
Interpretability-by-Design with Accurate Locally Additive Models and Conditional Feature Effects

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Generalized additive models (GAMs) offer interpretability through independent
2 univariate feature effects but underfit when interactions are present in data. GA^2Ms
3 add selected pairwise interactions which improves accuracy, but sacrifices inter-
4 pretability and limits model auditing. We propose *Conditionally Additive Local*
5 *Models* (CALMs), a new model class, that balances the interpretability of GAMs
6 with the accuracy of GA^2Ms . CALMs allow multiple univariate shape functions
7 per feature, each active in different regions of the input space. These regions are
8 defined independently for each feature as simple logical conditions (thresholds)
9 on the features it interacts with. As a result, effects remain locally additive while
10 varying across subregions to capture interactions. We further propose a principled
11 distillation-based training pipeline that identifies homogeneous regions with lim-
12 ited interactions and fits interpretable shape functions via region-aware backfitting.
13 Experiments on diverse classification and regression tasks show that CALMs con-
14 sistentlly outperform GAMs and achieve accuracy broadly comparable to GA^2Ms ,
15 while preserving the univariate auditability that GA^2Ms forfeit. Overall, CALMs
16 offer a favorable trade-off between predictive accuracy and interpretability.

17 1 Introduction

18 In high-stakes decision making, machine learning models must provide not only reliable predictions
19 but also human-understandable explanations [31, 10]. This requirement has motivated the develop-
20 ment of *interpretable-by-design* models, whose structure permits direct inspection of their behavior
21 [30]. However, interpretability is a continuum rather than a binary property: interpretable-by-design
22 models span a spectrum where gains in predictive accuracy often entail greater structural complexity
23 [38, 27].

24 This trade-off is particularly evident in GAMs [28] and their extension, GA^2Ms [7]. GAMs achieve
25 high interpretability by modeling predictions as a sum of independent univariate effects, but this
26 strict additivity prevents them from capturing feature interactions, limiting their predictive accuracy.
27 GA^2Ms add selected pairwise interactions, improving performance at the cost of crucial interpretabil-
28 ity properties [35]. Specifically, interaction terms (i) obscure the unique attribution of the prediction
29 to individual features because the pairwise interaction terms are generally not additively separable
30 (ii) complicate global auditing by requiring the simultaneous inspection of multiple (univariate and
31 bivariate) effects. Section 3.3 analyzes these limitations in detail.

32 We introduce *Conditionally Additive Local Models* (CALMs), a new model class that strikes a balance
33 between the predictive accuracy of GA^2Ms and the interpretability of GAMs. The core idea is to learn
34 *conditional feature effects*: multiple univariate effects per feature, each active within a distinct region
35 of the input space where the feature exhibits nearly additive behavior—that is, where it minimally
36 interacts with others. These regions are defined through threshold-based conditions on interacting

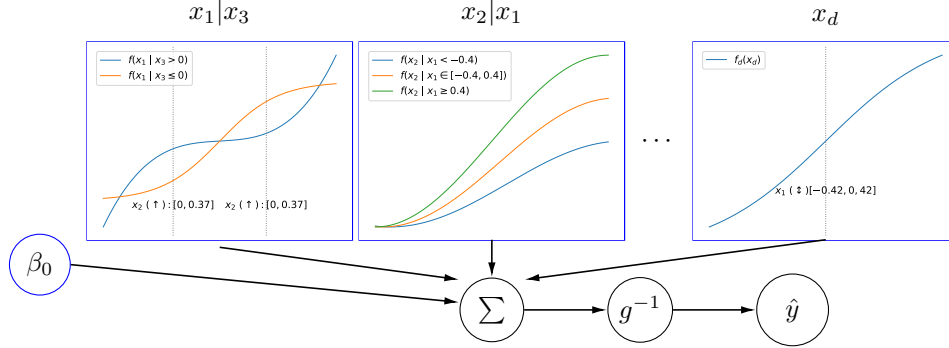


Figure 1: CALMs use *conditional feature effects*: every feature effect is expressed by a collection of 1D functions, each associated with a different region of the input space. The regions are defined conditioning on the interacting features. In the example, the effect of x_1 conditions on x_3 , the effect of x_2 on x_1 , while x_d does not interact with any other feature and thus has a single plot. CALMs are *interpretable-by-design*—summarized in d figures of 1D plots—and *accurate*, as they can model feature interactions.

37 features. For example, if x_i interacts with x_j , a CALM may learn separate effects for x_i ; one when
 38 $x_j < \tau$ and another when $x_j \geq \tau$ for some threshold τ . Figure 1 illustrates this representation.
 39 Conditional effects thus enable interaction-aware modeling while maintaining (i) transparent feature
 40 contributions and (ii) straightforward global model auditing.

41 To fit CALMs to data, we propose a three-step distillation-based pipeline: (1) Train a black-box
 42 reference model to capture complex interactions present in the data; (2) Partition the input space
 43 independently for each feature using CART-based splitters [17] optimized to minimize heterogeneity—
 44 a criterion that, when minimized, provably reduces feature interactions within the resulting
 45 regions [19]; (3) Fit region-specific shape functions using a region-aware extension of the standard
 46 backfitting algorithm. This procedure efficiently identifies interacting features (e.g., x_i interacts
 47 with x_j), optimal split points (e.g., τ such that separate effects are learned for $x_j < \tau$ and $x_j \geq \tau$),
 48 and requires minimal tuning—primarily the maximum allowed number of interactions per feature.
 49 Extensive evaluation on diverse regression and classification benchmarks demonstrates that CALMs
 50 consistently outperform GAMs in predictive accuracy, and often match GA²Ms. Furthermore, we
 51 formally show that CALMs satisfy key interpretability properties unsupported by GA²Ms. Code for
 52 reproducing all experiments is provided in the supplementary material.

53 **Contributions.** (i) We introduce *Conditionally Additive Local Models* (CALMs), a novel
 54 interpretable-by-design model class that balances GAM’s interpretability and GA²Ms accuracy
 55 via conditional feature effects. (ii) We develop a robust training algorithm based on distillation,
 56 heterogeneity-minimization and region-aware backfitting. (iii) We provide a formal analysis of
 57 the interpretability properties inherent in CALMs. (iv) We present extensive empirical evidence
 58 demonstrating improved accuracy–interpretability trade-offs over GAMs and GA²Ms. (v) A within-
 59 subjects user study (N=25) suggests that non-expert users interpret CALM plots more accurately than
 60 heatmap-based bivariate-interaction visualizations (Appendix D).

61 2 Background and Related Work

62 Interpretable-by-design models enforce transparency in their decision-making process. Examples
 63 include decision trees [4], rule lists [26, 2], and prototype-based classifiers [9, 3]. Among them,
 64 GAMs [16, 42] stand as a particularly popular candidate, mainly due to their simple global inter-
 65 pretability.

66 GAMs model the output as $g(\mathbb{E}[y|\mathbf{x}]) = \beta_0 + \sum_i f_i(x_i)$, where each f_i is a univariate shape function,
 67 β_0 is the global intercept, and g is a link function. The key benefit of GAMs is that feature effects
 68 can be *independently* visualized through a simple 1D plot. Methods for learning f_i for all i include
 69 spline-based approaches [28, 42], gradient boosting [11], and neural-based variants [1, 25, 35].
 70 However, standard GAMs assume that features contribute independently to the output, which limits

71 their accuracy when feature interactions are present in data. To capture these interactions, GA²Ms add
 72 pairwise terms: $g(\mathbb{E}[y|\mathbf{x}]) = \beta_0 + \sum_i f_i(x_i) + \sum_{i<j} f_{ij}(x_i, x_j)$ [29]. GA²Ms instances vary
 73 in (i) how they select the top- K interactions to maintain sparsity and (ii) how they learn f_i and
 74 f_{ij} . For example, [29] greedily selects interactions based on loss reduction, while, neural-based
 75 approaches [43, 23, 8] integrate pairwise interactions into standard neural-based architectures.

76 While GA²Ms improve accuracy, they obscure individual feature contributions and complicate model
 77 auditing (see Section 3.3). A key open question remains: *Can we approach GA²M-level accuracy*
 78 *while maintaining GAM-level interpretability?* To this end, we draw inspiration from *regional*
 79 *effect methods*, which handle interactions by partitioning the feature space into subregions where
 80 interactions are weak. To understand regional effects, consider first how standard global effects work.
 81 Global effect plots, such as PDP, ALE or SHAP-DP, explain a black box model by decomposing its
 82 complex d -dimensional function $f(\mathbf{x}) \rightarrow y$ into d one-dimensional plots $x_i \rightarrow y$. However, when
 83 strong interactions exist between x_i and other features, these plots can yield misleading explanations
 84 [12, 13]. Regional feature effect plots address that by partitioning the feature space into subregions
 85 where interactions are minimal [17, 19]. Within each subregion, simple univariate plots accurately
 86 represent feature effects without being confounded by interactions with other features. We adopt this
 87 strategy to identify low-interaction subregions and then we fit shape functions within each subregion.

88 Our approach relates to model distillation [20, 40] and surrogate modeling [15]. These approaches
 89 train an interpretable “student” model to mimic a complex “teacher”, either locally [36] or globally [15].
 90 However, instead of explaining the teacher, we use it to identify subregions with minimal
 91 feature interactions, on which we then fit shape functions. In a rough analogy, CALM serves as the
 92 interpretable-by-design “student” that replaces the black-box “teacher” as the final predictor. Other
 93 works that share the above idea are: SLIM [22] and related analyses [18] propose tree-based models
 94 with simple predictors in each region, mainly in the context of model distillation and explanation. In
 95 contrast, our goal is not to approximate a black-box model, but to construct an interpretable-by-design
 96 predictor with GAM-like structure. Related ideas also appear in [14], which leverages regional
 97 feature effects to build interpretable models; however, that work does not explicitly characterize the
 98 interpretability guarantees of the learned model or evaluate them systematically.

99 3 CALM: Conditionally Additive Local Model

100 CALM captures feature interactions via *conditional feature effects*, a set of univariate shape functions
 101 per feature, each active in a different region of the input space.

102 3.1 Model Formulation

103 Let $\mathbf{X} = (X_1, \dots, X_d) \in \mathcal{X} \subseteq \mathbb{R}^d$ be a random vector with joint distribution $P_{\mathbf{X}}$ over the input data,
 104 $\mathbf{x} = (x_1, \dots, x_d)$ where $\mathcal{X} = \mathcal{X}_1 \times \dots \times \mathcal{X}_d$. We use the subscript $-i$ to denote quantities excluding
 105 the i -th feature (e.g., \mathbf{X}_{-i} defined on space \mathcal{X}_{-i}). Let also Y be the output random variable, with
 106 $Y \in \mathbb{R}$ for regression and $Y \in \{0, 1\}$ for binary classification. CALM is then defined as

$$g(\mathbb{E}[Y | \mathbf{X} = \mathbf{x}]) = \beta_0 + \sum_{i=1}^d f_i^{(r_i(\mathbf{x}_{-i}))}(x_i) \quad (1)$$

107 where $\beta_0 \in \mathbb{R}$ is an intercept and g is a link function (identity for regression and logit for classification).
 108 For each feature i , CALM learns a set of univariate shape functions $\{f_i^{(r)}\}_{r=1}^{R_i}$. At prediction time,
 109 a region selection function $r_i(\mathbf{x}_{-i}) : \mathbb{R}^{d-1} \rightarrow \{1, \dots, R_i\}$ chooses the specific shape function for
 110 x_i based on the context of other features \mathbf{x}_{-i} . This allows the effect of x_i to vary across regions,
 111 effectively modeling feature interactions, while maintaining univariate interpretability within each
 112 region. The final prediction is given by:

$$f_{\text{CALM}}(\mathbf{x}) = g^{-1} \left(\beta_0 + \sum_{i=1}^d f_i^{(r_i(\mathbf{x}_{-i}))}(x_i) \right) \quad (2)$$

113 For brevity, we use $\hat{y}(\mathbf{x})$ to denote the prediction $f_{\text{CALM}}(\mathbf{x})$. As in standard GAMs, identifiability
 114 requires region-wise centering of each shape function; we adopt the centering constraint stated in
 115 Appendix A.3.1 (Eq. 15). The selection function $r_i(\mathbf{x}_{-i})$ partitions the input space (excluding x_i)

Algorithm 1 Training a CALM model

Require: Training data $\mathcal{D} = \{(\mathbf{x}^{(i)}, y^{(i)})\}_{i=1}^N$

Output: CALM predictor $f_{\text{CALM}}(\mathbf{x})$

- 1: **Step 1:** Train reference model f_{ref} on \mathcal{D} .
 - 2: **Step 2:** Learn feature-specific trees $T_i, i = 1, \dots, d$.
 - 3: **Step 3:** Estimate shape functions $\{f_i^{(r)}\}_{r=1}^{R_i}\}_{i=1}^d$.
-

116 into $\mathcal{P}_i := \{\mathcal{R}_i^{(r)}\}_{r=1}^{R_i}$, independently for each feature i , where each $\mathcal{R}_i^{(r)} \subseteq \mathcal{X}_{-i}$. The goal of these
117 partitions is to minimize interactions between x_i and other features within each region, ensuring the
118 effect of x_i is accurately captured by a single univariate shape function. To maintain interpretability,
119 each partition is represented by a binary decision tree T_i of maximum depth d_{max} , built over \mathbf{x}_{-i} .
120 Internal nodes apply axis-aligned splits, i.e., inequality thresholds $x_j < \tau$ or $x_j \geq \tau$ for continuous
121 features, and equality tests $x_j = \tau$ or $x_j \neq \tau$ for categorical ones. Each leaf of the tree defines a
122 region $\mathcal{R}_i^{(r)}$. Formally, each region is defined by a conjunction of at most $m_i^{(r)} \leq d_{\text{max}}$ rules:

$$\mathcal{R}_i^{(r)} = \left\{ \mathbf{x}_{-i} \mid \bigwedge_{k=1}^{m_i^{(r)}} (x_{j_k} \text{ op}_k \tau_k) \right\}. \quad (3)$$

123 where $j_k \in \{1, \dots, d\} \setminus \{i\}$ indexes an interacting feature, $\text{op}_k \in \{<, \geq, =, \neq\}$ is a comparison
124 operator, and $\tau_k \in \mathbb{R}$ is a threshold.

125 CALM is able to capture feature interactions involving up to $(d_{\text{max}} + 1)$ features, as each region
126 conditions on up to d_{max} features. While increasing d_{max} can improve accuracy, it comes at
127 the cost of interpretability. To balance this accuracy-interpretability tradeoff, CALM exposes two
128 hyperparameters: d_{max} (tree depth) bounds the number of shape functions per feature to $R_i \leq 2^{d_{\text{max}}}$,
129 while K limits the total number of shape functions across all features via $\sum_i R_i \leq K$. In our
130 experiments, we set $d_{\text{max}} = 2$ (yielding $R_i = 4$ regions per feature) and leave K unconstrained.

131 Notably, CALM learns, at most, $d \cdot 2^{d_{\text{max}}}$ univariate shape functions (up to $2^{d_{\text{max}}}$ per feature), yet
132 these combine to express up to $2^{d \cdot d_{\text{max}}}$ distinct additive models—an exponential increase in expressive
133 power.

134 3.2 Training algorithm for fitting CALM

135 Given a dataset $\mathcal{D} = \{(\mathbf{x}^{(i)}, y^{(i)})\}_{i=1}^N$, we fit a CALM predictor by minimizing the empirical
136 risk $\mathbb{E}_{(\mathbf{X}, Y) \sim \hat{P}_{\mathcal{D}}} [\mathcal{L}(Y, \hat{y}(\mathbf{X}))]$, where $\hat{y}(\mathbf{X})$ denotes a CALM model, \mathcal{L} a loss function and $\hat{P}_{\mathcal{D}}$ the
137 empirical distribution over \mathcal{D} . Fitting a CALM requires estimating: (a) a set of d partitions \mathcal{P}_i ,
138 $i \in \{1, \dots, d\}$, represented by binary trees T_i (one per feature); (b) region-specific shape functions
139 $\{f_i^{(r)}\}_{r=1}^{R_i}$ and a global intercept β_0 . We propose a three-step distillation-based pipeline summarized
140 in Algorithm 1.

141 **Step 1: Train a Reference Model.** We train a high-capacity black-box predictor $f_{\text{ref}} : \mathcal{X} \rightarrow \mathbb{R}$ on
142 \mathcal{D} . This model serves as functional proxy for $\mathbb{E}[Y|\mathbf{X}]$ and is used exclusively to detect interactions.
143 The reference model is discarded after training. The choice of f_{ref} is independent of later stages; any
144 accurate predictor can be used, such as gradient-boosted trees (our default), neural networks, random
145 forests or foundation models, like TabPFN [21].

146 **Step 2: Learn feature-specific partitioning trees.** This step learns d independent, feature-specific
147 partitions $\mathcal{P}_i := \{\mathcal{R}_i^{(r)}\}_{r=1}^{R_i}$, each represented by a binary tree T_i defined over \mathbf{x}_{-i} . The objective is
148 to identify near-additive regions $\mathcal{R}_i^{(r)}$, in which the effect of x_i on the output is well approximated by
149 a univariate function $f_i^{(r)}$, due to locally weak higher-order interactions.

150 To identify such regions, we use the interaction-related heterogeneity measure from [19]. Feature
151 effect methods (e.g., PDP or ALE) explain a black-box model (like $f_{\text{ref}}(\mathbf{x})$) by decomposing it
152 into univariate effects $f_i(x_i)$ for all i . To do so, they first define the local effects $h(x_i, \mathbf{x}_{-i}^{(j)})$ that
153 quantify the contribution of feature x_i on a specific instance $\mathbf{x}^{(j)}$. Then they compute $f_i(x_i)$ by
154 averaging the local effects. Heterogeneity is the variability of these local effects around their average,
155 directly measuring how much x_i 's effect depends on other features. High heterogeneity signals strong

156 interactions, while low heterogeneity identifies near-additive, interaction-free regions—making it an
 157 effective criterion for our purpose.

158 Following the above, we define the *pointwise* heterogeneity of x_i in a region \mathcal{R} (e.g., $\mathcal{R}_i^{(r)}$) as:

$$H_i^{\mathcal{R}}(x_i) = \mathbb{E}_{\mathbf{X}_{-i} | \mathbf{X}_{-i} \in \mathcal{R}} \left[(h(x_i, \mathbf{X}_{-i}) - \mu_i^{\mathcal{R}}(x_i))^2 \right] \quad (4)$$

159 where $\mu_i^{\mathcal{R}}(x_i) = \mathbb{E}_{\mathbf{X}_{-i} | \mathbf{X}_{-i} \in \mathcal{R}} [h(x_i, \mathbf{X}_{-i})]$. The corresponding *feature-level heterogeneity* is:

$$H_i^{\mathcal{R}} = \mathbb{E}_{X_i | \mathbf{X}_{-i} \in \mathcal{R}} [H_i^{\mathcal{R}}(X_i)] \quad (5)$$

160 In our experiments, heterogeneity is computed using PDP-based formulas; ALE- and SHAP-DP
 161 variants are valid alternatives [17, 19]. While PDP can suffer from extrapolation bias under feature
 162 correlations, here it is used only for region detection, not for prediction: even biased local effects
 163 yield elevated heterogeneity when genuine interactions are present. When correlations are strong,
 164 RHALE [13] is a drop-in alternative, with comparable results in Appendix C. A related caveat is that
 165 an interaction with one feature can also raise heterogeneity along its correlated proxies, so CALM
 166 may split on either; both choices recover the interaction structure, though the conditioning labels
 167 may not be uniquely identifiable under strong collinearity, a limitation shared by partition-based
 168 interaction detectors such as REPID.

169 As additional theoretical justification for using heterogeneity as the splitting criterion, we show
 170 that, under specific assumptions, the approximation error of CALM is bounded by the sum of the
 171 expected feature heterogeneities. Consequently, reducing heterogeneity of each feature through
 172 partitioning yields a CALM model f_{CALM} with lower mean squared error, assuming perfect fit in Step
 173 3 of Algorithm 1. The proof is provided in Appendix A (Proposition A.1).

174 Therefore, we use $H_i^{\mathcal{R}}$ as the splitting criterion and learn a binary decision tree T_i for each feature x_i
 175 using a greedy CART-style algorithm. At each node, we consider splits over all conditioning features
 176 $x_j \neq x_i$ and select the split that maximizes the reduction in heterogeneity. For numerical features,
 177 splits take the form $x_j \{ \leq, > \} \tau$, while for categorical features we use $x_j \{ =, \neq \} \tau$. A split is accepted
 178 only if the relative reduction in $H_i^{\mathcal{R}}$ exceeds a threshold ϵ ; otherwise, the node becomes a leaf. This
 179 procedure yields a partition $\mathcal{P}_i = \{ \mathcal{R}_i^{(r)} \}_{r=1}^{R_i}$ where heterogeneity is significantly reduced, ensuring
 180 the effect of x_i is less distorted by interactions than in the global space. Additional implementation
 181 details are provided in Appendix A. In all experiments, we use PDP-based heterogeneity, set the
 182 maximum tree depth to $d_{\text{max}} = 2$, and fix $\epsilon = 0.2$. It is important to emphasize that this CART-style
 183 algorithm is only used to determine the subregions for each feature and is not used for the final
 184 prediction.

185 **Step 3: Estimate Shape Functions.** Given the partitions $\{ \mathcal{P}_i \}_{i=1}^d$, we estimate the region-specific
 186 shape functions $\{ f_i^{(r)} \}$ by minimizing the empirical loss of the CALM predictor. We use a modified
 187 gradient boosting procedure in which, at each iteration, a single shape function $f_i^{(r)}$ is updated using
 188 only the observations whose \mathbf{x}_{-i} fall into the corresponding region $\mathcal{R}_i^{(r)}$ (see Appendix A for details).

189 Although each update is restricted to a single region, the resulting optimization problem is inherently
 190 *coupled*. The regions are defined separately for each feature, so the subsets of samples used to update
 191 different shape functions generally overlap. Consequently, updating one $f_i^{(r)}$ changes the residuals
 192 seen by all other shape functions. Therefore it can be understood as a coordinated optimization of a
 193 single additive predictor with region-gated components.

194 The following proposition formalizes this intuition by characterizing the target of Step 3 at the
 195 population level and its convergence behavior under idealized updates. The proof is in Appendix A.
 196 In practice, Step 3 is implemented using gradient boosting to approximate the exact regional updates.

197 **Proposition 3.1** (Optimality and convergence). *Let $m(\mathbf{x}) = \mathbb{E}[Y | \mathbf{X} = \mathbf{x}]$ denote the true regression
 198 function. Assume regression with squared loss, fixed partition trees $\{ T_i \}_{i=1}^d$ (as learned by Step 2)
 199 and $\mathbb{E}[Y^2] < \infty$. Let $\mathcal{H}(\{ T_i \})$ denote the fixed-tree CALM class (Appendix A.3.1), and assume
 200 $\mathcal{H}(\{ T_i \}) \subset L_2(P_X)$ is nonempty, closed, and convex. Then:*

- 201 1. *Optimality: Any $s^* \in \arg \min_{s \in \mathcal{H}(\{ T_i \})} \mathbb{E}[(Y - s(\mathbf{X}))^2]$ is the $L_2(P_X)$ -best approximation
 202 of m within $\mathcal{H}(\{ T_i \})$.*
- 203 2. *Convergence: The idealized exact cyclic regional backfitting converges to an empirical risk
 204 minimizer over $\mathcal{H}(\{ T_i \})$.*

205 **3.3 Interpretability of a CALM**

206 CALMs offer interpretability similar to GAMs, as both require inspecting d univariate plots; one per
 207 feature. CALMs are slightly more complex: unlike GAMs, each plot can contain up to $2^{d_{\max}}$ curves
 208 corresponding to a different region of the input space and with interaction-induced discontinuities
 209 marked by vertical lines. However, because regions are interpretable (specified by simple threshold
 210 conditions), these region-specific curves provide explanations which are suitable for model auditing
 211 and decision support. For example: *The effect of age on mortality_rate follows this curve when*
 212 *the patient is male and has a BMI above 30.*

213 **Interpreting a CALM plot.**

214 In Figure 2 each curve gives the contribution of x_1 to y
 215 in a specific region; the blue curve when $x_3 > 0$ and the
 216 orange curve when $x_3 \leq 0$. For example, at $x_1 = -0.5$,
 217 the contribution is approximately -0.2 (blue) or -0.75
 218 (orange), depending on x_3 . The plots also illustrate how
 219 altering x_1 to $x_1 \rightarrow x_1 + \Delta x$ impacts the prediction.
 220 Vertical dotted lines mark points of a hidden discontinuity
 221 which is due to x_1 participating as an interaction term
 222 for feature x_2 . As shown in Figure 1, the effect of x_2
 223 is conditioned by $x_1 \leq -0.4$, $-0.4 < x_1 \leq 0.4$ and
 224 $x_1 > 0.4$, therefore in Figure 2 we observe vertical lines
 225 in $x_1 \pm 0.4$. If a change in x_1 does not cross a vertical line,
 226 the change in the output (Δy) equals the curve difference
 227 (Δf_i). Crossing a line signifies a hidden jump, in the range
 228 $[\alpha, \beta]$, so $\Delta f_i + \alpha \leq \Delta y \leq \Delta f_i + \beta$. Arrows provide
 229 a fast understanding of the jump: \uparrow means $\Delta y > \Delta f_i$,
 230 \downarrow means $\Delta y < \Delta f_i$, \updownarrow means it depends. Below, we
 231 outline three crucial interpretability properties, along with
 232 discussion about their satisfiability by GAM, GA²M, and
 233 CALM. Complete proofs are provided in Appendix B.

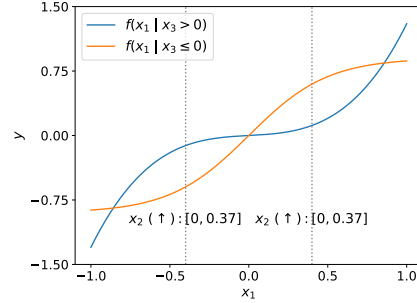


Figure 2: CALM plot for x_1 , from a synthetic function where x_1 interacts with x_2 and x_3 . Each curve shows the contribution of x_1 to y (P1) in a different region of the input space defined by $x_3 \leq 0$. Vertical lines indicate interaction-induced discontinuities: at $x_1 \approx \pm 0.4$, changing x_1 causes x_2 's region to switch, inducing a hidden jump of $[0, 0.37]$ that must be considered when assessing regional sensitivity (P2) or global properties (P3).

234 **P1. Local Feature Contribution:** What is the contribu-
 235 tion of each feature to the prediction?

236 *Formally:* Given an input \mathbf{x} , how much does each x_i con-
 237 tribute to $\hat{y}(\mathbf{x})$?

238 The explanation is *local*—it concerns a specific input \mathbf{x} .

239 In GAM, the contribution is $f_i(x_i)$. In GA²M, the contri-
 240 bution is neither explicit nor unique; it must be inferred post-hoc (e.g., via SHAP or LIME), with each
 241 method relying on different assumptions and yielding different results. In CALM, the contribution is
 242 $f_i^{(r(\mathbf{x}-i))}(x_i)$.

243 **P2. Regional Feature Sensitivity:** How does changing x_i change the prediction? *Formally:* Given
 244 x_i and $\Delta x > 0$, what is $\Delta \hat{y} = \hat{y}(\mathbf{x} + \mathbf{e}_i \Delta x) - \hat{y}(\mathbf{x})$, assuming only x_i is perturbed and \mathbf{x}_{-i} is fixed?
 245 The explanation is *regional*—it characterizes the effect of x_i on a specific region ($[x_i, x_i + \Delta x]$)
 246 *independently* of the values of other features. In GAM, the change is $\Delta \hat{y} = f_i(x_i + \Delta x) - f_i(x_i)$.
 247 In GA²M, the change cannot be determined without knowing the values of all features that interact
 248 with x_i . In CALM, an exact answer is possible only when the perturbation does not cross a vertical
 249 line. In this case, the resulting change is a set of values $\Delta y := \{\Delta f_i^{(r)}\}_{r=1}^{R_i}$, one for each curve in
 250 the plot: $\Delta f_i^{(r)} = f_i^{(r)}(x_i + \Delta x) - f_i^{(r)}(x_i)$. If a crossing occurs, an exact answer is not attainable,
 251 however, for less precise questions, such as whether the Δx change in x_i will have a positive impact
 252 on y , an answer is still feasible (see **P3.** below).

253 **P3. Global Feature Property:** Is the model globally monotonic increasing with respect to x_i ?

254 *Formally:* For all \mathbf{x} and $\Delta x > 0$, is $\Delta \hat{y} = \hat{y}(\mathbf{x} + \mathbf{e}_i \Delta x) - \hat{y}(\mathbf{x}) > 0$?

255 The explanation is *global*—it assesses the monotonicity of x_i across the entire input space. In GAM,
 256 monotonicity is easily determined by the shape of $f_i(x_i)$. In GA²M, verifying monotonicity requires
 257 a concurrent examination of the 1D shape of $f_i(x_i)$ along with all pairwise interactions f_{ij} along all
 258 j , an inspection which is infeasible for a human. In CALM, if no vertical lines are present, simply

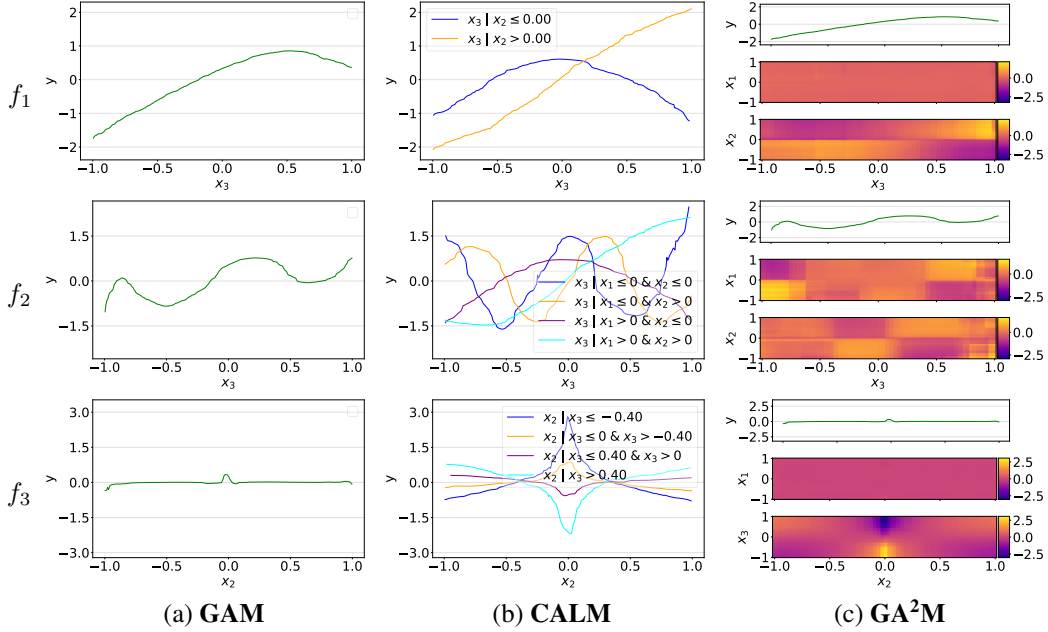


Figure 3: Explanatory plots for the three synthetic regression tasks: f_1 shows a two-region interaction; f_2 shows a four-region interaction; f_3 shows general interactions.

259 check whether all curves $f_i^{(r)}(x_i)$ for $r = 1, \dots, R_i$ are monotonically increasing. If vertical lines
 260 exist, simply verify that all arrows are positive.

261 3.4 Efficiency

262 CALM’s runtime is dominated by Step 2: fitting d binary trees costs $\mathcal{O}(d_{\max} d^2 N \log N)$, which
 263 reduces to $\mathcal{O}(d^2 N \log N)$ for the shallow trees we use—acceptable since d is on the order of tens
 264 for tabular data. Crucially, the local effects needed to score candidate splits are precomputed once
 265 and reused via indexed lookups, so each split evaluation avoids repeated black-box queries. Steps 1
 266 and 3 are standard model fits (XGBoost and gradient boosting in our default setup) and are fast in
 267 practice. Overall, CALM fits most tabular datasets within a few seconds; full runtimes are reported in
 268 Appendix C.

269 4 Empirical Evaluation

270 The empirical evaluation of CALM includes three synthetic regression datasets and 25 public real-
 271 world tabular datasets—10 for classification and 15 for regression. Across all experiments, CALM
 272 is applied with its default configuration: Step 1 uses XGBoost as the black-box model. Step 2 uses
 273 PDP-based heterogeneity with $d_{\max} = 2$, $\epsilon = 0.2$, and K remains unconstrained. Step 3 uses standard
 274 gradient boosting. All predictive results are reported as mean \pm standard deviation over standard
 275 5-fold cross-validation. Full experimental details are in Appendix C.

276 **Implementation details.** We used well-known open-source implementations for all baseline models:
 277 scikit-learn (random forests), XGBoost, tensorflow (neural networks and NAM), interpretml (EBM,
 278 EB²M), pygam (SPLINE), and official repositories for NODE-GA²M[8] and GAMI-Net[43].

279 4.1 Synthetic example

280 We compare CALM against EBM (as the GAM baseline) and EB²M (as the GA²M baseline) on three
 281 synthetic datasets. Both EBM and EB²M are used with their default parameters. Evaluation is based
 282 on R^2 performance. Each dataset contains 1000 samples with features drawn from $x_i \sim \mathcal{U}(-1, 1)$.
 283 Table 1 summarizes the accuracy of each approach.

Table 2: Classification accuracy vs. baselines (mean \pm std. dev. over 5-fold CV). \checkmark/\times denotes strictly higher/lower accuracy than CALM; = denotes exactly equal accuracy to CALM.

Dataset	BB		GAM		CALM	GA ² M		
	XGB	NAM	EBM	EB ² M		NODE	GAMI	
Adult	0.870 \pm 0.002 =	0.851 \pm 0.002 \times	0.870 \pm 0.002 =	0.870 \pm 0.002	0.871 \pm 0.003 \checkmark	0.850 \pm 0.002 \times	0.857 \pm 0.003 \times	
COMPAS	0.661 \pm 0.007 \times	0.682 \pm 0.016 \times	0.681 \pm 0.012 \times	0.684 \pm 0.013	0.684 \pm 0.011 =	0.667 \pm 0.013 \times	0.686 \pm 0.013 \checkmark	
HELOC	0.717 \pm 0.013 \times	0.723 \pm 0.011 \times	0.728 \pm 0.014 =	0.728 \pm 0.012	0.730 \pm 0.012 \checkmark	0.715 \pm 0.008 \times	0.725 \pm 0.012 \times	
MIMIC2	0.890 \pm 0.001 \checkmark	0.886 \pm 0.001 =	0.886 \pm 0.003 =	0.886 \pm 0.003	0.886 \pm 0.003 =	0.888 \pm 0.002 \checkmark	0.884 \pm 0.003 \times	
Appendicitis	0.868 \pm 0.069 \times	0.848 \pm 0.056 \times	0.877 \pm 0.077 \times	0.878 \pm 0.063	0.869 \pm 0.079 \times	0.849 \pm 0.016 \times	0.764 \pm 0.074 \times	
Phoneme	0.898 \pm 0.006 \checkmark	0.808 \pm 0.002 \times	0.821 \pm 0.007 \times	0.861 \pm 0.011	0.863 \pm 0.007 \checkmark	0.869 \pm 0.003 \checkmark	0.876 \pm 0.009 \checkmark	
SPECTF	0.862 \pm 0.029 \times	0.839 \pm 0.030 \times	0.894 \pm 0.015 =	0.894 \pm 0.015	0.882 \pm 0.017 \times	0.820 \pm 0.067 \times	0.874 \pm 0.016 \times	
Magic	0.885 \pm 0.004 \checkmark	0.850 \pm 0.006 \checkmark	0.857 \pm 0.005 \times	0.864 \pm 0.004	0.872 \pm 0.003 \checkmark	0.876 \pm 0.004 \checkmark	0.874 \pm 0.003 \checkmark	
Bank	0.908 \pm 0.003 \checkmark	0.901 \pm 0.003 \times	0.902 \pm 0.002 \times	0.905 \pm 0.002	0.909 \pm 0.002 \checkmark	0.903 \pm 0.001 \times	0.908 \pm 0.003 \checkmark	
Churn	0.958 \pm 0.004 \checkmark	0.885 \pm 0.009 \times	0.886 \pm 0.005 \times	0.946 \pm 0.003	0.957 \pm 0.009 \checkmark	0.954 \pm 0.009 \checkmark	0.952 \pm 0.007 \checkmark	
W/D/L	5/1/4	0/1/9	0/4/6		6/2/2	4/0/6	5/0/5	

284 **Case 1.** We define f_1 as $y = x_1^2 + \log(|x_2|) + 2 \sin(\frac{\pi}{2}x_3) \mathbf{1}_{x_2 \geq 0} + 2 \cos(\frac{\pi}{2}x_3) \mathbf{1}_{x_2 < 0}$. Since
285 GAM cannot model the $x_2 - x_3$ interaction, it simplifies the $x_3 \rightarrow y$ effect by averaging the
286 sine and cosine modes into one curve (Fig. 3a, f_1), resulting in $R^2 = 0.737$. Both GA²M and
287 CALM capture the interaction correctly, achieving near-perfect accuracy (0.974 and 0.995). However,
288 their interpretability differs. GA²M requires three views to interpret x_3 's effect: (i) the main effect
289 averaging the sine and cosine cases, (ii) a near-zero x_1-x_3 heatmap, and (iii) an x_2-x_3 heatmap
290 capturing the sign-dependent interaction (Fig. 3c, f_1). In contrast, CALM simplifies the interpretation
291 with two 1D plots and a clear separation of the regimes by conditioning on x_2 's sign (Fig. 3b, f_1).

292 **Case 2.** We define f_2 as: $y = x_1^2 + \log(|x_2|) + 2[\sin(\frac{\pi}{2}x_3) \mathbf{1}_{x_1 \geq 0, x_2 \geq 0} + \cos(\frac{\pi}{2}x_3) \mathbf{1}_{x_1 \geq 0, x_2 < 0} + \sin(2\pi x_3) \mathbf{1}_{x_1 < 0, x_2 \geq 0} + \cos(2\pi x_3) \mathbf{1}_{x_1 < 0, x_2 < 0}]$. GAM (Fig. 3a, f_2)
293 merges all four modes into a single curve, limiting
294 its accuracy to $R^2 = 0.479$. GA²M (Fig. 3c,
295 f_2) captures partial structure through two 2D
296 interactions (x_1-x_3 and x_2-x_3) but misses the
297 full 3-way interaction, reaching $R^2 = 0.712$. Furthermore, the model's behavior is hard to grasp,
298 as it requires integrating information from three distinct plots: the main effect and two interaction
299 heatmaps. CALM (Fig. 3b, f_2) separates the four regimes into distinct 1D curves, achieving both
300 high accuracy ($R^2 = 0.949$) and clear interpretability.

Table 1: Synthetic examples: R^2 comparison (mean \pm std. dev. over 5-fold CV).

Dataset	GAM	GA ² M	CALM
Case 1	0.737 \pm 0.028	0.974 \pm 0.011	0.995 \pm 0.002
Case 2	0.479 \pm 0.052	0.712 \pm 0.036	0.949 \pm 0.024
Case 3	0.527 \pm 0.023	0.961 \pm 0.009	0.975 \pm 0.006

304 **Case 3.** We define f_3 as $y = x_1^2 + \log(|x_2|) \sin(\frac{\pi}{2}x_3)$, where the logarithmic effect of x_2 , $\log |x_2|$, is
305 modulated by the sinusoidal term $\sin(\frac{\pi}{2}x_3)$, resulting in a general interaction structure. GAM (Fig. 3a,
306 f_3) fails to capture the interaction and instead fits a blurred sinusoidal curve, leading to low accuracy
307 ($R^2 = 0.527$). GA²M (Fig. 3c, f_3) achieves high accuracy ($R^2 = 0.961$) by capturing the interaction
308 in the x_2-x_3 heatmap. Still, interpretation remains difficult, as understanding the $x_2 \rightarrow y$ effect
309 requires integrating three views: the main effect plot and two interaction heatmaps. CALM (Fig. 3b,
310 f_3) cannot fully express the continuous x_2-x_3 interaction, but approximates it by partitioning x_3
311 and assigning each segment an average logarithmic response. While simplified, this yields a high
312 accuracy ($R^2 = 0.975$) and a concise, transparent explanation that remains *close* to the underlying
313 behavior.

314 4.2 Evaluation on Real Datasets

315 On each real dataset, we evaluate CALM against an XGBoost baseline (with neural networks and
316 random forests included in the appendix), two GAM models (NAM and EBM, as well as SPLINE
317 included in the appendix), and three GA²M models: EB²M (EBM with pairwise interactions enabled),
318 NODE-GA²M, and GAMI-Net. Evaluation is based on accuracy for classification tasks and RMSE
319 for regression datasets. In the main tables, CALM uses XGBoost as the Step-1 reference model
320 and EBM as the regional GAM; full results across all teacher \times regional-GAM combinations are in
321 Appendix C.

322 **Classification Results.** On classification datasets (Table 2), CALM: (i) outperforms both GAM
323 baselines in almost all cases (9/10 over NAM; 6/10 over EBM with 4 draws) (ii) is competitive

324 with the three GA²M models: CALM matches or beats EB²M on 4/10, NODE-GA²M on 6/10, and
 325 GAMI-Net on 5/10 datasets, and (iii) matches or beats the XGBoost black-box on 5/10 datasets.

Table 3: Regression RMSE vs. baselines (mean \pm std. dev. over 5-fold CV). \checkmark/\times denotes strictly lower/higher RMSE than CALM; = denotes exactly equal RMSE to CALM.

Dataset	BB		GAM		CALM	GA ² M		
	XGB	NAM	EBM	EB ² M		NODE	GAMI	
Bike Sharing	39.35 \pm 1.38 \checkmark	101.97 \pm 1.31 \times	100.21 \pm 1.01 \times	55.67 \pm 1.22	54.80 \pm 0.96 \checkmark	54.47 \pm 1.58 \checkmark	53.44 \pm 1.90 \checkmark	
California Housing	0.45 \pm 0.01 \checkmark	0.61 \pm 0.01 \times	0.55 \pm 0.01 \times	0.51 \pm 0.01	0.49 \pm 0.01 \checkmark	0.50 \pm 0.01 \checkmark	0.51 \pm 0.04 =	
Parkinsons Motor	1.44 \pm 0.09 \checkmark	6.11 \pm 0.16 \times	4.20 \pm 0.09 \times	2.24 \pm 0.13	2.35 \pm 0.06 \times	3.49 \pm 0.31 \times	2.76 \pm 0.31 \times	
Parkinsons Total	1.86 \pm 0.08 \checkmark	7.90 \pm 0.10 \times	4.85 \pm 0.11 \times	2.97 \pm 0.09	2.77 \pm 0.06 \checkmark	4.60 \pm 0.53 \times	3.81 \pm 0.66 \times	
Seoul Bike	209.6 \pm 3.47 \checkmark	320.2 \pm 4.38 \times	303.7 \pm 3.86 \times	238.9 \pm 1.64	235.2 \pm 1.58 \checkmark	231.0 \pm 4.23 \checkmark	245.82 \pm 8.45 \times	
Wine	0.62 \pm 0.01 \checkmark	0.72 \pm 0.01 \times	0.70 \pm 0.01 \times	0.69 \pm 0.02	0.68 \pm 0.01 \checkmark	0.69 \pm 0.01 =	0.71 \pm 0.00 \times	
Energy	67.97 \pm 2.58 \checkmark	89.80 \pm 2.67 \times	85.23 \pm 2.49 \times	83.09 \pm 1.97	77.33 \pm 2.46 \checkmark	82.18 \pm 2.53 \checkmark	180.1 \pm 84.09 \times	
CCPP	3.09 \pm 0.09 \checkmark	4.21 \pm 0.05 \times	3.44 \pm 0.08 \times	3.42 \pm 0.07	3.28 \pm 0.07 \checkmark	3.97 \pm 0.07 \times	3.92 \pm 0.07 \times	
Electrical	0.04 \pm 0.02 \times	0.04 \pm 0.01 \times	0.02 \pm 0.01 =	0.02 \pm 0.01 =	0.02 \pm 0.01 =	0.01 \pm 0.01 \checkmark	0.02 \pm 0.01 =	
Elevators	2.0 \times 10 ⁻³ \pm 0 =	2.0 \times 10 ⁻³ \pm 0 =	2.0 \times 10 ⁻³ \pm 0 =	2.0 \times 10 ⁻³ \pm 0 =	2.0 \times 10 ⁻³ \pm 0 =	2.0 \times 10 ⁻³ \pm 0 =	2.0 \times 10 ⁻³ \pm 0 =	
No2	0.47 \pm 0.03 \checkmark	0.50 \pm 0.02 \times	0.49 \pm 0.03 =	0.49 \pm 0.04	0.47 \pm 0.03 \checkmark	0.52 \pm 0.02 \times	0.50 \pm 0.03 \times	
Sensory	0.51 \pm 0.03 \times	0.48 \pm 0.01 \times	0.48 \pm 0.01 \times	0.45 \pm 0.02	0.45 \pm 0.02 =	0.49 \pm 0.01 \times	0.46 \pm 0.03 \times	
Airfoil	1.54 \pm 0.10 \checkmark	4.80 \pm 0.20 \times	4.57 \pm 0.15 \times	2.50 \pm 0.15	2.17 \pm 0.09 \checkmark	2.13 \pm 0.11 \checkmark	4.79 \pm 0.21 \times	
Skill Craft	0.94 \pm 0.02 \times	0.93 \pm 0.03 \times	0.90 \pm 0.02 =	0.90 \pm 0.02	0.90 \pm 0.03 =	0.91 \pm 0.02 \times	1.38 \pm 0.80 \times	
Ailerons	2.0 \times 10 ⁻⁴ \pm 0 =	2.0 \times 10 ⁻⁴ \pm 0 =	2.0 \times 10 ⁻⁴ \pm 0 =	2.0 \times 10 ⁻⁴ \pm 0 =	2.0 \times 10 ⁻⁴ \pm 0 =	2.0 \times 10 ⁻⁴ \pm 0 =	2.0 \times 10 ⁻⁴ \pm 0 =	
W/D/L	10/2/3	0/2/13	0/5/10		9/5/1	6/3/6	1/4/10	

326 **Regression Results.** On regression datasets (Table 3), CALM: (i) outperforms both GAM baselines
 327 in the large majority of cases (13/15 over NAM; 10/15 over EBM); (ii) matches NODE-GA²M
 328 (6W/3D/6L) and outperforms GAMI-Net (10W/4D/1L), while EB²M achieves lower RMSE in 9/15
 329 cases; (iii) matches or beats XGBoost on 5/15 datasets. CALM thus achieves clear improvements
 330 over GAMs while maintaining 1D shape functions per region; the gap relative to EB²M and XGBoost
 331 reflects the expected cost of stronger interpretability constraints.

332 **Model Complexity.** In addition to its strong predictive performance, CALM achieves its results using
 333 remarkably few pairwise feature interactions, 6.2 on average for classification and 15.1 for regression.
 334 On classification tasks, this is fewer than EB²M (15.2), GAMI-Net (17.6), and NODE-GA²M (97.7).
 335 On regression tasks, CALM remains sparser than GAMI-Net (16.3) and NODE-GA²M (87.3), and
 336 close to EB²M (14.1), as shown in Tables 6 and 7 of Appendix C. While most GA²M methods allow
 337 explicit control over the number of interactions, we apply them using their default configurations. In
 338 contrast, CALM does not impose an interaction cap, but still discovers a sparse structure through
 339 region-based conditioning.

340 **Runtime.** On classification tasks, CALM is slightly slower than EBM and EB²M, but noticeably
 341 faster than the other two GA²M models (GAMI-Net and NODE-GA²M), and even faster than NAM,
 342 despite the fact that NAM does not model interactions. In regression tasks, CALM exhibits a runtime
 343 similar to NAM and EB²M, and remains significantly faster than the remaining GA²M baselines.
 344 Average runtimes across datasets are reported in Tables 8 and 9 (Appendix C).

345 **Human Interpretability.** A within-subjects user study ($N = 25$) showed that CALM explanations
 346 are significantly easier to use than heatmap-based GA²M representations: participants achieved 54%
 347 vs. 34% accuracy overall ($p < 0.05$), with the largest gain on prediction-change tasks (40% vs. 12%).
 348 Full details are in Appendix D.

349 5 Conclusion

350 We introduced CALM, a novel class of interpretable-by-design models that bridge the gap between
 351 the interpretability of GAMs and the accuracy of GA²Ms. CALM’s model feature interactions
 352 using *conditional feature effects*, a set of univariate shape functions per feature, conditioned on its
 353 interacting features. Extensive evaluation shows that conditional effects often match the accuracy of
 354 GA²Ms, without resorting to 3D plots or heatmaps. CALM’s main limitation is that interpretability
 355 may harden as the number of interactions increases. While each feature is associated with up to up to
 356 $2^{d_{\max}} = 4$ shape functions and, in practice, CALM activates an average of only 6.2 interactions in
 357 total (Table 6), extensive dependencies with other features can result in plots with numerous vertical
 358 dashed lines, making interpretation more difficult.

References

- 359
- 360 [1] Rishabh Agarwal, Levi Melnick, Nicholas Frosst, Xuezhou Zhang, Ben Lengerich, Rich
361 Caruana, and Geoffrey E Hinton. Neural additive models: Interpretable machine learning
362 with neural nets. In *Advances in Neural Information Processing Systems*, volume 34, pages
363 4699–4711, 2021.
- 364 [2] Elaine Angelino, Nicholas Larus-Stone, Daniel Alabi, Margo Seltzer, and Cynthia Rudin.
365 Learning certifiably optimal rule lists for categorical data, 2018.
- 366 [3] Jacob Bien and Robert Tibshirani. Prototype selection for interpretable classification. *The*
367 *Annals of Applied Statistics*, 5(4), 2011.
- 368 [4] L. Breiman, J. Friedman, R.A. Olshen, and C.J. Stone. *Classification and Regression Trees*.
369 Chapman and Hall/CRC, 1984.
- 370 [5] Luis M. Candanedo. Appliances energy prediction dataset, 2017.
- 371 [6] Luis M. Candanedo, Véronique Feldheim, and Dominique Deramaix. Data driven prediction
372 models of energy use of appliances in a low-energy house. *Energy and Buildings*, 140:81–97,
373 2017.
- 374 [7] Rich Caruana, Yin Lou, Johannes Gehrke, Paul Koch, Marc Sturm, and Noemie Elhadad.
375 Intelligible models for healthcare: Predicting pneumonia risk and hospital 30-day readmission.
376 In *Proceedings of the 21th ACM SIGKDD international conference on knowledge discovery*
377 *and data mining*, pages 1721–1730, 2015.
- 378 [8] Chun-Hao Chang, Rich Caruana, and Anna Goldenberg. Node-gam: Neural generalized additive
379 model for interpretable deep learning. In *International Conference on Learning Representations*,
380 2022.
- 381 [9] Chaofan Chen, Oscar Li, Chaofan Tao, Alina Jade Barnett, Jonathan Su, and Cynthia Rudin.
382 This looks like that: Deep learning for interpretable image recognition. In *NeurIPS*, 2019.
- 383 [10] Finale Doshi-Velez and Been Kim. Towards a rigorous science of interpretable machine learning.
384 *arXiv preprint arXiv:1702.08608*, 2017.
- 385 [11] Jerome H Friedman. Greedy function approximation: A gradient boosting machine. *Annals of*
386 *Statistics*, 29(5):1189–1232, 2001.
- 387 [12] Vasilis Gkolemis, Theodore Dalamagas, and Christos Diou. Dale: Differential accumulated
388 local effects for efficient and accurate global explanations. In *Asian Conference on Machine*
389 *Learning*, pages 375–390. PMLR, 2023.
- 390 [13] Vasilis Gkolemis, Theodore Dalamagas, Eirini Ntoutsis, and Christos Diou. Rhale: robust and
391 heterogeneity-aware accumulated local effects. In *ECAI 2023*, pages 859–866. IOS Press, 2023.
- 392 [14] Vasilis Gkolemis, Anargiros Tzerefos, Theodore Dalamagas, Eirini Ntoutsis, and Christos Diou.
393 Regionally additive models: Explainable-by-design models minimizing feature interactions.
394 In *Joint European Conference on Machine Learning and Knowledge Discovery in Databases*,
395 pages 433–447. Springer, 2023.
- 396 [15] Riccardo Guidotti, Anna Monreale, Salvatore Ruggieri, Franco Turini, Fosca Giannotti, and
397 Dino Pedreschi. A survey of methods for explaining black box models. *ACM Computing*
398 *Surveys*, 51(5):1–42, 2018.
- 399 [16] Trevor Hastie and Robert Tibshirani. *Generalized Additive Models*. Chapman and Hall/CRC,
400 1990.
- 401 [17] Julia Herbringer, Bernd Bischl, and Giuseppe Casalicchio. Repid: Regional effect plots with
402 implicit interaction detection. In Gustau Camps-Valls, Francisco J. R. Ruiz, and Isabel Valera,
403 editors, *Proceedings of The 25th International Conference on Artificial Intelligence and Statis-*
404 *tics*, volume 151 of *Proceedings of Machine Learning Research*, pages 10209–10233. PMLR,
405 March 2022.

- 406 [18] Julia Herbinger, Susanne Dandl, Fiona K Ewald, Sofia Loibl, and Giuseppe Casalicchio.
407 Leveraging model-based trees as interpretable surrogate models for model distillation. In
408 *European Conference on Artificial Intelligence*, pages 232–249. Springer, 2023.
- 409 [19] Julia Herbinger, Marvin N Wright, Thomas Nagler, Bernd Bischl, and Giuseppe Casalicchio.
410 Decomposing global feature effects based on feature interactions. *Journal of Machine Learning*
411 *Research*, 25(381):1–65, 2024.
- 412 [20] Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. Distilling the knowledge in a neural network. In
413 *NeurIPS Deep Learning and Representation Learning Workshop*, 2015.
- 414 [21] Noah Hollmann, Samuel Müller, Katharina Eggersperger, and Frank Hutter. Tabpfn: A
415 transformer that solves small tabular classification problems in a second. *arXiv preprint*
416 *arXiv:2207.01848*, 2022.
- 417 [22] Linwei Hu, Jie Chen, Vijayan N Nair, and Agus Sudjianto. Surrogate locally-interpretable
418 models with supervised machine learning algorithms. *arXiv preprint arXiv:2007.14528*, 2020.
- 419 [23] Shibal Ibrahim, Gabriel Afriat, Kayhan Behdin, and Rahul Mazumder. GRAND-SLAMIN’
420 Interpretable Additive Modeling with Structural Constraints. *Advances in Neural Information*
421 *Processing Systems*, 36:61158–61186, December 2023.
- 422 [24] Markelle Kelly, Rachel Longjohn, and Kolby Nottingham. The uci machine learning repository,
423 2025.
- 424 [25] Mathias Kraus, Daniel Tschernutter, Sven Weinzierl, and Patrick Zschech. Interpretable
425 generalized additive neural networks. *European Journal of Operational Research*, 317(2):303–
426 316, 2024.
- 427 [26] Benjamin Letham, Cynthia Rudin, Tyler H. McCormick, and David Madigan. Interpretable
428 classifiers using rules and bayesian analysis: Building a better stroke prediction model. *The*
429 *Annals of Applied Statistics*, 9(3), 2015.
- 430 [27] Zachary C Lipton. The mythos of model interpretability: In machine learning, the concept of
431 interpretability is both important and slippery. *Queue*, 16(3):31–57, 2018.
- 432 [28] Yin Lou, Rich Caruana, and Johannes Gehrke. Intelligible models for classification and
433 regression. In *Proceedings of the 18th ACM SIGKDD international conference on Knowledge*
434 *discovery and data mining*, pages 150–158, 2012.
- 435 [29] Yin Lou, Rich Caruana, Johannes Gehrke, and Giles Hooker. Accurate intelligible models with
436 pairwise interactions. In *Proceedings of the 19th ACM SIGKDD international conference on*
437 *Knowledge discovery and data mining*, pages 623–631, 2013.
- 438 [30] Christoph Molnar. *Interpretable Machine Learning*. 2020.
- 439 [31] W James Murdoch, Chandan Singh, Karl Kumbier, Reza Abbasi-Asl, and Bin Yu. Interpretable
440 machine learning: definitions, methods, and applications. *arXiv preprint arXiv:1901.04592*,
441 2019.
- 442 [32] Averkiy Oliabev. Home equity line of credit (HELOC), 2018.
- 443 [33] R. Kelley Pace and Ronald Barry. Sparse spatial autoregressions. *Statistics and Probability*
444 *Letters*, 33:291–297, 1997.
- 445 [34] ProPublica. Compas data and analysis for ‘machine bias’, 2016.
- 446 [35] Filip Radenovic, Abhimanyu Dubey, and Dhruv Mahajan. Neural basis models for interpretabil-
447 ity. *Advances in Neural Information Processing Systems*, 35:8414–8426, 2022.
- 448 [36] Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin. “why should i trust you?”: Explaining
449 the predictions of any classifier. In *Proceedings of the 22nd ACM SIGKDD International*
450 *Conference on Knowledge Discovery and Data Mining*, pages 1135–1144. ACM, 2016.

- 451 [37] Joseph D Romano, Trang T Le, William La Cava, John T Gregg, Daniel J Goldberg, Praneel
452 Chakraborty, Natasha L Ray, Daniel Himmelstein, Weixuan Fu, and Jason H Moore. Pmlb v1.0:
453 an open source dataset collection for benchmarking machine learning methods. *arXiv preprint*
454 *arXiv:2012.00058v2*, 2021.
- 455 [38] Cynthia Rudin. Stop explaining black box machine learning models for high stakes decisions
456 and use interpretable models instead. *Nature machine intelligence*, 1(5):206–215, 2019.
- 457 [39] Mohammed Saeed, Mauricio Villarroel, Andrew T. Reisner, Gari Clifford, Li-Wei Lehman,
458 George Moody, Thomas Heldt, Tin H. Kyaw, Benjamin Moody, and Roger G. Mark. Multipar-
459 ameter intelligent monitoring in intensive care ii (mimic-ii): A public-access intensive care
460 unit database. *Critical Care Medicine*, 39(5):952–960, 2011.
- 461 [40] Sarah Tan, Rich Caruana, Giles Hooker, and Yin Lou. Distill-and-compare: Auditing black-box
462 models using transparent model distillation. In *Proceedings of the 2018 AAAI/ACM Conference*
463 *on AI, Ethics, and Society*, pages 303–310, 2018.
- 464 [41] Joaquin Vanschoren, Jan N. van Rijn, Bernd Bischl, and Luis Torgo. Openml: Networked
465 science in machine learning. *SIGKDD Explorations*, 15(2):49–60, 2013.
- 466 [42] Simon N Wood. *Generalized additive models: an introduction with R*. CRC press, 2017.
- 467 [43] Zebin Yang, Aijun Zhang, and Agus Sudjianto. Gami-net: An explainable neural network based
468 on generalized additive models with structured interactions. *arXiv preprint arXiv:2003.07132*,
469 2020.

470 A Conditionally Additive Local Models (CALMs)

471 In the main paper (Algorithm 1), we summarize the procedure of fitting a Conditionally Additive
472 Local Model (CALM):

$$f_{\text{CALM}}(\mathbf{x}) = g^{-1} \left(\beta_0 + \sum_{i=1}^d f_i^{r_i(\mathbf{x}_{-i})}(x_i) \right),$$

473 as defined in Eq.(1), to a dataset $\mathcal{D} = \{\mathbf{x}^{(i)}, y^{(i)}\}_{i=1}^N$ in three main steps: (i) fit a high-capacity
474 reference model f_{ref} on \mathcal{D} , (ii) fit a partition tree T_i using an interaction-related heterogeneity
475 measure H_i for each feature $i = 1, \dots, d$, and (iii) estimate region-specific effects $\{f_i^{(r)}\}_{r=1}^{R_i}$
476 for each feature $i = 1, \dots, d$. We provide additional details for each step below.

477 A.1 Step 1: Fit a high-capacity reference model f_{ref} on \mathcal{D}

478 The initial stage involves fitting a high-capacity predictive model, denoted as $f_{\text{ref}} : \mathbb{R}^d \rightarrow \mathbb{R}$ to
479 dataset $\mathcal{D} = \{(\mathbf{x}^{(i)}, y^{(i)})\}_{i=1}^N$. The reference model serves as an accurate functional approximation
480 of the underlying relationship between the input features \mathbf{x} and the response variable y . The choice
481 of f_{ref} is flexible; any sufficiently expressive and accurate model can be employed. Commonly used
482 options include gradient-boosted decision trees (e.g., XGBoost, which we adopt by default), deep
483 neural networks, random forests, or more recent architectures such as TabPFN [21].

484 A.2 Step 2: Fit a partition tree T_i using heterogeneity H_i for each feature

485 We first define a suitable heterogeneity measure H_i (see Section A.2.1), and then explore the
486 relationship between heterogeneity and CALM’s error in approximating f_{ref} . Motivated by these
487 results, we fit a partition tree T_i for each feature i (see Section A.2.4). Finally, if the user specifies a
488 constraint K on the maximum number of interactions, we apply a pruning step to retain the K most
489 significant interactions (see Section A.2.5).

490 A.2.1 Heterogeneity Measures H_i

491 Interaction-related heterogeneity H_i quantifies the extent to which a feature x_i interacts with all other
492 features x_j , where $j \in \{1, \dots, d\} \setminus \{i\}$. The computation of H_i is based on the variance of local
493 effects: $h(x_i, \mathbf{x}_{-i}^{(j)})$.

494 **Local Effect.** We define $h(x_i, \mathbf{x}_{-i}^{(j)})$ as the local effect of feature x_i when applied to the j -th
495 background observation $\mathbf{x}_{-i}^{(j)}$. That is, if we take the j -th observation and substitute the i -th feature
496 with the value x_i , the resulting change (local effect) in the model output is captured by $h(x_i, \mathbf{x}_{-i}^{(j)})$.
497 The computation of $h(\cdot)$ relies on two components: (i) a reference teacher model f_{ref} and (ii) a
498 background dataset $\{\mathbf{x}^{(i)}\}_{i=1}^N$ containing input features only. The specific form of $h(\cdot)$ depends on
499 the chosen feature effect method. Several methods can be used, including Partial Dependence Plot
500 (PDP), SHAP-Dependence Plot (SHAP-DP) and Robust and Heterogeneity-aware ALE (RHALE).
501 Below, we provide the definition of h in the case of PDP, which we adopt as the default method in our
502 experiments. Let f_{ref} be the reference model and $\{\mathbf{x}^{(j)}\}_{j=1}^N$ the background dataset. The PDP-based
503 heterogeneity is defined as:

$$h(x_i, \mathbf{x}_{-i}^{(j)}) = f_{\text{ref}}(x_i, \mathbf{x}_{-i}^{(j)}) - c(\mathbf{x}_{-i}) \quad (6)$$

504 where $c(\mathbf{x}_{-i}^{(j)}) = \mathbb{E}_{X_i} [f_{\text{ref}}(X_i, \mathbf{x}_{-i}^{(j)})]$ is a centering constant. For details on alternative approaches,
505 we refer the reader to [19, 17].

506 **Point-wise heterogeneity $H_i(x_i)$ and Heterogeneity H_i .** The pointwise heterogeneity is then
507 defined as

$$H_i(x_i) = \mathbb{E}_{\mathbf{X}_{-i}} \left[(h(x_i, \mathbf{X}_{-i}) - \mu_i(x_i))^2 \right] \quad (7)$$

508 where $\mu_i(x_i) = \mathbb{E}_{\mathbf{X}_{-i}} [h(x_i, \mathbf{X}_{-i})]$ is the mean local effect. Inside a region, this heterogeneity takes
509 the form of Eq. (4). Feature-level heterogeneity is

$$H_i = \mathbb{E}_{X_i} [H_i(X_i)] = \mathbb{E}_{X_i} [\text{Var}_{\mathbf{X}_{-i}} [h(X_i, \mathbf{X}_{-i})]] \quad (8)$$

510 Inside a region this definition takes the form of Eq. (5).

511 A more formal motivation for the use of heterogeneity in CALM comes from the proposition of the
512 following section.

513 A.2.2 Heterogeneity and CALM approximation error

514 **Proposition A.1.** Let $\mathcal{E} = \mathbb{E}_{\mathbf{X}} \left[(f_{\text{ref}}(\mathbf{X}) - f_{\text{CALM}}(\mathbf{X}))^2 \right]$ be the mean squared error of the CALM
515 approximation of f_{ref} . Assume that $f_{\text{ref}} : \mathcal{X} \rightarrow \mathbb{R}$ admits a multivariate regionally additive
516 approximation

$$f_{\text{ref}}(\mathbf{x}) \approx \beta_0 + \sum_{i=1}^d h_i^{r_i(\mathbf{x}_{-i})}(\mathbf{x})$$

517 such that each $h_i^{r_i(\mathbf{x}_{-i})}$ captures the contribution of the i -th feature to the function's output (including
518 interactions). Also, assume that there is no error in fitting of the univariate CALM shape functions.
519 For clarity of exposition, the proof treats this approximation as an exact equality; the bound (9) holds
520 up to the residual approximation variance. Then,

$$\min_{f_{\text{CALM}}} \mathcal{E} \leq d \sum_{i=1}^d \left(\sum_{r=1}^{R_i} P_{\mathcal{R}_i^{(r)}} H_i^{\mathcal{R}_i^{(r)}} \right) \quad (9)$$

521 where region $\mathcal{R}_i^{(r)}$ is the r -th region of the i -th feature, $P_{\mathcal{R}_i^{(r)}}$ is the probability mass of the region in
522 the data and $H_i^{\mathcal{R}_i^{(r)}}$ is the heterogeneity of $h_i^{\mathcal{R}_i^{(r)}}(\mathbf{x})$.

523 *Proof.* For the approximation error we have

$$\begin{aligned} \mathcal{E} &= \mathbb{E}_{\mathbf{X}} \left[(f_{\text{ref}}(\mathbf{X}) - f_{\text{CALM}}(\mathbf{X}))^2 \right] = \mathbb{E}_{\mathbf{X}} \left[\left(\left(\beta_0 + \sum_{i=1}^d h_i^{r_i(\mathbf{x}_{-i})}(\mathbf{x}) \right) - \left(\beta_0 + \sum_{i=1}^d f_i^{r_i(\mathbf{x}_{-i})}(x_i) \right) \right)^2 \right] \\ &= \mathbb{E}_{\mathbf{X}} \left[\left(\sum_{i=1}^d \left(h_i^{r_i(\mathbf{x}_{-i})}(\mathbf{x}) - f_i^{r_i(\mathbf{x}_{-i})}(x_i) \right) \right)^2 \right] \leq d \sum_{i=1}^d \mathbb{E}_{\mathbf{X}} \left[\left(h_i^{r_i(\mathbf{x}_{-i})}(\mathbf{x}) - f_i^{r_i(\mathbf{x}_{-i})}(x_i) \right)^2 \right] \end{aligned}$$

524 where the last step results from the Jensen's inequality. We therefore have

$$\mathcal{E} \leq d \sum_{i=1}^d \mathcal{E}_i^{r_i(\mathbf{x}_{-i})} \quad (10)$$

525 Next, we show that the minimum error achieved for each region is equal to the heterogeneity of
 526 the local effects. For any measurable function $g_i(x_i)$ and region \mathcal{R} , decompose, using $\mu_i^{\mathcal{R}}(x_i) =$
 527 $\mathbb{E}_{\mathbf{X}_{-i} | \mathbf{X}_{-i} \in \mathcal{R}} [h_i^{\mathcal{R}}(x_i, \mathbf{X}_{-i})]$:

$$h_i(X_i, \mathbf{X}_{-i}) - g_i(X_i) = (h_i(X_i, \mathbf{X}_{-i}) - \mu_i^{\mathcal{R}}(X_i)) + (\mu_i^{\mathcal{R}}(X_i) - g_i(X_i))$$

528 Taking expectations inside \mathcal{R} and squaring:

$$\begin{aligned} & \mathbb{E}_{X_i, \mathbf{X}_{-i} | \mathbf{X}_{-i} \in \mathcal{R}} \left[(h_i^{\mathcal{R}}(X_i, \mathbf{X}_{-i}) - g_i(X_i))^2 \right] \\ &= \mathbb{E}_{X_i, \mathbf{X}_{-i} | \mathbf{X}_{-i} \in \mathcal{R}} \left[(h_i^{\mathcal{R}}(X_i, \mathbf{X}_{-i}) - \mu_i^{\mathcal{R}}(X_i))^2 \right] \\ & \quad + \mathbb{E}_{X_i} \left[(\mu_i^{\mathcal{R}}(X_i) - g_i(X_i))^2 \right] \\ & \quad + 2\mathbb{E}_{X_i} \left[\mathbb{E}_{\mathbf{X}_{-i} | \mathbf{X}_{-i} \in \mathcal{R}} [h_i^{\mathcal{R}}(X_i, \mathbf{X}_{-i}) - \mu_i^{\mathcal{R}}(X_i)] (\mu_i^{\mathcal{R}}(X_i) - g_i(X_i)) \right] \end{aligned}$$

529 The cross-term equals zero because:

$$\mathbb{E}_{\mathbf{X}_{-i} | \mathbf{X}_{-i} \in \mathcal{R}} [h_i^{\mathcal{R}}(x_i, \mathbf{X}_{-i}) - \mu_i^{\mathcal{R}}(x_i)] = \mu_i^{\mathcal{R}}(x_i) - \mu_i^{\mathcal{R}}(x_i) = 0$$

530 by definition of $\mu_i^{\mathcal{R}}$. Therefore:

$$\mathbb{E}_{X_i, \mathbf{X}_{-i} | \mathbf{X}_{-i} \in \mathcal{R}} \left[(h_i^{\mathcal{R}}(X_i, \mathbf{X}_{-i}) - g_i(X_i))^2 \right] = \mathcal{E}_i^{\mathcal{R}} + \mathbb{E}_{X_i} \left[(\mu_i^{\mathcal{R}}(X_i) - g_i(X_i))^2 \right] \geq \mathcal{E}_i^{\mathcal{R}}$$

531 with equality if and only if $g_i = \mu_i^{\mathcal{R}}$ almost surely. This indicates that within each region, $\mu_i^{\mathcal{R}}$
 532 minimizes the approximation error.

533 Moreover, from the definition of pointwise heterogeneity:

$$H_i^{\mathcal{R}}(x_i) := \mathbb{E}_{\mathbf{X}_{-i} | \mathbf{X}_{-i} \in \mathcal{R}} \left[(h_i^{\mathcal{R}}(x_i, \mathbf{X}_{-i}) - \mu_i^{\mathcal{R}}(x_i))^2 \right]$$

534 Taking the expectation over X_i :

$$H_i^{\mathcal{R}} = \mathbb{E}_{X_i} [H_i^{\mathcal{R}}(X_i)] = \mathbb{E}_{X_i} \left[\mathbb{E}_{\mathbf{X}_{-i} | \mathbf{X}_{-i} \in \mathcal{R}} \left[(h_i^{\mathcal{R}}(X_i, \mathbf{X}_{-i}) - \mu_i^{\mathcal{R}}(X_i))^2 \right] \right]$$

535 which gives

$$H_i^{\mathcal{R}} = \mathbb{E}_{X_i, \mathbf{X}_{-i} | \mathbf{X}_{-i} \in \mathcal{R}} \left[(h_i^{\mathcal{R}}(X_i, \mathbf{X}_{-i}) - \mu_i^{\mathcal{R}}(X_i))^2 \right] = \min \mathcal{E}_i^{\mathcal{R}}$$

536 i.e., heterogeneity is the minimum approximation error.

537 Since heterogeneity minimizes the error for each region we can estimate the feature-level error across
 538 regions as

$$\min \mathcal{E}_i = \sum_{r=1}^{R_i} P_{\mathcal{R}_i^{(r)}} H_i^{\mathcal{R}_i^{(r)}}$$

539 where $P_{\mathcal{R}_i^{(r)}}$ is the probability mass of region $\mathcal{R}_i^{(r)}$.

540 Using this result with Eq. (10), we have

$$\min_{f_{\text{CALM}}} \mathcal{E} \leq d \sum_{i=1}^d \left(\sum_{r=1}^{R_i} P_{\mathcal{R}_i^{(r)}} H_i^{\mathcal{R}_i^{(r)}} \right)$$

541

□

542 *Remark A.2.* The switch from the conditional expectation $\mathbb{E}_{X_i|\mathbf{X}_{-i} \in \mathcal{R}}$ to the marginal \mathbb{E}_{X_i} in the
 543 second term of the decomposition above implicitly assumes $X_i \perp \mathbf{X}_{-i} \mid \mathbf{X}_{-i} \in \mathcal{R}$ within each
 544 region. This is the standard assumption underlying PDP-style decompositions and is shared by all
 545 feature-effect methods that use the marginal distribution of X_i .

546 This result links heterogeneity with the estimation error of the local effect inside a region and motivates
 547 the CALM approach for region splitting. The initial assumption about the additive approximation
 548 of f_{ref} is a common approach followed by feature effect methods (where the effect of each feature
 549 is computed independently). Moreover, the selected local effect function and heterogeneity are
 550 solely used for identifying the splitting regions and not as estimators of f_{ref} . Using this result, and
 551 especially Eq. (9), the following two sections present the CALM approach for using heterogeneity to
 552 define a partition for each feature.

553 A.2.3 Heterogeneity estimation

554 For heterogeneity estimation, denote with $\mathcal{I} \subseteq \{1, \dots, N\}$ the index set of active background
 555 instances considered in the heterogeneity calculation i.e., those that reside in region $\mathcal{R}_i^{(r_i(\mathbf{x}_{-i}))}$ of the
 556 i -th instance. Then, the point-wise heterogeneity at feature value x_i over \mathcal{I} can be estimated by

$$\hat{H}_i^{\mathcal{I}}(x_i) = \frac{1}{|\mathcal{I}|} \sum_{j \in \mathcal{I}} \left(h(x_i, \mathbf{x}_{-i}^{(j)}) - \mu(x_i) \right)^2, \text{ where } \mu(x_i) = \frac{1}{|\mathcal{I}|} \sum_{j \in \mathcal{I}} h(x_i, \mathbf{x}_{-i}^{(j)}) \quad (11)$$

557 Eq.(11) quantifies the strength of the interactions of the i -th feature, at position x_i , as the variance of
 558 the local effects at x_i . To obtain a global measure of these interactions, we average the pointwise
 559 heterogeneity over a grid of M values $\{\tilde{x}_i^{(m)}\}_{m=1}^M$ sampled from the domain of x_i :

$$\hat{H}_i^{\mathcal{I}} = \frac{1}{M} \sum_{m=1}^M H_i^{\mathcal{I}}(\tilde{x}_i^{(m)}). \quad (12)$$

560 A.2.4 Computing the partition tree T_i for each feature.

561 For each feature x_i , we construct a binary tree T_i of maximum depth d_{\max} . At each internal node
 562 of the tree, we evaluate candidate binary splits based on all features x_j for $j \in \{1, \dots, d\} \setminus \{i\}$.
 563 For each candidate splitting feature x_j , we consider a fixed number T of candidate threshold values
 564 τ . These candidate thresholds are selected as T equally spaced values over the range of x_j in the
 565 training data, i.e., $[\min_k x_j^{(k)}, \max_k x_j^{(k)}]$, where $x_j^{(k)}$ denotes the value of feature x_j for the k -th
 566 training instance. Alternatively, users may choose to evaluate all unique observed values of x_j , i.e.,
 567 $\{x_j^{(k)}\}_{k=1}^N$.

568 To determine the best split at each node, we compute the heterogeneity drop, which quantifies
 569 the decrease in interaction-related heterogeneity of the target feature x_i after performing the split.
 570 Specifically, for a candidate split, we partition the current set of instances \mathcal{I} into left and right subsets
 571 \mathcal{I}_L and \mathcal{I}_R , based on whether the splitting feature x_j falls below or above the threshold τ . The
 572 heterogeneity drop is defined as:

$$\Delta H_i^- = \frac{H_i^{\mathcal{I}} - \left(\frac{|\mathcal{I}_L|}{|\mathcal{I}|} H_i^{\mathcal{I}_L} + \frac{|\mathcal{I}_R|}{|\mathcal{I}|} H_i^{\mathcal{I}_R} \right)}{H_i^{\mathcal{I}}} \quad (13)$$

573 where $H_i^{\mathcal{I}}$ denotes the interaction-related heterogeneity of feature x_i over the current set of instances
 574 \mathcal{I} ; \mathcal{I}_L and \mathcal{I}_R are the subsets of \mathcal{I} resulting from the candidate split; $H_i^{\mathcal{I}_L}$ and $H_i^{\mathcal{I}_R}$ are the hetero-
 575 geneities of x_i computed on the left and right subsets, respectively; and $|\mathcal{I}_L|/|\mathcal{I}|$ and $|\mathcal{I}_R|/|\mathcal{I}|$ are the
 576 proportions of instances in each subset, which serve as weights in the weighted average.

577 This normalized metric enhances interpretability by allowing users to define a meaningful threshold,
 578 ϵ , for split acceptance. For instance, a user may require a candidate split to achieve a minimum
 579 heterogeneity reduction of $\epsilon = 0.2$ (representing a 20% drop) to be considered significant. In our
 580 experimental setup, we utilize a default threshold of $\epsilon = 0.2$.

581 The procedure for fitting the tree T_i is given in Algorithm 2.

Algorithm 2 Fitting a Partition Tree T_i for Feature x_i

```
1: Input:  $f_{\text{ref}}, \mathcal{D} = \{(\mathbf{x}^{(k)}, y^{(k)})\}_{k=1}^N$ ; Parameters: max depth  $d_{\text{max}}$ , threshold  $\epsilon$ , grid size  $M = 20$ 
2: Output: Binary tree  $T_i$ 
3:
4: function BUILDTREE(node  $\nu$ , depth  $\ell$ , active indices  $\mathcal{I}$ )
5: if  $\ell = d_{\text{max}}$  then
6:   return  $\nu$  {Stop splitting}
7: end if
8: Compute  $H_i^{\mathcal{I}}$  using Eq.(12)
9: Initialize:  $\Delta H_{\text{max}} \leftarrow 0$ 
10: for each feature  $j \neq i$  do
11:   Determine candidate thresholds  $\mathcal{T}_j$  {Default: Unique values if  $x_j$  categorical, else  $M$ -grid (default:  $M = 20$ )}
12:   for each  $\tau \in \mathcal{T}_j$  do
13:     if  $x_j$  is numerical then
14:        $\mathcal{I}_L \leftarrow \{k \in \mathcal{I} : x_j^{(k)} \leq \tau\}, \mathcal{I}_R \leftarrow \mathcal{I} \setminus \mathcal{I}_L$ 
15:     else { $x_j$  is categorical}
16:        $\mathcal{I}_L \leftarrow \{k \in \mathcal{I} : x_j^{(k)} = \tau\}, \mathcal{I}_R \leftarrow \mathcal{I} \setminus \mathcal{I}_L$ 
17:     end if
18:     Compute  $\Delta H_i$  using Eq. (13)
19:     if  $\Delta H_i > \Delta H_{\text{max}}$  then
20:        $\Delta H_{\text{max}} \leftarrow \Delta H_i$ , Store  $(j, \tau)$  as optimal split
21:     end if
22:   end for
23: end for
24: if  $\Delta H_{\text{max}} > \epsilon$  then
25:   Split node  $\nu$  into  $\nu_L, \nu_R$  using optimal  $(j, \tau)$ 
26:    $\nu_L \leftarrow \text{BUILDTREE}(\nu_L, \ell + 1, \mathcal{I}_L)$ 
27:    $\nu_R \leftarrow \text{BUILDTREE}(\nu_R, \ell + 1, \mathcal{I}_R)$ 
28: end if
29: return  $\nu$ 
30: end function
```

582 A.2.5 Prune trees to keep top K interactions

583 If the user sets an upper threshold K on the number of interactions, we want to select the K most
584 important splits. Given the trees T_i for $i = 1, \dots, d$, we denote by ΔH_i^ν the heterogeneity reduction
585 associated with node ν in tree T_i . We collect all such splits ΔH_i^ν across all features i and nodes ν , and
586 then sort them in descending order. From this ordered list, we select the top K nodes that correspond
587 to the largest heterogeneity decreases, while ensuring that no child node is retained without its parent
588 node also being included. Based on this selection, we prune each tree T_i by removing nodes outside
589 the retained set, thereby preserving only the most significant interactions.

590 A.3 Step 3: Estimating the region-specific effects $\{f_i(r)\}_{r=1}^{R_i}$ for each feature.

591 The original gradient boosting procedure for fitting a standard, global GAM, $f_{\text{GAM}}(\mathbf{x}) =$
592 $g^{-1}(\beta_0 + \sum_i f_i(x_i))$, follows the round-robin approach of [11]. At each boosting iteration, one
593 univariate shape function f_i is updated by fitting to the current residuals over the entire dataset,
594 thereby greedily reducing the overall loss. The procedure as describe by [28] is summarized in
595 Algorithm 3.

596 To fit a CALM: $f_{\text{CALM}}(\mathbf{x}) = g^{-1}(\beta_0 + \sum_{i=1}^d f_i^{r_i(\mathbf{x}-i)}(x_i))$ we customize the standard gradient
597 boosting to accommodate for region-specific effects, i.e., we adapt this scheme so that each shape
598 function $f_i^{(r)}$ is trained only on the subset of instances that lie within its assigned subregion $\mathcal{R}_i^{(r)}$.

Algorithm 3 Gradient Boosting for GAM

```
1: Initialize  $f_i \leftarrow 0$  for all  $i = 1, \dots, d$ 
2: for  $m = 1$  to  $M$  do
3:   for  $i = 1$  to  $d$  do
4:      $\mathcal{E} \leftarrow \{(x_i^{(j)}, y^{(j)} - f_{\text{GAM}}(\mathbf{x}^{(j)}))\}_{j=1}^N$  {Compute partial residuals}
5:     Learn shaping function  $S : x_i \rightarrow \mathbb{R}$  using  $\mathcal{E}$  {Fit shape function}
6:      $f_i \leftarrow f_i + \eta S$  {Update with learning rate  $\eta$ }
7:   end for
8: end for
```

Algorithm 4 Gradient Boosting for CALM

```
1: Initialize  $\beta_0 \leftarrow \frac{1}{N} \sum_{j=1}^N y^{(j)}$ ;  $f_i^{(r)} \leftarrow 0$  for all  $i = 1, \dots, d$  and  $r = 1, \dots, R_i$ 
2: for  $m = 1$  to  $M$  do
3:   for  $i = 1$  to  $d$  do
4:     for  $r = 1$  to  $R_i$  do
5:        $\mathcal{E} \leftarrow \{(x_i^{(j)}, y^{(j)} - f_{\text{CALM}}(\mathbf{x}^{(j)})) : \mathbf{x}_{-i}^{(j)} \in \mathcal{R}_i^{(r)}\}$  {Filter by region  $\mathcal{R}_i^{(r)}$ }
6:       Learn shaping function  $S : x_i \rightarrow \mathbb{R}$  using  $\mathcal{E}$ 
7:        $f_i^{(r)} \leftarrow f_i^{(r)} + \eta S$  {Update regional shape function}
8:     end for
9:   end for
10:  Center: for all  $i, r$ , let  $c_i^{(r)} \leftarrow \frac{1}{|\mathcal{R}_i^{(r)}|} \sum_{\mathbf{x}^{(j)} \in \mathcal{R}_i^{(r)}} f_i^{(r)}(x_i^{(j)})$ ; set  $f_i^{(r)} \leftarrow f_i^{(r)} - c_i^{(r)}$  and
     $\beta_0 \leftarrow \beta_0 + \sum_i \sum_r \mathbb{P}(\mathcal{R}_i^{(r)}) c_i^{(r)}$  {Enforce identifiability}
11: end for
```

599 This preserves the additive, boosting framework while enforcing that each $f_i^{(r)}$ captures only the
600 behavior of x_i within its corresponding region. The procedure is summarized in Algorithm 4.

601 A.3.1 Theoretical analysis of Step 3

602 This section formalizes the claims made in Step 3 of the main text. We show that, under squared loss
603 and fixed region-selection trees, (i) the population target of Step 3 is an $L_2(P_{\mathbf{X}})$ -best approximation
604 of the true regression function onto the CALM function class, and (ii) an idealized exact version of
605 the algorithm converges to the corresponding empirical risk minimizer.

606 Consider the CALM score function

$$s(\mathbf{x}) = \beta_0 + \sum_{i=1}^d f_i^{(r_i(\mathbf{X}_{-i}))}(x_i), \quad (14)$$

607 where the region-selection functions $\{r_i\}$ (or equivalently the trees $\{T_i\}$) are treated as fixed. Let

$$m(\mathbf{x}) = \mathbb{E}[Y \mid \mathbf{X} = \mathbf{x}]$$

608 denote the true regression function. We assume regression with identity link and squared loss, and
609 $\mathbb{E}[Y^2] < \infty$. We restrict attention to scores $s \in L_2(P_{\mathbf{X}})$ (hence $m \in L_2(P_{\mathbf{X}})$).

610 To ensure identifiability, we impose the centering condition

$$\mathbb{E}\left[f_i^{(r)}(X_i) \mid r_i(\mathbf{X}_{-i}) = r\right] = 0 \quad \text{for all } i, r \text{ with } \mathbb{P}(r_i(\mathbf{X}_{-i}) = r) > 0. \quad (15)$$

611 This convention assigns all region-wise constants to the intercept β_0 and ensures uniqueness of the
612 decomposition.

613 Let $\mathcal{H}(\{T_i\})$ denote the class of all CALM score functions of the form (14) satisfying (15), and
614 assume $\mathcal{H}(\{T_i\}) \subset L_2(P_{\mathbf{X}})$ is a nonempty closed *convex* set.

615 A.3.2 Population target

616 **Proposition A.3** (Population optimality). *There exists at least one minimizer*

$$s^* \in \arg \min_{s \in \mathcal{H}(\{T_i\})} \mathbb{E}[(Y - s(\mathbf{X}))^2],$$

617 and it satisfies

$$s^* \in \arg \min_{s \in \mathcal{H}(\{T_i\})} \mathbb{E}[(m(\mathbf{X}) - s(\mathbf{X}))^2].$$

618 In particular,

$$\mathbb{E}[(Y - s^*(\mathbf{X}))^2] - \mathbb{E}[(Y - m(\mathbf{X}))^2] = \mathbb{E}[(m(\mathbf{X}) - s^*(\mathbf{X}))^2].$$

619 *Proof.* For any measurable function s ,

$$Y - s(\mathbf{X}) = (Y - m(\mathbf{X})) + (m(\mathbf{X}) - s(\mathbf{X})).$$

620 Taking squares and expectations yields

$$\mathbb{E}[(Y - s(\mathbf{X}))^2] = \mathbb{E}[(Y - m(\mathbf{X}))^2] + \mathbb{E}[(m(\mathbf{X}) - s(\mathbf{X}))^2],$$

621 since $\mathbb{E}[Y - m(\mathbf{X}) \mid \mathbf{X}] = 0$ eliminates the cross-term. The first term does not depend on s ,
 622 so minimizing prediction error over $\mathcal{H}(\{T_i\})$ is equivalent to minimizing $\mathbb{E}[(m(\mathbf{X}) - s(\mathbf{X}))^2]$.
 623 Existence of s^* follows since $\mathcal{H}(\{T_i\}) \subset L_2(P_{\mathbf{X}})$ is nonempty, closed, and convex. Condition (15)
 624 ensures uniqueness of the representation $(\beta_0, \{f_i^{(r)}\})$ for any given $s \in \mathcal{H}(\{T_i\})$, assuming no exact
 625 functional dependence (concurvity) among features within each region. \square

626 **Interpretation.** Proposition A.3 shows that, with fixed regions, Step 3 computes an $L_2(P_{\mathbf{X}})$ -best
 627 approximation of the true regression function m within the CALM function class. Any remaining error
 628 is therefore purely due to the structural limitations of the chosen regions and univariate components,
 629 not to the optimization procedure.

630 A.3.3 Empirical convergence of exact Step 3

631 Given i.i.d. data $\{(x^{(k)}, y^{(k)})\}_{k=1}^N$, define the empirical risk

$$\widehat{\mathcal{R}}_N(s) = \frac{1}{N} \sum_{k=1}^N (y^{(k)} - s(x^{(k)}))^2.$$

632 **Proposition A.4** (Convergence of exact cyclic regional backfitting). *Assume fixed trees and squared*
 633 *loss. Consider an idealized version of Step 3 that cyclically updates each $f_i^{(r)}$ by exactly minimizing*
 634 *$\widehat{\mathcal{R}}_N$ over that function using only samples with $r_i(\mathbf{x}_{-i}^{(k)}) = r$, followed by empirical centering*
 635 *implemented by shifting the region-wise sample mean from $f_i^{(r)}$ into the intercept β_0 , so that fitted*
 636 *values (and hence $\widehat{\mathcal{R}}_N$) are unchanged. Then $\widehat{\mathcal{R}}_N$ is non-increasing along the iterates, and the*
 637 *procedure converges (in empirical L_2) to a limit $\hat{s} \in \arg \min_{s \in \mathcal{H}(\{T_i\})} \widehat{\mathcal{R}}_N(s)$. Moreover, the*
 638 *fitted-value vector $(\hat{s}(\mathbf{x}^{(1)}), \dots, \hat{s}(\mathbf{x}^{(N)}))$ is unique.*

639 *Proof.* With fixed trees, the CALM score is linear in the collection of shape functions $\{f_i^{(r)}\}$
 640 evaluated on the data, and $\widehat{\mathcal{R}}_N$ is a convex quadratic function of these parameters. On the finite
 641 sample, $\widehat{\mathcal{R}}_N$ depends on each $f_i^{(r)}$ only through the finite vector $\{f_i^{(r)}(x_i^{(k)}) : r_i(\mathbf{x}_{-i}^{(k)}) = r\}$, so
 642 the problem can be viewed as a finite-dimensional least-squares problem. Each regional update
 643 solves the exact least-squares problem for one block of parameters while holding the others fixed,
 644 and therefore cannot increase $\widehat{\mathcal{R}}_N$. Although regions overlap across features, all updates jointly
 645 optimize a single global objective. Standard results for cyclic block coordinate descent on convex
 646 quadratics imply convergence to a global minimizer, and the minimizing fitted values are unique by
 647 strict convexity of the quadratic loss in the fitted-value vector. The centering condition (15) ensures a
 648 unique decomposition into components. \square

649 **Interpretation.** This result shows that Step 3 is not a collection of independent local fits. Instead,
 650 it performs a coupled optimization of a single additive predictor with region-gated components,
 651 analogous to classical backfitting for GAMs. Under exact updates, this procedure is guaranteed to
 652 converge to the best CALM fit for the given data and fixed regions.

653 **B Formal Proofs for Interpretability Properties**

654 This appendix provides formal characterizations and proofs for the interpretability properties dis-
 655 cussed in Section 3.3. For each property, we define it precisely and analyze whether it can be
 656 answered exactly under GAM, GA²M, and CALM.

657 **Note on link functions.** Properties P1–P3 are stated and proved for the additive linear predictor
 658 $\eta(\mathbf{x}) = \beta_0 + \sum_i f_i^{(r_i(\mathbf{x}-i))}(x_i)$. When a non-identity link function g is used (Eq. 1), these properties
 659 apply to η and translate to the final output $\hat{y} = g^{-1}(\eta)$ for any monotone link (e.g., logit, log).

660 **Proposition B.1 (Local Feature Contribution).** *Let $\hat{y} : \mathbb{R}^d \rightarrow \mathbb{R}$ be a prediction function, and fix*
 661 *an input $\mathbf{x} = (x_1, \dots, x_d) \in \mathbb{R}^d$. For each index $i \in \{1, \dots, d\}$, we define the local contribution*
 662 *$\phi_i(\mathbf{x}) \in \mathbb{R}$ as a function intended to represent the contribution of feature x_i to the output $\hat{y}(\mathbf{x})$. We*
 663 *analyze whether such a decomposition*

$$\hat{y}(\mathbf{x}) = \sum_{i=1}^d \phi_i(\mathbf{x})$$

664 *is uniquely determined by the structure of the model.*

665 **(A) In GAM.** *Assume the model has the additive form*

$$\hat{y}(\mathbf{x}) = \sum_{j=1}^d f_j(x_j),$$

666 *with each $f_j : \mathbb{R} \rightarrow \mathbb{R}$. Then for all $i \in \{1, \dots, d\}$,*

$$\phi_i(\mathbf{x}) := f_i(x_i)$$

667 *is well-defined, and*

$$\hat{y}(\mathbf{x}) = \sum_{i=1}^d \phi_i(\mathbf{x}).$$

668 **(B) In GA²M.** *Assume the model includes pairwise interactions:*

$$\hat{y}(\mathbf{x}) = \sum_{j=1}^d f_j(x_j) + \sum_{j < k} f_{jk}(x_j, x_k),$$

669 *where each $f_j : \mathbb{R} \rightarrow \mathbb{R}$ and $f_{jk} : \mathbb{R}^2 \rightarrow \mathbb{R}$. Then in general, there does not exist a unique additive*
 670 *decomposition*

$$\hat{y}(\mathbf{x}) = \sum_{i=1}^d \phi_i(\mathbf{x})$$

671 *with each $\phi_i(\mathbf{x})$ depending only on x_i .*

672 **(C) In CALM.** *Assume the model has the form*

$$\hat{y}(\mathbf{x}) = \sum_{j=1}^d f_j^{(r_j(\mathbf{x}-j))}(x_j),$$

673 *where for each j , $r_j : \mathbb{R}^{d-1} \rightarrow \{1, \dots, R_j\}$ is a region selector and $f_j^{(r)} : \mathbb{R} \rightarrow \mathbb{R}$ is a univariate*
 674 *shape function for region r . Here $\mathbf{x}_{-j} := \mathbf{x} \setminus x_j \in \mathbb{R}^{d-1}$.*

675 *Then for all $i \in \{1, \dots, d\}$,*

$$\phi_i(\mathbf{x}) := f_i^{(r_i(\mathbf{x}-i))}(x_i)$$

676 *is well-defined, and*

$$\hat{y}(\mathbf{x}) = \sum_{i=1}^d \phi_i(\mathbf{x}).$$

677 *Proof.* We examine each case separately.

678 **(A) In GAM.** By direct substitution, the prediction is an additive sum of univariate functions, each
 679 depending only on a single feature x_j . Thus the contribution of feature x_i is uniquely defined as
 680 $f_i(x_i)$, and the decomposition is exact.

681 **(B) In GA²M.** Suppose, for contradiction, that there exists a decomposition

$$\hat{y}(\mathbf{x}) = \sum_{i=1}^d \phi_i(\mathbf{x}),$$

682 where each $\phi_i : \mathbb{R}^d \rightarrow \mathbb{R}$ represents the contribution of feature x_i and depends only on x_i , i.e.,
 683 $\phi_i(\mathbf{x}) = \phi_i(x_i)$.

684 Then each interaction term $f_{jk}(x_j, x_k)$ must be split between ϕ_j and ϕ_k in such a way that the total
 685 sum remains correct and additive over individual features. This is only possible if f_{jk} is additively
 686 separable, i.e., if there exist functions $g_j, g_k : \mathbb{R} \rightarrow \mathbb{R}$ such that

$$f_{jk}(x_j, x_k) = g_j(x_j) + g_k(x_k).$$

687 However, the model does not constrain f_{jk} to be separable. In the general case, f_{jk} is non-separable
 688 and depends jointly on x_j and x_k .

689 Therefore, no decomposition $\sum_i \phi_i(x_i)$ can reproduce the prediction $\hat{y}(\mathbf{x})$ using only univariate
 690 terms. The contributions $\phi_i(\mathbf{x})$ are not uniquely determined by the model, and must rely on external
 691 assumptions or post-hoc attribution methods.

692 **(C) In CALM.** For a fixed input \mathbf{x} , each region function $r_j(\mathbf{x}_{-j})$ returns a unique region index in
 693 $\{1, \dots, R_j\}$. The corresponding function $f_j^{(r_j(\mathbf{x}_{-j}))}$ is applied to x_j , yielding a scalar. Thus, the
 694 contribution of feature x_i is precisely defined as:

$$\phi_i(\mathbf{x}) := f_i^{(r_i(\mathbf{x}_{-i}))}(x_i),$$

695 and the model prediction is recovered by summing over all features:

$$\hat{y}(\mathbf{x}) = \sum_{i=1}^d f_i^{(r_i(\mathbf{x}_{-i}))}(x_i) = \sum_{i=1}^d \phi_i(\mathbf{x}).$$

696

□

697 The second property asks: *How does the prediction change if we perturb a single feature x_i while*
 698 *keeping all others fixed?* The answer differs across GAM, GA²M, and CALM, as shown below.

699 **Proposition B.2** (Regional Feature Sensitivity). *Let $\hat{y} : \mathbb{R}^d \rightarrow \mathbb{R}$ be a prediction function. Fix an*
 700 *input $\mathbf{x} \in \mathbb{R}^d$, an index $i \in \{1, \dots, d\}$, and a perturbation $\varepsilon > 0$. Define the prediction change*
 701 *under perturbation of x_i as*

$$\Delta \hat{y} := \hat{y}(\mathbf{x} + \varepsilon \mathbf{e}_i) - \hat{y}(\mathbf{x}).$$

702 *We analyze how $\Delta \hat{y}$ can be computed in GAM, GA²M, and CALM.*

703 **(A) In GAM.** *Assume the model has the additive form*

$$\hat{y}(\mathbf{x}) = \sum_{j=1}^d f_j(x_j),$$

704 *with each $f_j : \mathbb{R} \rightarrow \mathbb{R}$ univariate. Then*

$$\Delta \hat{y} = f_i(x_i + \varepsilon) - f_i(x_i).$$

705 **(B) In GA²M.** *We show that in GA²M, the change $\Delta \hat{y}$ caused by perturbing x_i depends on the values*
 706 *of all interacting features x_j for $j \neq i$, and thus cannot be computed from x_i alone.*

707 **(C) In CALM:** *We show that the change in prediction $\Delta \hat{y} := \hat{y}(\mathbf{x} + \varepsilon \mathbf{e}_i) - \hat{y}(\mathbf{x})$ is computable from*
 708 *x_i alone if and only if no region transitions are triggered in other features. Otherwise, the change*
 709 *depends on \mathbf{x}_{-i} , and regional sensitivity is not determined solely by x_i .*

710 *Proof. Note:* The proposition is stated for a fixed input \mathbf{x} (single-valued); the set-valued definition of
 711 P2 in the main text arises by ranging over all possible values of \mathbf{x}_{-i} , each of which selects a region
 712 $r_i(\mathbf{x}_{-i}) \in \{1, \dots, R_i\}$ and thus one element of the set $\{\Delta f_i^{(r)}\}_{r=1}^{R_i}$.

713 **(A) In GAM.** Since all other features remain unchanged, and their contributions are independent of
 714 x_i , we have

$$\hat{y}(\mathbf{x} + \varepsilon \mathbf{e}_i) = \sum_{j \neq i} f_j(x_j) + f_i(x_i + \varepsilon),$$

$$\hat{y}(\mathbf{x}) = \sum_{j \neq i} f_j(x_j) + f_i(x_i),$$

716 and therefore

$$\Delta \hat{y} = f_i(x_i + \varepsilon) - f_i(x_i).$$

717 **(B) In GA²M.** Let $\mathcal{J} := \{j \neq i \mid f_{ij}(x_i, x_j) \text{ is present in the model}\}$, i.e., the set of features that
 718 interact with x_i . Since only x_i is perturbed, all other features remain unchanged. The only affected
 719 terms are the univariate function $f_i(x_i)$, and the interaction terms $f_{ij}(x_i, x_j)$ for each $j \in \mathcal{J}$. All
 720 other terms in the model remain constant.

721 The prediction after perturbation is:

$$\hat{y}(\mathbf{x} + \varepsilon \mathbf{e}_i) = f_i(x_i + \varepsilon) + \sum_{j \in \mathcal{J}} f_j(x_j) + \sum_{j \in \mathcal{J}} f_{ij}(x_i + \varepsilon, x_j) + C,$$

722 and before perturbation:

$$\hat{y}(\mathbf{x}) = f_i(x_i) + \sum_{j \in \mathcal{J}} f_j(x_j) + \sum_{j \in \mathcal{J}} f_{ij}(x_i, x_j) + C,$$

723 where $C := \sum_{j < k, j, k \neq i} f_{jk}(x_j, x_k)$ denotes the interaction terms not involving x_i , which are
 724 unaffected.

725 Subtracting, we obtain:

$$\Delta \hat{y} = f_i(x_i + \varepsilon) - f_i(x_i) + \sum_{j \in \mathcal{J}} [f_{ij}(x_i + \varepsilon, x_j) - f_{ij}(x_i, x_j)].$$

726 Since this expression depends on interacting x_j for $j \neq i$, it cannot be computed from x_i alone.

727 **(C) In CALM.** Assume the model has the form

$$\hat{y}(\mathbf{x}) = \sum_{j=1}^d f_j^{(r_j(\mathbf{x}_{-j}))}(x_j),$$

728 where $r_j : \mathbb{R}^{d-1} \rightarrow \{1, \dots, R_j\}$ assigns a region index to each feature j based on the context
 729 $\mathbf{x}_{-j} := \mathbf{x} \setminus x_j$, and $f_j^{(r)} : \mathbb{R} \rightarrow \mathbb{R}$ is the shape function for region r .

730 Let $\mathbf{x}_{-i} = \mathbf{x} \setminus x_i$, and define the perturbation $x_i \mapsto x_i + \varepsilon$. Then

$$\Delta \hat{y} = f_i^{(r_i(\mathbf{x}_{-i}))}(x_i + \varepsilon) - f_i^{(r_i(\mathbf{x}_{-i}))}(x_i) + \sum_{j \neq i} \left[f_j^{(r_j(\mathbf{x}_{-j}^{+\varepsilon}))}(x_j) - f_j^{(r_j(\mathbf{x}_{-j}))}(x_j) \right],$$

731 where $\mathbf{x}_{-j}^{+\varepsilon} := \mathbf{x}_{-j}$ with $x_i \mapsto x_i + \varepsilon$ (since $x_i \in \mathbf{x}_{-j}$ for all $j \neq i$).

732 *Case 1: No region transitions in other features.* Assume that for all $j \neq i$, we have $r_j(\mathbf{x}_{-j}^{+\varepsilon}) =$
 733 $r_j(\mathbf{x}_{-j})$. Then:

$$\Delta \hat{y} = f_i^{(r_i(\mathbf{x}_{-i}))}(x_i + \varepsilon) - f_i^{(r_i(\mathbf{x}_{-i}))}(x_i).$$

734 This case is fully interpretable from the function $f_i^{(r)}$ alone.

735 *Case 2: Region transitions occur in other features.* If for some $j \neq i$, the region function changes:
 736 $r_j(\mathbf{x}_{-j}^{+\varepsilon}) \neq r_j(\mathbf{x}_{-j})$, then the additional terms contribute:

$$\Delta f_j := f_j^{(r_j(\mathbf{x}_{-j}^{+\varepsilon}))}(x_j) - f_j^{(r_j(\mathbf{x}_{-j}))}(x_j).$$

737 Hence, the total change includes not only the direct shift in f_i , but also discrete region-based changes
 738 in other features and the value of $\Delta \hat{y}$ cannot be recovered from x_i alone. \square

739 The third property asks whether increasing a feature always increases the model's prediction, re-
 740 gardless of the values of other features. In other words, does the model treat the feature as globally
 741 monotonic?

742 **Proposition B.3** (Global Feature Monotonicity). *Let $\hat{y} : \mathbb{R}^d \rightarrow \mathbb{R}$ be a prediction function, and fix a*
 743 *feature index $i \in \{1, \dots, d\}$. Define the global monotonicity condition for feature x_i as follows: the*
 744 *model is said to be globally increasing in x_i if for all $\mathbf{x} \in \mathbb{R}^d$ and all $\delta > 0$,*

$$\hat{y}(\mathbf{x} + \delta \mathbf{e}_i) \geq \hat{y}(\mathbf{x}),$$

745 *and strictly increasing if the inequality is strict.*

746 *We analyze whether this property can be verified from the structure of the model.*

747 **(A) In GAM.** *Assume the model has the form:*

$$\hat{y}(\mathbf{x}) = \sum_{j=1}^d f_j(x_j),$$

748 *with each $f_j : \mathbb{R} \rightarrow \mathbb{R}$.*

749 *Then the model is globally increasing in x_i if and only if f_i is monotonically increasing.*

750 **(B) In GA²M.** *Assume the model has the form:*

$$\hat{y}(\mathbf{x}) = \sum_{j=1}^d f_j(x_j) + \sum_{j < k} f_{jk}(x_j, x_k).$$

751 *Then global monotonicity in x_i cannot be determined from f_i alone.*

752 **(C) In CALM.** *Assume the model has the form:*

$$\hat{y}(\mathbf{x}) = \sum_{j=1}^d f_j^{(r_j(\mathbf{x}_{-j}))}(x_j),$$

753 *where $r_j : \mathbb{R}^{d-1} \rightarrow \{1, \dots, R_j\}$ is a region selector, and $\mathbf{x}_{-j} := \mathbf{x} \setminus x_j$.*

754 *Then global monotonicity in x_i holds if the following two conditions are satisfied:*

- 755 1. *For all regions $r \in \{1, \dots, R_i\}$, the function $f_i^{(r)} : \mathbb{R} \rightarrow \mathbb{R}$ is increasing.*
- 756 2. *For all $j \neq i$, and all $x_j \in \mathbb{R}$, the function $x_i \mapsto f_j^{(r_j(\mathbf{x}_{-j}))}(x_j)$ is non-decreasing; i.e.,*
 757 *region transitions caused by changing x_i do not decrease f_j 's contribution.*

758 *Proof. (A) In GAM.* Since the model is additive and each term depends only on a single variable,
 759 we have:

$$\hat{y}(\mathbf{x} + \delta \mathbf{e}_i) - \hat{y}(\mathbf{x}) = f_i(x_i + \delta) - f_i(x_i).$$

760 Thus, $\hat{y}(\mathbf{x} + \delta \mathbf{e}_i) \geq \hat{y}(\mathbf{x})$ if and only if f_i is increasing.

761 **(B) In GA²M.** As in case (A), $f_i(x_i + \delta) - f_i(x_i)$ gives the direct contribution. However, for each
 762 $j \neq i$, the term $f_{ij}(x_i, x_j)$ also contributes. The total change is:

$$\hat{y}(\mathbf{x} + \delta \mathbf{e}_i) - \hat{y}(\mathbf{x}) = f_i(x_i + \delta) - f_i(x_i) + \sum_{j \neq i} [f_{ij}(x_i + \delta, x_j) - f_{ij}(x_i, x_j)].$$

763 The sign of this expression depends on the values of x_j , and therefore cannot be determined from
 764 f_i alone. Consequently, verifying global monotonicity in x_i requires knowing the full interaction
 765 structure and input values.

766 **(C) In CALM** Let $\mathbf{x} \in \mathbb{R}^d$ and $\delta > 0$. Define:

$$\Delta \hat{y} := \hat{y}(\mathbf{x} + \delta \mathbf{e}_i) - \hat{y}(\mathbf{x}).$$

767 Then:

$$\Delta \hat{y} = f_i^{(r_i(\mathbf{x}_{-i}))}(x_i + \delta) - f_i^{(r_i(\mathbf{x}_{-i}))}(x_i) + \sum_{j \neq i} \left[f_j^{(r_j(\mathbf{x}_{-j}^{+\delta}))}(x_j) - f_j^{(r_j(\mathbf{x}_{-j}))}(x_j) \right],$$

768 where $\mathbf{x}_{-j}^{+\delta}$ is obtained by replacing $x_i \mapsto x_i + \delta$ in \mathbf{x}_{-j} .

769 The first term is non-negative if $f_i^{(r)}$ is increasing for all regions r , and the region $r_i(\mathbf{x}_{-i})$ is fixed.

770 The second term is a sum over changes in other features' contributions due to possible changes
771 in their regions r_j . For the model to be globally increasing in x_i , all such contributions must be
772 non-negative. This requires that increasing x_i does not cause a decrease in the contribution of any
773 other feature x_j , i.e., region transitions must preserve or increase each f_j .

774 Therefore, global monotonicity in x_i requires both conditions above. Conversely, if both conditions
775 are satisfied, then each term in the sum is non-negative, implying $\Delta \hat{y} \geq 0$.

776

□

777 **C Detailed Experiments**

778 **C.1 Experimental Setup Details**

779 **Reporting convention.** As in Section 4, we report mean \pm standard deviation over 5-fold cross-
 780 validation.

Table 4: Classification Datasets

Dataset	Size	Attributes
Adult	45222	13
COMPAS	6167	9
HELOC	10459	23
MIMIC2	24508	17
Appendicitis	106	7
Phoneme	5404	5
SPECTF	349	44
Magic	19020	10
Bank	45211	16
Churn	5000	19

Table 5: Regression Datasets

Dataset	Size	Attributes
Bike Sharing	17379	11
California Housing	20640	8
Parkinsons Motor	5875	19
Parkinsons Total	5875	19
Seoul Bike	8465	14
Wine	6497	11
Energy	19735	28
CCPP	9568	4
Electrical	10000	13
Elevators	16599	18
No2	500	7
Sensory	576	11
Airfoil	1503	5
Skill Craft	3338	19
Ailerons	13750	39

781 **Datasets details.** Tables 4 and 5 summarize the 10 classification and 15 regression datasets used in
 782 our experiments. These datasets are sourced from various publicly available repositories. The UCI
 783 Machine Learning Repository [24] provides datasets including Adult, Magic, Bank, Bike Sharing,
 784 Parkinson’s Motor, Parkinson’s Total, Seoul Bike, Wine, CCPP and Skill Craft. OpenML [41]
 785 offers datasets such as Electrical, Elevators, No2, Sensory, Airfoil, and Ailerons. The Penn Machine
 786 Learning Benchmarks (PMLB) [37] includes datasets like Appendicitis, Phoneme, SPECTF, and
 787 Churn. Additional datasets used in our experiments include California Housing [33], MIMIC-II [39],
 788 COMPAS [34], HELOC [32], and Energy [5, 6].

789 **Data Preprocessing.** All input features are standardized using z-score normalization. For regression
 790 tasks, the target variable y is also standard scaled. The only exception is the GAMI-Net model, which
 791 requires input features to be scaled using min-max normalization. For RMSE, predictions and targets
 792 are inverse-transformed to the original scale prior to evaluation.

793 **Model Configurations.** Our black box models comprise a fully connected neural network, a
 794 random forest and an XGBoost ensemble. The deep network contains two hidden layers with 50 units
 795 each; ReLU activations are used for regression whereas a final sigmoid unit closes the classification
 796 variant. Training is carried out for 200 epochs with the Adam optimiser (learning rate=0.001), a batch
 797 size of 200 and either mean-squared error or binary cross-entropy loss, depending on the task. The
 798 random-forest baseline consists of 500 trees grown to a maximum depth of 25 with a minimum of
 799 three samples required at each leaf. For boosted trees we employ XGBoost with 300 boosting rounds,
 800 a learning rate of 0.1 and the log-loss evaluation metric for classification

801 For GAM models we consider NAM and EBM without interactions and Spline. The Neural Additive
 802 Model (NAM) builds one independent multilayer perceptron per feature; each sub-network has three
 803 hidden layers of 100,100 and 10 ReLU units followed by a linear output, and the additive sum is
 804 passed through a sigmoid for binary classification. Training uses Adam (learning rate 0.001), for 10
 805 epochs and a mini-batch size of 32. Spline-GAMs are implemented with PyGAM’s LinearGAM or
 806 LogisticGAM, allocating a single spline term to every feature while relying on PyGAM’s automatic
 807 smoothing. Finally, the Explainable Boosting Machine (EBM) is used without interaction terms by
 808 explicitly setting $n_interactions = 0$.

809 To capture pairwise interactions we experiment with three GA²M variants. EB²M extends the EBM
 810 by enabling interactions with a default strength of 0.9; NODE-GA²M activates its interaction mode

Table 6: Number of Interactions per Model in Classification Datasets

Dataset	CALM	EB ² M	NodeGA ² M	GAMINet
Adult	4.2±1.4	12.0±0.0	72.6±3.2	20.0±0.0
COMPAS	1.4±1.7	9.0±0.0	36.0±0.0	20.0±0.0
HELOC	1.1±1.6	21.0±0.0	163.8±4.4	20.0±0.0
MIMIC2	11.5±2.1	16.0±0.0	115.0±3.3	20.0±0.0
Appendicitis	4.0±1.7	7.0±0.0	21.0±0.0	6.2±8.1
Phoneme	8.6±1.6	5.0±0.0	10.0±0.0	10.0±0.0
SPECTF	0.0±0.0	40.0±0.0	273.0±14.4	20.0±0.0
Magic	8.5±2.4	9.0±0.0	44.0±0.6	20.0±0.0
Bank	9.2±2.0	15.0±0.0	104.2±3.2	20.0±0.0
Churn	14.0±0.0	18.0±0.0	137.4±2.9	20.0±0.0
Avg.	6.2	15.2	97.7	17.6

811 via the `ga2m=1` flag and restricts training to a five-minute time limit to ensure parity with the other
812 models; this cap was chosen to match the practical training budget of competing methods, and we
813 observed that NODE-GA²M’s training loss had largely plateaued within this window on most datasets
814 where the cap was reached; GAMI-Net is run with the default hyper-parameters.

815 For CALM, we configure the heterogeneity drop threshold to 0.2, determining whether a split is
816 considered statistically significant. The region detector, which measures heterogeneity, relies on
817 Partial Dependence Plots (PDP) for RF and XGBoost black-box models, while for DNNs we consider
818 both PDP and RHALE to measure heterogeneity. The detector evaluates 20 candidate split points per
819 feature. The GAM used within regions is an EBM without interactions and the CALM is applied on
820 top of all three black-box models (DNN, RF, XGBoost).

821 **Compatibility.** The GAMI-Net baseline depends on a native binary (`lib_ebmcore_mac_x64.dylib`)
822 compiled for `x86_64` architecture. As a result, it is not compatible with Apple Silicon Macs, and
823 running it on such systems will lead to an architecture mismatch error. This issue is limited to this
824 external model and does not affect any of the proposed methods or other baselines. We provide
825 instructions in the README file to guide such users on running all other methods except GAMI-Net.

826 **Computer Resources.** All experiments were conducted on an in-house server with cloud infras-
827 tructure equipped with an Intel(R) Core(TM) i9-10900X CPU @ 3.70GHz, 128 GB of RAM. No
828 GPU acceleration was utilized during these experiments.

829 C.2 Number of Interactions

830 **Number of Interactions (Classification).** For CALM, the number of interactions is computed
831 as $\sum_i |S_i|$, where S_i is the set of conditioning features appearing in the region tree T_i for feature
832 x_i (i.e., the set of features used as split variables in T_i). Table 6 reports the number of feature
833 interactions selected or used by each model across classification datasets. CALM consistently uses
834 significantly fewer interactions than GA²M-style baselines, often by a wide margin. In most datasets,
835 CALM activates less than a third of the interactions compared to NodeGA²M, and even fewer than
836 GAMINet’s default of 20. The numbers shown for CALM reflect actual detected regions with
837 interaction-specific behavior, confirming that its compact regionalization mechanism results in sparse
838 and interpretable models.

839 **Number of Interactions (Regression).** As shown in Table 7, CALM activates fewer feature
840 interactions than NodeGA²M and GAMI-Net in regression tasks, and is comparable to EB²M (15.1
841 vs. 14.1 on average). While GAMI-Net and EB²M use a fixed or near-fixed interaction set, and
842 NodeGA²M often selects over 100 interactions, CALM remains sparse across datasets.

Table 7: Number of Interactions per Model in Regression Datasets

Dataset	CALM	EB ² M	NodeGA ² M	GAMINet
Bike Sharing	19.3±2.7	10.0±0.0	53.6±1.4	20.0±0.0
California Housing	10.7±2.0	8.0±0.0	28.0±0.0	20.0±0.0
Parkinsons Motor	13.8±4.4	18.0±0.0	126.2±1.7	20.0±0.0
Parkinsons Total	14.4±5.8	18.0±0.0	129.2±3.9	20.0±0.0
Seoul Bike	20.4±1.0	13.0±0.0	83.2±0.7	20.0±0.0
Wine	13.9±3.3	10.0±0.0	54.6±0.5	20.0±0.0
Energy	49.2±4.4	26.0±0.0	192.0±4.5	4.0±8.0
CCPP	1.9±2.2	4.0±0.0	6.0±0.0	5.8±0.4
Electrical	0.0±0.0	12.0±0.0	68.6±3.0	12.0±0.0
Elevators	11.7±3.0	17.0±0.0	115.4±2.1	20.0±0.0
No2	12.8±1.9	7.0±0.0	21.0±0.0	20.0±0.0
Sensory	6.9±4.2	10.0±0.0	54.0±0.9	16.0±4.9
Airfoil	9.3±0.5	5.0±0.0	10.0±0.0	10.0±0.0
Skill Craft	0.0±0.0	18.0±0.0	135.6±2.4	16.0±8.0
Ailerons	41.7±2.6	36.0±0.0	232.4±5.1	20.0±0.0
Avg.	15.1	14.1	87.3	16.3

Table 8: Runtime (Seconds) for Classification Datasets

Dataset	BlackBox	GAM		CALM	GA ² M		
	XGB	NAM	EBM		EB ² M	NodeGA ² M	GAMINet
Adult	0.2±0.01	38±2	10±1	36±1	13±2	97±9	718±122
COMPAS	0.1±0.004	10±0.1	2±2	2±0.1	1±0.04	51±1	129±5
HELOC	0.2±0.01	31±6	2±2	17±1	2±0.2	57±0.1	249±60
MIMIC2	0.2±0.01	31±0.3	4±3	26±3	5±1	69±9	359±32
Appendicitis	0.1±0.001	5±0.1	2±4	2±0.1	1±0.5	24±0.1	6±5
Phoneme	0.1±0.01	6±1	3±4	3±3	4±1	61±4	166±38
SPECTF	0.1±0.01	32±0.3	3±4	3±3	4±5	25±0.2	180±13
Magic	0.2±0.003	28±18	5±5	17±3	7±4	98±7	489±63
Bank	0.2±0.04	44±1	5±1	42±3	13±1	98±17	994±265
Churn	0.2±0.01	19±3	2±1	10±4	3±2	56±2	257±32
Avg.	0.2	24	4	16	5	64	355

843 C.3 Practical Runtime

844 We report the total runtime of each method across all datasets in Tables 8 and 9. For CALM, the
845 reported time includes not only the regions detection fitting steps, but also the training time of the
846 underlying black-box model and the GAM component used within regions. While this inclusion
847 gives a complete view of end-to-end cost, it somewhat disadvantages CALM in comparison to other
848 models, whose runtimes reflect only their own training processes. In practice, the black-box or GAM
849 components can be selected to be lightweight and the blackbox model can be reused or pretrained
850 independently, making CALM’s region discovery step lightweight and modular.

851 Despite this conservative accounting, CALM’s runtime remains competitive. For example, it often
852 trains faster than or comparably to GA²M methods such as NODE-GA²M and GAMI-Net, which tend
853 to have high computational overhead. On small and medium-sized datasets like COMPAS, HELOC,
854 California Housing, and Wine, CALM runs in under a minute, showing that its interpretability gains
855 come with reasonable computational cost. On larger datasets such as Bank or Ailerons, CALM’s
856 runtime scales moderately but remains practical.

Table 9: Runtime (Seconds) for Regression Datasets

Dataset	BlackBox	GAM		CALM	GA ² M		
	XGB	NAM	EBM		EB ² M	NodeGA ² M	GAMINet
Bike Sharing	0.1±0.01	18±1	2±1	15±1	19±3	187±31	299±33
California Housing	0.2±0.01	15±0.1	4±2	13±1	16±0.5	200±46	677±147
Parkinsons Motor	0.3±0.003	21±6	5±2	17±2	44±4	193±34	536±88
Parkinsons Total	0.3±0.01	18±0.2	5±3	16±2	40±5	227±76	543±79
Seoul Bike	0.2±0.01	15±0.2	3±3	13±3	12±4	127±18	268±53
Wine	0.2±0.01	12±1	2±3	9±3	3±0.4	62±10	187±20
Energy	0.3±0.01	43±0.4	8±3	82±5	97±15	206±57	253±214
CCPP	0.2±0.004	6±0.1	5±4	8±5	14±2	133±18	88±8
Electrical	0.1±0.002	25±20	9±4	11±5	8±3	112±28	279±43
Elevators	0.3±0.1	25±0.4	7±4	35±5	43±8	181±38	773±269
No2	0.1±0.003	5±0.04	2±0.1	2±0.1	2±0.1	26±0.1	78±1
Sensory	0.1±0.001	10±1	4±6	7±5	2±0.4	27±0.5	74±16
Airfoil	0.1±0.001	4±0.03	3±6	5±6	7±1	51±12	76±34
Skill Craft	0.3±0.003	17±0.2	4±6	9±6	1±0.1	44±0.1	141±52
Ailerons	0.3±0.02	53±1	5±7	70±8	3±0.1	126±36	619±134
Avg.	0.2	19	5	21	21	127	326

857 C.4 Performance Metrics for all Experiments

858 **Performance Summary.** We evaluated CALM on 25 datasets using three model classes: DNN,
859 XGB, and RF. For each model and dataset, we compare CALM to both its corresponding GAM
860 baseline and the original black-box model.

861 In comparisons with GAMs, we match the model structure: for example, CALM-NAM is compared
862 directly to NAM, so the performance difference reflects only the benefit of introducing regional
863 modeling. In comparisons with the black-box, we report the accuracy of the *best* CALM-GAM.
864 For XGB and RF, this is selected across multiple GAM types (NAM, EBM, Spline) using a fixed
865 heterogeneity threshold of 0.2. For DNNs, we first identify the best heterogeneity modeling method
866 (PDP or RHALE) for each GAM, and then select the best-performing GAM among all types. While
867 the threshold is fixed in our experiments, we note that alternative thresholds could yield further
868 improvements.

869 The results cover six tables: Tables 10–12 report performance on regression datasets using DNN, XGB,
870 and RF respectively, while Tables 13–15 report on the same three model classes for classification
871 datasets. Across these settings, CALM consistently outperforms the corresponding GAM baseline. It
872 matches or exceeds the performance of its GAM counterpart in 97.8% and 100% of cases for DNNs
873 (Tables 10, 13), 97.8% and 96.7% for XGB (Tables 11, 14), and 95.6% and 96.7% for RF (Tables 12,
874 15). These results reflect direct, structure-matched comparisons between CALM and GAMs of the
875 same architecture.

876 When compared to the black-box models, CALM also achieves strong results. For regression, the
877 best CALM-GAM variant strictly outperforms the black-box on 9/15 datasets with DNN (with 2
878 ties; Table 10), on 3/15 with XGB (2 ties; Table 11), and on 3/15 with RF (1 tie; Table 12). In
879 classification, it beats the black-box in 90.0% of DNN cases (Table 13), 50.0% for XGB (Table 14),
880 and 60.0% for RF (Table 15).

881 Overall, CALM outperforms or matches the corresponding GAM in 219 out of 225 cases (97.3%) and
882 surpasses the black-box in 40 out of 75 cases (53.3%). The rare cases where CALM does not improve
883 over the GAM occur in both regression and classification tasks and almost always correspond to
884 scenarios where the GAM itself performs as well as or better than the black-box. These cases suggest
885 that the underlying function is already well captured by global additive effects, leaving little room for
886 further refinement through regional modeling.

Table 10: RMSE Score for Regression Datasets, DNN

Dataset	BlackBox	GAM			CALM - NAM		CALM - EBM		CALM - Spline	
	DNN	NAM	EBM	Spline	PDP	RHALE	PDP	RHALE	PDP	RHALE
Bike Sharing	42.952±1.460	101.968±1.309	100.211±1.010	100.876±0.868	63.034±0.883	64.941±2.566	58.444±1.637	60.966±2.983	59.868±1.564	62.108±2.989
California Housing	0.519±0.017	0.611±0.011	0.554±0.008	0.632±0.009	0.580±0.012	0.591±0.016	0.522±0.010	0.526±0.009	0.593±0.018	0.597±0.006
Parkinsons Motor	3.654±0.162	6.112±0.158	4.204±0.087	6.027±0.070	4.220±0.192	5.303±0.312	2.681±0.105	3.658±0.324	4.149±0.211	5.318±0.430
Parkinsons Total	5.139±0.102	7.897±0.102	4.847±0.107	7.519±0.072	5.456±0.182	6.776±0.381	3.587±0.030	4.319±0.222	5.408±0.365	6.827±0.209
Seoul Bike	241.967±5.081	320.225±4.377	303.685±3.855	302.020±4.701	269.277±6.414	273.241±8.379	249.457±2.407	245.880±3.153	246.855±3.936	243.926±4.433
Wine	0.700±0.015	0.719±0.009	0.703±0.012	0.736±0.039	0.712±0.005	0.711±0.012	0.695±0.011	0.695±0.012	0.729±0.027	0.734±0.035
Energy	77.139±2.916	89.799±2.665	85.225±2.489	86.498±2.384	89.304±2.480	88.600±2.179	84.347±2.358	83.972±2.557	84.929±1.868	83.526±2.653
CCPP	3.916±0.079	4.213±0.050	3.435±0.076	4.144±0.054	4.205±0.067	4.188±0.079	3.427±0.058	3.379±0.083	3.995±0.061	4.006±0.058
Electrical	0.051±0.009	0.040±0.005	0.018±0.011	0.095±0.005	0.040±0.005	0.040±0.005	0.018±0.011	0.012±0.007	0.094±0.005	0.088±0.003
Elevators	0.002±0.000	0.003±0.000	0.002±0.000	0.002±0.000	0.002±0.000	0.002±0.000	0.002±0.000	0.002±0.000	0.002±0.000	0.002±0.000
No2	0.522±0.012	0.503±0.021	0.492±0.034	0.480±0.032	0.507±0.016	0.501±0.025	0.492±0.030	0.498±0.034	0.485±0.028	0.492±0.033
Sensory	0.522±0.022	0.483±0.014	0.477±0.009	0.482±0.012	0.461±0.025	0.449±0.019	0.447±0.021	0.444±0.012	0.453±0.025	0.446±0.020
Airfoil	2.131±0.131	4.801±0.198	4.565±0.148	4.542±0.121	3.275±0.196	3.143±0.150	2.647±0.194	2.449±0.227	2.697±0.161	2.588±0.248
Skill Craft	1.231±0.195	0.925±0.025	0.901±0.021	2.515±3.092	0.923±0.028	0.918±0.023	0.901±0.021	0.905±0.018	1.652±1.323	2.345±2.792
Ailerons	0.0002±0.000	0.0002±0.000	0.0002±0.000	0.0002±0.000	0.0002±0.000	0.0002±0.000	0.0002±0.000	0.0002±0.000	0.0002±0.000	0.001±0.001

Table 11: RMSE Score for Regression Datasets, XGB

Dataset	BlackBox	GAM			CALM		
	XGB	NAM	EBM	Spline	NAM	EBM	Spline
Bike Sharing	39.346±1.375	101.968±1.309	100.211±1.010	100.876±0.868	59.732±1.414	55.667±1.220	56.656±1.320
California Housing	0.454±0.010	0.611±0.011	0.554±0.008	0.632±0.009	0.566±0.010	0.511±0.010	0.594±0.014
Parkinsons Motor	1.443±0.094	6.112±0.158	4.204±0.087	6.027±0.070	4.574±0.137	2.243±0.126	4.432±0.185
Parkinsons Total	1.863±0.079	7.897±0.102	4.847±0.107	7.519±0.072	6.109±0.361	2.968±0.089	5.638±0.288
Seoul Bike	209.587±3.474	320.225±4.377	303.685±3.855	302.020±4.701	270.365±7.049	238.912±1.643	237.705±2.658
Wine	0.622±0.012	0.719±0.009	0.703±0.012	0.736±0.039	0.703±0.011	0.693±0.016	0.737±0.043
Energy	67.967±2.579	89.799±2.665	85.225±2.489	86.498±2.384	87.416±2.894	83.088±1.974	82.747±2.421
CCPP	3.086±0.090	4.213±0.050	3.435±0.076	4.144±0.054	4.171±0.042	3.419±0.066	4.039±0.036
Electrical	0.037±0.015	0.040±0.005	0.018±0.011	0.095±0.005	0.040±0.005	0.018±0.011	0.095±0.005
Elevators	0.002±0.000	0.003±0.000	0.002±0.000	0.002±0.000	0.002±0.000	0.002±0.000	0.002±0.000
No2	0.473±0.026	0.503±0.021	0.492±0.034	0.480±0.032	0.498±0.022	0.492±0.036	0.479±0.031
Sensory	0.508±0.028	0.483±0.014	0.477±0.009	0.482±0.012	0.462±0.016	0.450±0.020	0.451±0.029
Airfoil	1.544±0.100	4.801±0.198	4.565±0.148	4.542±0.121	3.146±0.137	2.503±0.147	2.582±0.130
Skill Craft	0.938±0.019	0.925±0.025	0.901±0.021	2.515±3.092	0.925±0.025	0.901±0.021	2.167±2.261
Ailerons	0.0002±0.000	0.0002±0.000	0.0002±0.000	0.0002±0.000	0.0002±0.000	0.0002±0.000	0.0002±0.000

Table 12: RMSE Score for Regression Datasets, RF

Dataset	BlackBox	GAM			CALM		
	RF	NAM	EBM	Spline	NAM	EBM	Spline
Bike Sharing	43.247±1.037	101.968±1.309	100.211±1.010	100.876±0.868	61.174±1.544	57.158±1.288	57.823±1.393
California Housing	0.502±0.012	0.611±0.011	0.554±0.008	0.632±0.009	0.570±0.015	0.515±0.012	0.611±0.025
Parkinsons Motor	1.477±0.155	6.112±0.158	4.204±0.087	6.027±0.070	4.515±0.094	2.230±0.081	4.349±0.061
Parkinsons Total	1.869±0.243	7.897±0.102	4.847±0.107	7.519±0.072	5.761±0.259	3.044±0.068	5.393±0.162
Seoul Bike	225.572±3.086	320.225±4.377	303.685±3.855	302.020±4.701	270.514±8.146	243.910±2.580	240.471±4.337
Wine	0.618±0.012	0.719±0.009	0.703±0.012	0.736±0.039	0.708±0.006	0.691±0.011	0.726±0.036
Energy	69.406±2.781	89.799±2.665	85.225±2.489	86.498±2.384	88.448±2.607	84.741±2.520	84.690±2.546
CCPP	3.378±0.078	4.213±0.050	3.435±0.076	4.144±0.054	4.129±0.057	3.422±0.112	4.029±0.052
Electrical	0.009±0.009	0.040±0.005	0.018±0.011	0.095±0.005	0.040±0.005	0.018±0.011	0.095±0.005
Elevators	0.003±0.000	0.003±0.000	0.002±0.000	0.002±0.000	0.002±0.000	0.002±0.000	0.002±0.000
No2	0.466±0.033	0.503±0.021	0.492±0.034	0.480±0.032	0.481±0.023	0.482±0.031	0.484±0.029
Sensory	0.446±0.007	0.483±0.014	0.477±0.009	0.482±0.012	0.451±0.022	0.439±0.020	0.442±0.025
Airfoil	2.196±0.170	4.801±0.198	4.565±0.148	4.542±0.121	3.073±0.143	2.533±0.137	2.627±0.108
Skill Craft	0.916±0.019	0.925±0.025	0.901±0.021	2.515±3.092	0.926±0.037	0.901±0.022	2.320±2.696
Ailerons	0.0002±0.000	0.0002±0.000	0.0002±0.000	0.0002±0.000	0.0002±0.000	0.0002±0.000	0.0002±0.000

Table 13: Accuracy Score for Classification Datasets, DNN

Dataset	BlackBox	GAM			CALM - NAM		CALM - EBM		CALM - Spline	
	DNN	NAM	EBM	Spline	PDP	RHALE	PDP	RHALE	PDP	RHALE
Adult	0.849±0.002	0.851±0.002	0.870±0.002	0.856±0.001	0.853±0.001	0.854±0.003	0.870±0.003	0.871±0.003	0.857±0.003	0.859±0.003
COMPAS	0.682±0.011	0.682±0.016	0.681±0.012	0.685±0.012	0.682±0.007	0.681±0.014	0.683±0.011	0.679±0.012	0.683±0.015	0.683±0.011
HELOC	0.721±0.009	0.723±0.011	0.728±0.014	0.726±0.010	0.723±0.011	0.723±0.013	0.727±0.015	0.728±0.014	0.721±0.012	0.726±0.011
MIMIC2	0.885±0.003	0.886±0.001	0.886±0.003	0.886±0.003	0.887±0.001	0.886±0.001	0.886±0.003	0.886±0.002	0.886±0.003	0.886±0.003
Appendicitis	0.877±0.072	0.848±0.056	0.877±0.077	0.887±0.064	0.839±0.025	0.848±0.056	0.897±0.055	0.877±0.077	0.887±0.064	0.887±0.064
Phoneme	0.815±0.009	0.808±0.002	0.821±0.007	0.832±0.006	0.839±0.012	0.834±0.008	0.858±0.010	0.856±0.008	0.855±0.009	0.852±0.006
SPECTF	0.814±0.016	0.839±0.030	0.894±0.015	0.845±0.040	0.848±0.036	0.882±0.025	0.885±0.021	0.903±0.014	0.848±0.032	0.782±0.036
Magic	0.870±0.004	0.850±0.006	0.857±0.005	0.855±0.004	0.860±0.007	0.858±0.004	0.862±0.006	0.860±0.006	0.865±0.005	0.860±0.005
Bank	0.903±0.003	0.901±0.003	0.902±0.002	0.903±0.002	0.902±0.001	0.903±0.005	0.905±0.003	0.904±0.003	0.906±0.003	0.905±0.005
Churn	0.938±0.005	0.885±0.009	0.886±0.005	0.889±0.005	0.957±0.010	0.949±0.004	0.953±0.006	0.946±0.006	0.959±0.006	0.941±0.007

887 **D User Study**

888 **Study Overview.** We conducted a within-subjects user study ($N = 25$) to evaluate the practical
 889 interpretability of CALMs. Participants ranged from novices to experts in machine learning and data

Table 14: Accuracy Score for Classification Datasets, XGB

Dataset	BlackBox	GAM			CALM		
	XGB	NAM	EBM	Spline	NAM	EBM	Spline
Adult	0.870±0.002	0.851±0.002	0.870±0.002	0.856±0.001	0.853±0.002	0.870±0.002	0.858±0.001
COMPAS	0.661±0.007	0.682±0.016	0.681±0.012	0.685±0.012	0.686±0.014	0.684±0.013	0.686±0.016
HELOC	0.717±0.013	0.723±0.011	0.728±0.014	0.726±0.010	0.724±0.011	0.728±0.012	0.726±0.011
MIMIC2	0.890±0.001	0.886±0.001	0.886±0.003	0.886±0.003	0.886±0.001	0.886±0.003	0.886±0.003
Appendicitis	0.868±0.069	0.848±0.056	0.877±0.077	0.887±0.064	0.868±0.069	0.878±0.063	0.878±0.063
Phoneme	0.898±0.006	0.808±0.002	0.821±0.007	0.832±0.006	0.843±0.006	0.861±0.011	0.860±0.010
SPECTF	0.862±0.029	0.839±0.030	0.894±0.015	0.845±0.040	0.857±0.043	0.894±0.015	0.845±0.040
Magic	0.885±0.004	0.850±0.006	0.857±0.005	0.855±0.004	0.863±0.007	0.864±0.004	0.868±0.006
Bank	0.908±0.003	0.901±0.003	0.902±0.002	0.903±0.002	0.903±0.002	0.905±0.002	0.907±0.002
Churn	0.958±0.004	0.885±0.009	0.886±0.005	0.889±0.005	0.954±0.007	0.946±0.003	0.953±0.008

Table 15: Accuracy Score for Classification Datasets, RF

Dataset	BlackBox	GAM			CALM		
	RF	NAM	EBM	Spline	NAM	EBM	Spline
Adult	0.843±0.001	0.851±0.002	0.870±0.002	0.856±0.001	0.853±0.002	0.871±0.003	0.859±0.001
COMPAS	0.662±0.016	0.682±0.016	0.681±0.012	0.685±0.012	0.682±0.016	0.683±0.012	0.685±0.012
HELOC	0.724±0.014	0.723±0.011	0.728±0.014	0.726±0.010	0.723±0.014	0.728±0.014	0.726±0.010
MIMIC2	0.888±0.002	0.886±0.001	0.886±0.003	0.886±0.003	0.886±0.002	0.887±0.002	0.886±0.003
Appendicitis	0.859±0.064	0.848±0.056	0.877±0.077	0.887±0.064	0.848±0.056	0.877±0.077	0.887±0.064
Phoneme	0.898±0.005	0.808±0.002	0.821±0.007	0.832±0.006	0.838±0.005	0.859±0.005	0.859±0.009
SPECTF	0.877±0.027	0.839±0.030	0.894±0.015	0.845±0.040	0.848±0.026	0.894±0.018	0.842±0.048
Magic	0.878±0.004	0.850±0.006	0.857±0.005	0.855±0.004	0.861±0.004	0.864±0.004	0.866±0.005
Bank	0.901±0.004	0.901±0.003	0.902±0.002	0.903±0.002	0.905±0.003	0.907±0.002	0.907±0.002
Churn	0.957±0.005	0.885±0.009	0.886±0.005	0.889±0.005	0.952±0.005	0.951±0.007	0.955±0.006

890 analysis, and were asked to estimate feature contributions for specific inputs and reason about how
 891 predictions change under feature perturbations. Participants were not financially compensated.

892 **Participant Expertise.** Participants self-reported their level of expertise using a 5-point scale
 893 ranging from no experience to expert.

Figure 4: Self-reported expertise question used to assess participant background.

894 **Interpretation Protocol.** Before starting the tasks, participants were provided with standardized
 895 instructions on how to interpret conditional feature plots. These instructions explained how to: (i)
 896 select the appropriate curve based on feature conditions, (ii) read feature contributions directly from
 897 the plot, and (iii) account for interaction-induced discontinuities when estimating changes.

898 **Task Design.** Participants completed tasks targeting two aspects of interpretability: estimating
 899 feature contributions for specific inputs and reasoning about how predictions change under feature
 900 perturbations.

901 **Example Questions.** Figure 6 shows a representative example of the questions used in the study.
 902 The full questionnaire is available at <https://forms.gle/tTippi5eJ6cZHrfy9>.



Figure 5: Instruction panel provided to participants for interpreting conditional feature plots.

903 **Evaluation Metrics.** We evaluate performance using three metrics: task accuracy (percentage of
 904 correct answers), accuracy on prediction change tasks, and user experience, measured via self-reported
 905 ease of use and confidence.

906 **Results.** Participants achieved an average validation accuracy of 83.2% using CALM. On real-world
 907 tasks, CALM significantly outperformed heatmap-based representations in overall task performance
 908 (54% vs. 34%, $p < 0.05$), driven primarily by improvements in prediction change tasks (40% vs.
 909 12%, $p < 0.05$).

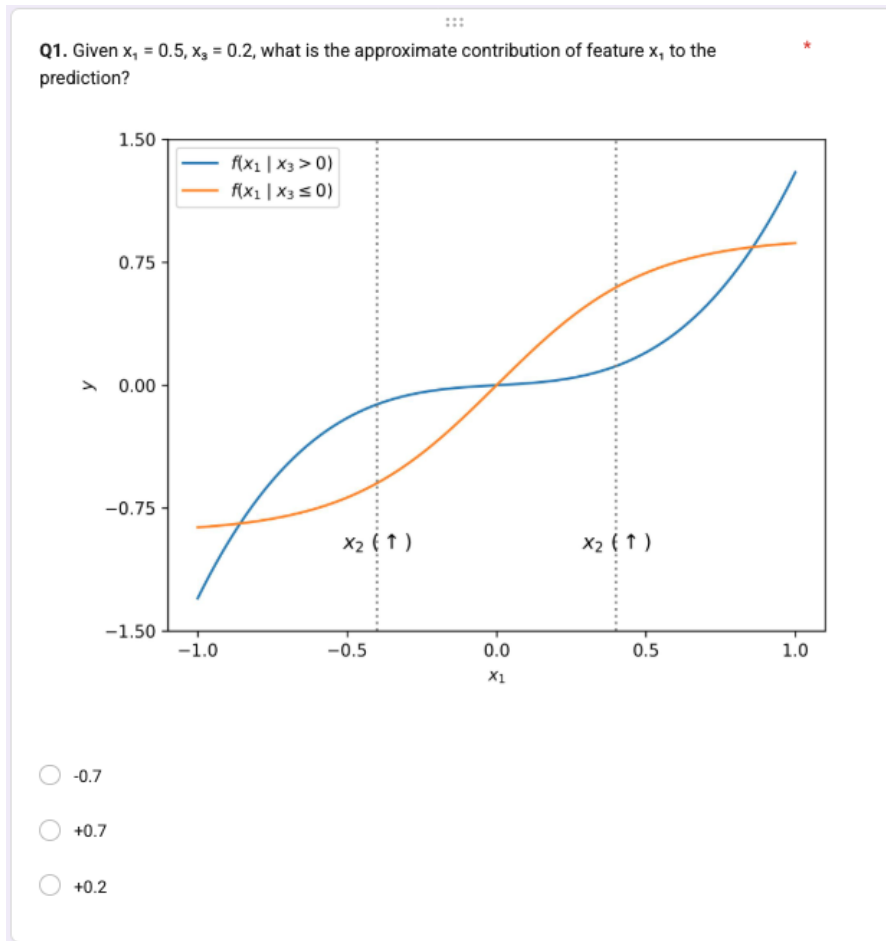


Figure 6: Example question requiring estimation of feature contribution from a conditional plot.

910 Participants also rated CALM as easier to use ($p < 0.01$) and more confidence-inducing ($p < 0.001$),
 911 indicating improved usability alongside stronger objective performance.

912 **NeurIPS Paper Checklist**

913 **1. Claims**

914 Question: Do the main claims made in the abstract and introduction accurately reflect the
915 paper’s contributions and scope?

916 Answer: [Yes]

917 Justification: The abstract and introduction accurately summarize the paper’s contributions,
918 including the proposed CALM model, its training procedure, theoretical analysis, and
919 empirical evaluation, while appropriately reflecting the scope of the results and supporting
920 claims on interpretability with both quantitative benchmarks and a user study.

921 Guidelines:

- 922 • The answer [N/A] means that the abstract and introduction do not include the claims
923 made in the paper.
- 924 • The abstract and/or introduction should clearly state the claims made, including the
925 contributions made in the paper and important assumptions and limitations. A [No] or
926 [N/A] answer to this question will not be perceived well by the reviewers.
- 927 • The claims made should match theoretical and experimental results, and reflect how
928 much the results can be expected to generalize to other settings.
- 929 • It is fine to include aspirational goals as motivation as long as it is clear that these goals
930 are not attained by the paper.

931 **2. Limitations**

932 Question: Does the paper discuss the limitations of the work performed by the authors?

933 Answer: [Yes]

934 Justification: The paper explicitly discusses key limitations, including the dependence on
935 the quality of the reference model, the potential reduction in interpretability with increasing
936 interactions, and the trade-off between model complexity and clarity.

937 Guidelines:

- 938 • The answer [N/A] means that the paper has no limitation while the answer [No] means
939 that the paper has limitations, but those are not discussed in the paper.
- 940 • The authors are encouraged to create a separate “Limitations” section in their paper.
- 941 • The paper should point out any strong assumptions and how robust the results are to
942 violations of these assumptions (e.g., independence assumptions, noiseless settings,
943 model well-specification, asymptotic approximations only holding locally). The authors
944 should reflect on how these assumptions might be violated in practice and what the
945 implications would be.
- 946 • The authors should reflect on the scope of the claims made, e.g., if the approach was
947 only tested on a few datasets or with a few runs. In general, empirical results often
948 depend on implicit assumptions, which should be articulated.
- 949 • The authors should reflect on the factors that influence the performance of the approach.
950 For example, a facial recognition algorithm may perform poorly when image resolution
951 is low or images are taken in low lighting. Or a speech-to-text system might not be
952 used reliably to provide closed captions for online lectures because it fails to handle
953 technical jargon.
- 954 • The authors should discuss the computational efficiency of the proposed algorithms
955 and how they scale with dataset size.
- 956 • If applicable, the authors should discuss possible limitations of their approach to
957 address problems of privacy and fairness.
- 958 • While the authors might fear that complete honesty about limitations might be used by
959 reviewers as grounds for rejection, a worse outcome might be that reviewers discover
960 limitations that aren’t acknowledged in the paper. The authors should use their best
961 judgment and recognize that individual actions in favor of transparency play an impor-
962 tant role in developing norms that preserve the integrity of the community. Reviewers
963 will be specifically instructed to not penalize honesty concerning limitations.

964 **3. Theory assumptions and proofs**

965 Question: For each theoretical result, does the paper provide the full set of assumptions and
966 a complete (and correct) proof?

967 Answer: [Yes]

968 Justification: The paper states the assumptions for its theoretical results and provides
969 complete proofs in the appendix, with references from the main text.

970 Guidelines:

- 971 • The answer [N/A] means that the paper does not include theoretical results.
- 972 • All the theorems, formulas, and proofs in the paper should be numbered and cross-
973 referenced.
- 974 • All assumptions should be clearly stated or referenced in the statement of any theorems.
- 975 • The proofs can either appear in the main paper or the supplemental material, but if
976 they appear in the supplemental material, the authors are encouraged to provide a short
977 proof sketch to provide intuition.
- 978 • Inversely, any informal proof provided in the core of the paper should be complemented
979 by formal proofs provided in appendix or supplemental material.
- 980 • Theorems and Lemmas that the proof relies upon should be properly referenced.

981 4. Experimental result reproducibility

982 Question: Does the paper fully disclose all the information needed to reproduce the main ex-
983 perimental results of the paper to the extent that it affects the main claims and/or conclusions
984 of the paper (regardless of whether the code and data are provided or not)?

985 Answer: [Yes]

986 Justification: The paper provides detailed descriptions of the model, training pipeline,
987 datasets, evaluation protocols, and hyperparameters, and includes code for reproducing all
988 experiments in the supplementary material.

989 Guidelines:

- 990 • The answer [N/A] means that the paper does not include experiments.
- 991 • If the paper includes experiments, a [No] answer to this question will not be perceived
992 well by the reviewers: Making the paper reproducible is important, regardless of
993 whether the code and data are provided or not.
- 994 • If the contribution is a dataset and/or model, the authors should describe the steps taken
995 to make their results reproducible or verifiable.
- 996 • Depending on the contribution, reproducibility can be accomplished in various ways.
997 For example, if the contribution is a novel architecture, describing the architecture fully
998 might suffice, or if the contribution is a specific model and empirical evaluation, it may
999 be necessary to either make it possible for others to replicate the model with the same
1000 dataset, or provide access to the model. In general, releasing code and data is often
1001 one good way to accomplish this, but reproducibility can also be provided via detailed
1002 instructions for how to replicate the results, access to a hosted model (e.g., in the case
1003 of a large language model), releasing of a model checkpoint, or other means that are
1004 appropriate to the research performed.
- 1005 • While NeurIPS does not require releasing code, the conference does require all submis-
1006 sions to provide some reasonable avenue for reproducibility, which may depend on the
1007 nature of the contribution. For example
 - 1008 (a) If the contribution is primarily a new algorithm, the paper should make it clear how
1009 to reproduce that algorithm.
 - 1010 (b) If the contribution is primarily a new model architecture, the paper should describe
1011 the architecture clearly and fully.
 - 1012 (c) If the contribution is a new model (e.g., a large language model), then there should
1013 either be a way to access this model for reproducing the results or a way to reproduce
1014 the model (e.g., with an open-source dataset or instructions for how to construct
1015 the dataset).
 - 1016 (d) We recognize that reproducibility may be tricky in some cases, in which case
1017 authors are welcome to describe the particular way they provide for reproducibility.
1018 In the case of closed-source models, it may be that access to the model is limited in

1019 some way (e.g., to registered users), but it should be possible for other researchers
1020 to have some path to reproducing or verifying the results.

1021 5. Open access to data and code

1022 Question: Does the paper provide open access to the data and code, with sufficient instruc-
1023 tions to faithfully reproduce the main experimental results, as described in supplemental
1024 material?

1025 Answer: [Yes]

1026 Justification: The paper provides open access to code and includes instructions for reproduc-
1027 ing the main experimental results, as described in the supplementary material.

1028 Guidelines:

- 1029 • The answer [N/A] means that paper does not include experiments requiring code.
- 1030 • Please see the NeurIPS code and data submission guidelines ([https://neurips.cc/
1031 public/guides/CodeSubmissionPolicy](https://neurips.cc/public/guides/CodeSubmissionPolicy)) for more details.
- 1032 • While we encourage the release of code and data, we understand that this might not
1033 be possible, so [No] is an acceptable answer. Papers cannot be rejected simply for not
1034 including code, unless this is central to the contribution (e.g., for a new open-source
1035 benchmark).
- 1036 • The instructions should contain the exact command and environment needed to run to
1037 reproduce the results. See the NeurIPS code and data submission guidelines ([https://
1038 neurips.cc/public/guides/CodeSubmissionPolicy](https://neurips.cc/public/guides/CodeSubmissionPolicy)) for more details.
- 1039 • The authors should provide instructions on data access and preparation, including how
1040 to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- 1041 • The authors should provide scripts to reproduce all experimental results for the new
1042 proposed method and baselines. If only a subset of experiments are reproducible, they
1043 should state which ones are omitted from the script and why.
- 1044 • At submission time, to preserve anonymity, the authors should release anonymized
1045 versions (if applicable).
- 1046 • Providing as much information as possible in supplemental material (appended to the
1047 paper) is recommended, but including URLs to data and code is permitted.

1048 6. Experimental setting/details

1049 Question: Does the paper specify all the training and test details (e.g., data splits, hyperpa-
1050 rameters, how they were chosen, type of optimizer) necessary to understand the results?

1051 Answer: [Yes]

1052 Justification: The paper specifies the experimental setup, including datasets, evaluation
1053 protocols, and key training details, with additional implementation and hyperparameter
1054 information provided in the appendix.

1055 Guidelines:

- 1056 • The answer [N/A] means that the paper does not include experiments.
- 1057 • The experimental setting should be presented in the core of the paper to a level of detail
1058 that is necessary to appreciate the results and make sense of them.
- 1059 • The full details can be provided either with the code, in appendix, or as supplemental
1060 material.

1061 7. Experiment statistical significance

1062 Question: Does the paper report error bars suitably and correctly defined or other appropriate
1063 information about the statistical significance of the experiments?

1064 Answer: [Yes]

1065 Justification: The paper reports results as mean \pm standard deviation over 5-fold cross-
1066 validation, clearly indicating variability due to data splits; statistical significance is addition-
1067 ally reported for the user study.

1068 Guidelines:

- 1069 • The answer [N/A] means that the paper does not include experiments.

- 1070 • The authors should answer [Yes] if the results are accompanied by error bars, confidence
1071 intervals, or statistical significance tests, at least for the experiments that support the
1072 main claims of the paper.
- 1073 • The factors of variability that the error bars are capturing should be clearly stated (for
1074 example, train/test split, initialization, random drawing of some parameter, or overall
1075 run with given experimental conditions).
- 1076 • The method for calculating the error bars should be explained (closed form formula,
1077 call to a library function, bootstrap, etc.)
- 1078 • The assumptions made should be given (e.g., Normally distributed errors).
- 1079 • It should be clear whether the error bar is the standard deviation or the standard error
1080 of the mean.
- 1081 • It is OK to report 1-sigma error bars, but one should state it. The authors should
1082 preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis
1083 of Normality of errors is not verified.
- 1084 • For asymmetric distributions, the authors should be careful not to show in tables or
1085 figures symmetric error bars that would yield results that are out of range (e.g., negative
1086 error rates).
- 1087 • If error bars are reported in tables or plots, the authors should explain in the text how
1088 they were calculated and reference the corresponding figures or tables in the text.

1089 8. Experiments compute resources

1090 Question: For each experiment, does the paper provide sufficient information on the com-
1091 puter resources (type of compute workers, memory, time of execution) needed to reproduce
1092 the experiments?

1093 Answer: [Yes]

1094 Justification: The paper reports the compute environment used for all experiments, including
1095 CPU type and memory, and clarifies that no GPU acceleration was used.

1096 Guidelines:

- 1097 • The answer [N/A] means that the paper does not include experiments.
- 1098 • The paper should indicate the type of compute workers CPU or GPU, internal cluster,
1099 or cloud provider, including relevant memory and storage.
- 1100 • The paper should provide the amount of compute required for each of the individual
1101 experimental runs as well as estimate the total compute.
- 1102 • The paper should disclose whether the full research project required more compute
1103 than the experiments reported in the paper (e.g., preliminary or failed experiments that
1104 didn't make it into the paper).

1105 9. Code of ethics

1106 Question: Does the research conducted in the paper conform, in every respect, with the
1107 NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines>?

1108 Answer: [Yes]

1109 Justification: The research complies with the NeurIPS Code of Ethics, involving standard
1110 machine learning methodology, use of publicly available datasets, and an anonymized user
1111 study conducted with informed consent.

1112 Guidelines:

- 1113 • The answer [N/A] means that the authors have not reviewed the NeurIPS Code of
1114 Ethics.
- 1115 • If the authors answer [No], they should explain the special circumstances that require a
1116 deviation from the Code of Ethics.
- 1117 • The authors should make sure to preserve anonymity (e.g., if there is a special consid-
1118 eration due to laws or regulations in their jurisdiction).

1119 10. Broader impacts

1120 Question: Does the paper discuss both potential positive societal impacts and negative
1121 societal impacts of the work performed?

1122
1123
1124
1125
1126
1127
1128
1129
1130
1131
1132
1133
1134
1135
1136
1137
1138
1139
1140
1141
1142
1143
1144
1145
1146
1147
1148
1149
1150
1151
1152
1153
1154
1155
1156
1157
1158
1159
1160
1161
1162
1163
1164
1165
1166
1167
1168
1169
1170
1171
1172
1173
1174

Answer: [Yes]

Justification: The paper discusses the positive impact of improved interpretability and outlines technical limitations of the approach, though broader societal risks are not extensively explored.

Guidelines:

- The answer [N/A] means that there is no societal impact of the work performed.
- If the authors answer [N/A] or [No], they should explain why their work has no societal impact or why the paper does not address societal impact.
- Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.
- The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate Deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
- The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
- If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

11. Safeguards

Question: Does the paper describe safeguards that have been put in place for responsible release of data or models that have a high risk for misuse (e.g., pre-trained language models, image generators, or scraped datasets)?

Answer: [N/A]

Justification: The paper does not involve the release of models or datasets that pose a high risk of misuse, as it focuses on a method for interpretable modeling on standard tabular data.

Guidelines:

- The answer [N/A] means that the paper poses no such risks.
- Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
- Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
- We recognize that providing effective safeguards is challenging, and many papers do not require this, but we encourage authors to take this into account and make a best faith effort.

12. Licenses for existing assets

Question: Are the creators or original owners of assets (e.g., code, data, models), used in the paper, properly credited and are the license and terms of use explicitly mentioned and properly respected?

Answer: [Yes]

Justification: The paper uses standard publicly available datasets and implementations, which are properly cited; while licenses are not explicitly listed, all assets are used in accordance with their standard terms.

1175
1176
1177
1178
1179
1180
1181
1182
1183
1184
1185
1186
1187
1188
1189
1190
1191
1192
1193
1194
1195
1196
1197
1198
1199
1200
1201
1202
1203
1204
1205
1206
1207
1208
1209
1210
1211
1212
1213
1214
1215
1216
1217
1218
1219
1220
1221
1222
1223
1224

Guidelines:

- The answer [N/A] means that the paper does not use existing assets.
- The authors should cite the original paper that produced the code package or dataset.
- The authors should state which version of the asset is used and, if possible, include a URL.
- The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.
- If assets are released, the license, copyright information, and terms of use in the package should be provided. For popular datasets, `paperswithcode.com/datasets` has curated licenses for some datasets. Their licensing guide can help determine the license of a dataset.
- For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.
- If this information is not available online, the authors are encouraged to reach out to the asset’s creators.

13. New assets

Question: Are new assets introduced in the paper well documented and is the documentation provided alongside the assets?

Answer: [Yes]

Justification: The paper introduces code for the proposed method, which is documented and provided in the supplementary material.

Guidelines:

- The answer [N/A] means that the paper does not release new assets.
- Researchers should communicate the details of the dataset/code/model as part of their submissions via structured templates. This includes details about training, license, limitations, etc.
- The paper should discuss whether and how consent was obtained from people whose asset is used.
- At submission time, remember to anonymize your assets (if applicable). You can either create an anonymized URL or include an anonymized zip file.

14. Crowdsourcing and research with human subjects

Question: For crowdsourcing experiments and research with human subjects, does the paper include the full text of instructions given to participants and screenshots, if applicable, as well as details about compensation (if any)?

Answer:[Yes]

Justification: The paper includes detailed descriptions of the user study, with instructions, screenshots, and representative questions provided in the appendix, along with access to the full questionnaire and includes details about participant compensation.

Guidelines:

- The answer [N/A] means that the paper does not involve crowdsourcing nor research with human subjects.
- Including this information in the supplemental material is fine, but if the main contribution of the paper involves human subjects, then as much detail as possible should be included in the main paper.
- According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other labor should be paid at least the minimum wage in the country of the data collector.

15. Institutional review board (IRB) approvals or equivalent for research with human subjects

1225 Question: Does the paper describe potential risks incurred by study participants, whether
1226 such risks were disclosed to the subjects, and whether Institutional Review Board (IRB)
1227 approvals (or an equivalent approval/review based on the requirements of your country or
1228 institution) were obtained?

1229 Answer: [No]

1230 Justification: The paper includes a user study but does not explicitly discuss IRB approval
1231 or risk disclosure procedures.

1232 Guidelines:

- 1233 • The answer [N/A] means that the paper does not involve crowdsourcing nor research
1234 with human subjects.
- 1235 • Depending on the country in which research is conducted, IRB approval (or equivalent)
1236 may be required for any human subjects research. If you obtained IRB approval, you
1237 should clearly state this in the paper.
- 1238 • We recognize that the procedures for this may vary significantly between institutions
1239 and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the
1240 guidelines for their institution.
- 1241 • For initial submissions, do not include any information that would break anonymity (if
1242 applicable), such as the institution conducting the review.

1243 16. Declaration of LLM usage

1244 Question: Does the paper describe the usage of LLMs if it is an important, original, or
1245 non-standard component of the core methods in this research? Note that if the LLM is used
1246 only for writing, editing, or formatting purposes and does *not* impact the core methodology,
1247 scientific rigor, or originality of the research, declaration is not required.

1248 Answer: [N/A]

1249 Justification: The research does not involve LLMs as part of the core methodology, experi-
1250 ments, or analysis.

1251 Guidelines:

- 1252 • The answer [N/A] means that the core method development in this research does not
1253 involve LLMs as any important, original, or non-standard components.
- 1254 • Please refer to our LLM policy in the NeurIPS handbook for what should or should not
1255 be described.