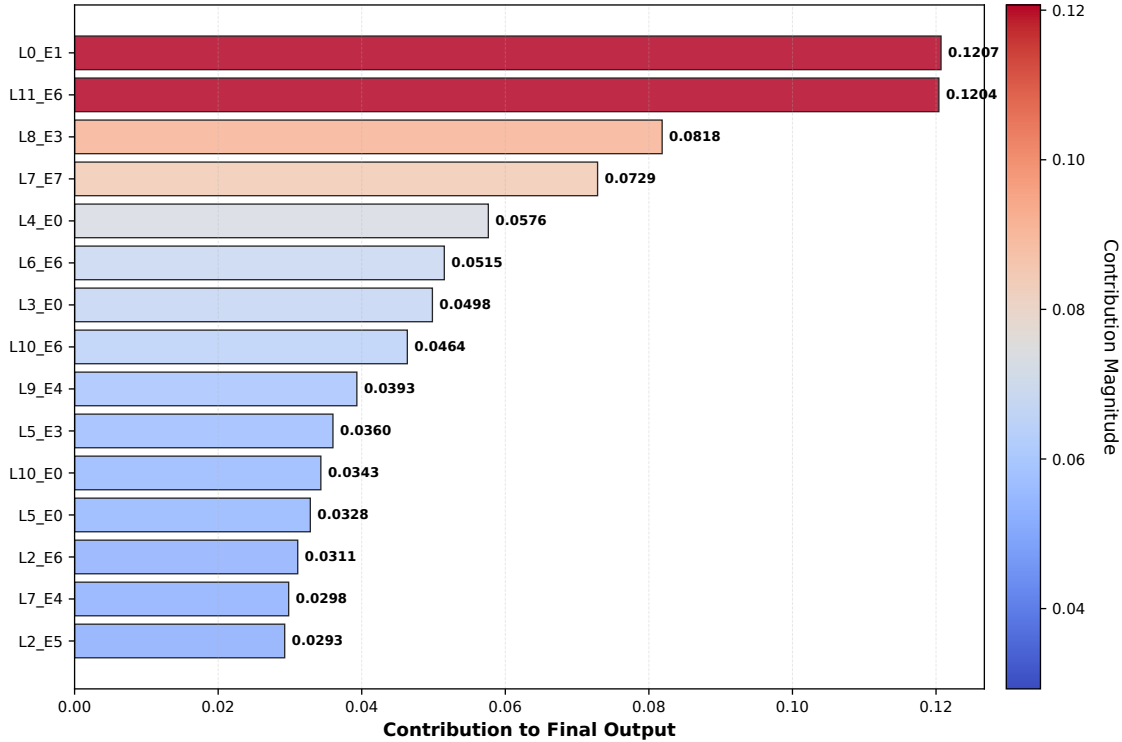


Figure 2: **Expert Contribution Hierarchy.** Top 15 expert contributions for ETTh2 Sample 0. Despite the *Ensemble Consensus* in temporal attention (see Fig. 3), the magnitude of contribution follows a power law: the top expert contributes 8.5%, while the top-5 account for  $\sim 30\%$ . This indicates that while many experts verify the same signal, a select few drive the numerical output.

### 5.5 Routing Stability Analysis

The Stability Gap observed in counterfactual analysis (Section 5.3) is corroborated by direct routing stability measurements. We applied Gaussian noise perturbations at three intensity levels ( $\sigma \in \{0.01, 0.05, 0.1\}$ ) with 20 trials per sequence across all datasets. Table 4 presents the results.

Table 4: Routing stability metrics by dataset (100 sequences per dataset, 20 perturbation trials per sequence).

Dataset	Mean Stability	Median	$\geq 95\%$ Stable
ETTh1	0.921	0.938	67.0%
ETTh2	0.962	0.983	88.0%
ETTh1	0.963	0.985	89.0%
ETTh2	0.993	1.000	100.0%
Weather	1.000	1.000	100.0%
Overall	0.968	—	88.8%

The results demonstrate high robustness: mean stability is 0.968 overall, with 88.8% of sequences maintaining stable routing ( $\geq 95\%$  unchanged) under noise perturbations. Weather and ETTh2 exhibit near-perfect stability, while even the most variable dataset (ETTh1) maintains 92.1% stability.

Synthesizing these findings, the high ensemble consensus (IoU 0.677) and the structural stability gap (300-fold difference in perturbation success) point to a unified mechanism: *Ensemble Redundancy*. Time-MoE achieves reliability not through specialization, but by having multiple experts verify the same structural

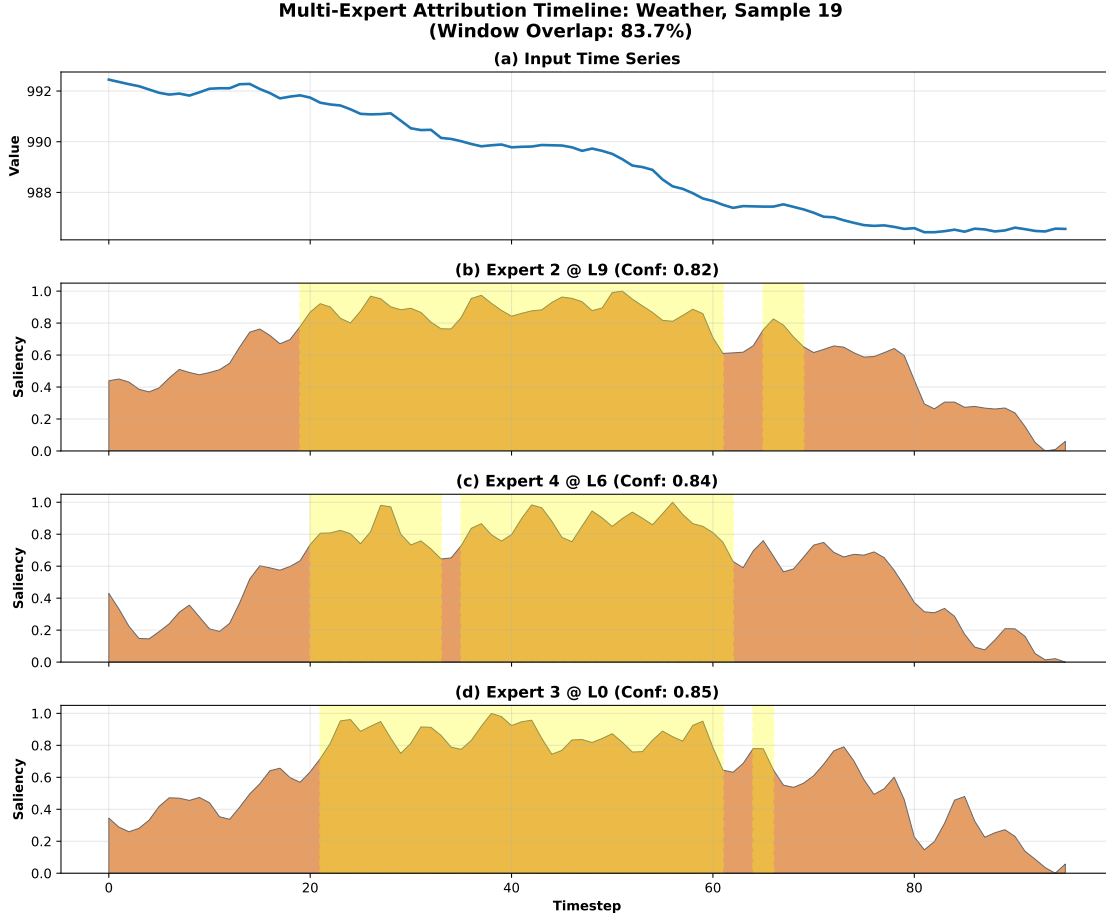


Figure 3: **Visualizing Ensemble Consensus.** (a) Input time series from the Weather dataset. (b)–(d) Saliency maps for three experts at different layers (L9, L6, L0) analyzing the same input. The orange filled area shows saliency magnitude at each timestep; yellow bands highlight the critical windows identified by each expert. Experts converge on overlapping temporal regions (83.7% pairwise overlap) with high confidence (0.82–0.85), demonstrating that they do not partition the timeline but instead independently identify similar structural anchors. This redundancy provides a robust “verification committee” against input noise.

anchors, creating a consensus that is impervious to superficial noise but responsive to fundamental signal changes.

## 5.6 Cross-Architecture Analysis

To assess whether our findings generalize beyond Time-MoE, we conduct three analyses: (1) scaling comparison across Time-MoE model sizes, (2) routing ablations that replace learned routing with random or uniform expert selection, and (3) validation on Moirai-MoE Liu et al. (2024a), a different MoE architecture for time series forecasting that uses distance-based routing rather than learned gates.

We evaluate Time-MoE-200M ( $4\times$  parameters, 24 layers) alongside the 50M variant, and test Moirai-MoE-Small (117M parameters, 6 layers, 32 experts per layer). These models employ different routing mechanisms: Time-MoE uses *learned gates*, where router networks are trained end-to-end to produce expert selection logits Shazeer et al. (2017), while Moirai-MoE uses *distance-based routing*, selecting experts based on similarity between input representations and learned expert centroids Liu et al. (2024a). For routing ablations, we replace learned routing decisions at inference time with either random expert selection or uniform weighting

across all experts. All experiments use identical evaluation protocols with 300 expert-sample pairs across five datasets. Table 5 presents the results.

Table 5: Cross-architecture comparison. Saliency IoU measures expert attention overlap (higher indicates stronger consensus). Stability columns show the percentage of inputs where routing decisions changed under perturbation.

Model	Routing	IoU	Gentle Flip (%)	Aggr. Flip (%)
Time-MoE-50M	Learned	0.673	0.3	99.7
Time-MoE-200M	Learned	<b>0.804</b>	0.5	99.5
Moirai-MoE	Distance	0.417	100.0*	100.0
<i>Random</i>	Stochastic	0.000 <sup>†</sup>	–	–
<i>Uniform</i>	Fixed	0.000 <sup>†</sup>	–	–

<sup>†</sup>Ablated routing yields no consistent saliency pattern.

\*High frequency, low severity: routing changes on every input, but affects only  $\sim 20\%$  of experts selected per input.

**Ensemble Consensus is a general MoE property.** Both Time-MoE and Moirai-MoE exhibit Ensemble Consensus, with experts converging on similar temporal regions ( $\text{IoU} > 0.4$ ). This confirms that consensus is not an artifact of Time-MoE’s specific architecture. However, the strength of consensus varies: Time-MoE’s learned gates produce stronger agreement ( $\text{IoU}$  0.67–0.80) than Moirai-MoE’s distance-based routing ( $\text{IoU}$  0.42). Notably, consensus strengthens with scale: the 200M model shows 19% higher IoU than the 50M model, suggesting that larger models develop more coherent expert behavior.

**Ensemble Consensus requires learned routing.** Ablating learned routing eliminates consensus. When we replace Time-MoE’s learned gates with random or uniform expert selection, saliency IoU drops from 0.673 to 0.000, indicating a complete absence of agreement between experts. This result serves as a negative control: it confirms that high IoU is not merely a result of experts reacting to prominent input features, but is rather a coordinated behavior orchestrated by the learned router. Experts no longer converge on which temporal regions matter because routing decisions become arbitrary. This demonstrates that Ensemble Consensus emerges from learned routing behavior, not from the MoE architecture itself.

**Routing mechanism dictates stability profile.** Time-MoE exhibits a distinct Stability Gap: gentle perturbations rarely change routing ( $< 1\%$ ), while aggressive perturbations succeed ( $> 99\%$ ). In contrast, Moirai-MoE is highly sensitive, with routing changes occurring in 100% of inputs even under gentle noise. This divergence isolates the impact of the routing mechanism: Time-MoE’s *learned gates* develop a decision margin that filters noise, whereas Moirai-MoE’s *distance-based routing* (comparing queries to centroids) lacks this margin and responds to minute input variations. Note that while Moirai-MoE exhibits high *frequency* of change (100% of inputs), the *severity* is low: only  $\sim 20\%$  of individual expert selections change per input, compared to Time-MoE’s near-zero change.

## 6 Discussion

Our experimental validation reveals properties of Time-MoE’s routing architecture that have broader implications for understanding MoE models in time series forecasting.

**Ensemble Consensus vs. Expert Specialization.** A common assumption in MoE interpretability is that experts develop distinct specializations, such as one expert for trends and another for seasonality. Our saliency IoU analysis (mean 0.677) challenges this assumption: experts at different layers converge on the same temporal windows rather than partitioning the input space. This “verification committee” architecture may explain MoE robustness: multiple independent assessors attending to the same features provide a stable consensus even when individual assessments are noisy. This redundancy functions as a learned variance

reduction mechanism; in noisy time series, multiple experts verifying the same temporal patterns allows the model to average out individual errors, producing more stable routing decisions.

**The Stability Gap.** The 300-fold difference between gentle (0.3%) and aggressive (99.7%) perturbation success rates reveals that Time-MoE’s routing is robust. The router ignores superficial signal modifications and responds only when structural anchors, the patterns identified by saliency maps, are destroyed. This robustness is desirable for deployment but creates challenges for counterfactual-based explainability: meaningful “what-if” questions require substantial input modifications.

**Distributed Focus Scores.** The mean focus score of 0.563 indicates distributed attention across temporal windows. Rather than interpreting this as “low confidence,” we understand it as reflecting the consensus architecture: when multiple experts attend to overlapping regions, the combined saliency is naturally distributed rather than concentrated on a single peak.

**Generalizability across architectures.** Our cross-architecture analysis reveals that Ensemble Consensus generalizes to other MoE models, though its strength varies with the routing mechanism. This finding suggests that consensus may be an emergent property of sparse expert selection under learned routing, rather than specific to Time-MoE’s design. The absence of a Stability Gap in Moirai-MoE indicates that routing robustness is not inherent to MoE architectures but develops when gates learn to identify invariant signal features. These results have practical implications: practitioners seeking robust routing behavior should prefer learned gate mechanisms over distance-based approaches.

**Implications for MoE Interpretability.** Our findings suggest several reconsiderations for MoE explainability. First, interpretability frameworks should not assume expert specialization; our results indicate that MoE routing may rely on ensemble consensus rather than task partitioning. Second, for highly stable routing architectures, counterfactual methods face a tension: gentle perturbations that preserve signal semantics do not alter routing, while aggressive perturbations that change routing may destroy meaningful information. Third, rather than interpreting saliency maps as “importance scores,” our framework positions them as identifying structural anchors, regions whose integrity is essential for routing stability.

**Limitations.** Several limitations merit acknowledgment. Our cross-architecture comparison shows that Moirai-MoE differs from Time-MoE in multiple dimensions (routing mechanism, number of experts, layer count), making it difficult to isolate which factor causes the observed differences in stability profiles. Future work could examine models that vary only in routing mechanism while holding other factors constant. Additionally, the perturbations that successfully alter routing (zero-out, noise injection) may not correspond to semantically meaningful input modifications; future work could explore optimization-based counterfactuals that minimize semantic distance while achieving routing change. Finally, our analysis operates at the input level; some routing decisions may be better explained by intermediate representations, and extending the framework to layer-wise attribution could address this limitation.

## 7 Conclusion and Future Work

We have presented RPATH, a post-hoc explainability framework for time series Mixture-of-Experts models that combines temporal saliency mapping, counterfactual generation, and uncertainty quantification. Our approach operates on frozen models without requiring gradient access or architectural modifications.

Experimental validation on Time-MoE-50M across 300 expert-sample pairs reveals two properties of the architecture. First, we demonstrate *Ensemble Consensus*: experts at different layers independently converge on the same critical temporal windows (mean saliency IoU 0.677), challenging the assumption that MoE models achieve performance through distinct expert specializations. Second, we identify *Structural Robustness* through a Stability Gap, where gentle perturbations alter routing in only 0.3% of cases while aggressive perturbations succeed in 99.7%, indicating that routing decisions reflect structural anchors rather than superficial signal characteristics.

Together, these findings demonstrate that Time-MoE’s performance stems from *Ensemble Redundancy*: multiple experts verify the same structural anchors, providing a robust consensus that is insensitive to noise but responsive to meaningful signal changes. Cross-architecture analysis on Moirai-MoE confirms that Ensemble Consensus is a general property of MoE models, though learned gate routing produces stronger consensus than distance-based approaches. The Stability Gap, however, appears specific to learned routing mechanisms. Our framework provides practitioners with tools to visualize expert attention, identify critical input regions, and quantify routing stability.

Several research directions emerge from this work. Understanding how ensemble consensus emerges during training could inform MoE architecture design and whether explicit consensus regularization improves robustness. If routing depends primarily on structural anchors, models might be compressed by reducing expert redundancy while preserving consensus coverage. Finally, the stability gap suggests a path toward certified routing robustness: inputs where gentle perturbations never alter routing may be flagged as high-confidence predictions. We provide complete open-source implementation to facilitate further research and practical applications in time series forecasting.

## Broader Impact Statement

## Author Contributions

## Acknowledgments

## References

- Jacob Andreas, Marcus Rohrbach, Trevor Darrell, and Dan Klein. Neural module networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 39–48, 2016.
- Yonatan Belinkov. Probing classifiers: Promises, shortcomings, and advances. *Computational Linguistics*, 48(1):207–219, 2022. doi: 10.1162/coli\_a\_00422.
- Zixuan Chen, Xiaolin Wang, and Yu Zhang. Tf-lime: Interpretation method for time-series models based on time-frequency features. *Sensors*, 25(9):2845, 2025. doi: 10.3390/s25092845.
- Susanne Dandl, Christoph Molnar, Martin Binder, and Bernd Bischl. Multi-objective counterfactual explanations. In *Parallel Problem Solving from Nature – PPSN XVI*, pp. 448–469. Springer, 2020.
- Nan Du, Yanping Huang, Andrew M Dai, Simon Tong, Dmitry Lepikhin, Yuanzhong Xu, Maxim Krikun, Yanqi Zhou, Adams Wei Yu, Orhan Firat, et al. GLaM: Efficient scaling of language models with mixture-of-experts. In *International Conference on Machine Learning*, pp. 5547–5569. PMLR, 2022.
- William Fedus, Barret Zoph, and Noam Shazeer. Switch transformers: Scaling to trillion parameter models with simple and efficient sparsity. *Journal of Machine Learning Research*, 23(120):1–39, 2022.
- Yash Goyal, Ziyang Wu, Jan Ernst, Dhruv Batra, Devi Parikh, and Stefan Lee. Counterfactual visual explanations. In *International Conference on Machine Learning*, pp. 2376–2384. PMLR, 2019.
- Ronghang Hu, Jacob Andreas, Marcus Rohrbach, Trevor Darrell, and Kate Saenko. Learning to reason: End-to-end module networks for visual question answering. In *Proceedings of the IEEE International Conference on Computer Vision*, pp. 804–813, 2017.
- Sarthak Jain and Byron C Wallace. Attention is not explanation. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pp. 3543–3556, 2019.
- Been Kim, Martin Wattenberg, Justin Gilmer, Carrie Cai, James Wexler, Fernanda Viegas, and Rory Sayres. Interpretability beyond feature attribution: Quantitative testing with concept activation vectors (TCAV). In *International Conference on Machine Learning*, pp. 2668–2677. PMLR, 2018.

- Pang Wei Koh, Thao Nguyen, Yew Siang Tang, Stephen Mussmann, Emma Pierson, Been Kim, and Percy Liang. Concept bottleneck models. In *International Conference on Machine Learning*, pp. 5338–5348. PMLR, 2020.
- Peiyu Li, Omar Bahri, Soukaina Filali Boubrahimi, and Shah Muhammad Hamdi. M-cels: Counterfactual explanation for multivariate time series data guided by learned saliency maps. In *2024 International Conference on Machine Learning and Applications (ICMLA)*, pp. 713–718. IEEE, 2024a.
- Peiyu Li, Omar Bahri, Pouya Hosseinzadeh, Soukaina Filali Boubrahimi, and Shah Muhammad Hamdi. Info-cels: Informative saliency map guided counterfactual explanation. *arXiv preprint arXiv:2410.20539*, 2024b.
- Xu Liu, Juncheng Liu, Gerald Woo, Taha Aksu, Yuxuan Liang, Roger Zimmermann, Chenghao Liu, Silvio Savarese, Caiming Xiong, and Doyen Sahoo. Moirai-moe: Empowering time series foundation models with sparse mixture of experts. *arXiv preprint arXiv:2410.10469*, 2024a.
- Zichuan Liu, Yingying Zhang, Tianchun Wang, Zelong Wang, Yitian Chen, Zhu Zhao, and Wei Wang. Explaining time series via contrastive and locally sparse perturbations. In *International Conference on Learning Representations*, 2024b.
- Ka Man Lo, Zeyu Huang, Zihan Qiu, Zili Wang, and Jie Fu. A closer look into mixture-of-experts in large language models. *arXiv preprint arXiv:2406.18219*, 2024. URL <https://arxiv.org/abs/2406.18219>. NAACL 2025 Findings.
- Scott M Lundberg and Su-In Lee. A unified approach to interpreting model predictions. In *Advances in Neural Information Processing Systems*, volume 30, pp. 4765–4774, 2017.
- Geoffrey J McLachlan and David Peel. *Finite Mixture Models*. John Wiley & Sons, New York, 2000.
- Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin. Why should i trust you? explaining the predictions of any classifier. In *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pp. 1135–1144, 2016.
- Alexis Ross, Ana Marasović, and Matthew E Peters. Explaining NLP models via minimal contrastive editing (MiCE). In *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*, pp. 3840–3852, 2021.
- Sofia Serrano and Noah A Smith. Is attention interpretable? In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pp. 2931–2951, 2019.
- Noam Shazeer, Azalia Mirhoseini, Krzysztof Maziarczyk, Andy Davis, Quoc Le, Geoffrey Hinton, and Jeff Dean. Outrageously large neural networks: The sparsely-gated mixture-of-experts layer. In *International Conference on Learning Representations*, 2017.
- Xiaoming Shi, Shiyu Ma, Haitao Cheng, Zhicheng Wu, Kenny Q Zhu, and Jiang Bian. Time-MoE: Billion-scale time series foundation models with mixture of experts. *arXiv preprint arXiv:2409.16040*, 2024. URL <https://arxiv.org/abs/2409.16040>.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Advances in Neural Information Processing Systems*, volume 30, pp. 5998–6008, 2017.
- Sandra Wachter, Brent Mittelstadt, and Chris Russell. Counterfactual explanations without opening the black box: Automated decisions and the GDPR. *Harvard Journal of Law & Technology*, 31(2):841–887, 2018.