

# INTERPRETABLE CLASSIFICATION VIA A RULE NETWORK WITH SELECTIVE LOGICAL OPERATORS

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

We introduce the Rule Network with Selective Logical Operators (RNS), a novel neural architecture that employs **selective logical operators** to adaptively choose between AND and OR operations at each neuron during training. Unlike existing approaches that rely on fixed architectural designs with predetermined logical operations, our selective logical operators treat weight parameters as hard selectors, enabling the network to automatically discover optimal logical structures while learning rules. The core innovation lies in our **selective logical operators** implemented through specialized Logic Selection Layers (LSLs) with adaptable AND/OR neurons, a Negation Layer for input negations, and a Normal Form Constraint (NFC) to streamline neuron connections. We demonstrate that this selective logical operator framework can be effectively optimized using adaptive gradient updates with the Straight-Through Estimator to overcome gradient vanishing challenges. Through extensive experiments on 13 datasets, RNS demonstrates superior classification performance, rule quality, and efficiency compared to 25 state-of-the-art alternatives, showcasing the power of RNS in rule learning. Code and data are available at [https://anonymous.4open.science/r/RNS\\_-4A67/](https://anonymous.4open.science/r/RNS_-4A67/).

## 1 INTRODUCTION

Unlike explainable models (Rudin, 2019), which clarify black-box predictions by analyzing input feature contributions, interpretable models (Molnar, 2020) are inherently transparent, enabling direct human comprehension of their inference process (e.g., decision trees). This transparency is crucial for ensuring reliability, safety, and trust, especially in high-stakes domains like healthcare, finance, and law, where justifications for predictions are as important as the outcomes.

Rule-based models (Yin & Han, 2003; Frank & Witten, 1998; Cohen, 1995; Quinlan, 2014; Yang et al., 2017; Marton et al., 2023) are widely recognized for their inherent interpretability. A range of traditional approaches has been proposed, including example-based rule learning algorithms (Michalski, 1973), systems for learning first-order Horn clauses and Conjunctive Normal Form (CNF) or Disjunctive Normal Form (DNF) rules (Quinlan, 1990; Cohen, 1995; Frank & Witten, 1998; Michalski, 1973; Clark & Niblett, 1989; Mooney, 1995; Beck et al., 2023), ensemble methods and fuzzy rule systems (Ke et al., 2017; Breiman, 2001), and Bayesian frameworks (Letham et al., 2015; Wang et al., 2017; Yang et al., 2017). Despite extensive development, these traditional models often struggle with limitations in prediction accuracy (Quinlan, 2014; Loh, 2011; Cohen, 1995), suboptimal interpretability (Ke et al., 2017; Breiman, 2001), or poor scalability to large datasets (Letham et al., 2015; Wang et al., 2017; Yang et al., 2017).

Neural networks offer significant potential for representing and learning interpretable logical rules due to their expressiveness, generalization capability, robustness, and data-driven nature (Wang et al., 2020; 2021; 2024; Yu et al., 2023). This enables the automatic and efficient learning of complex rules at scale, combining the computational power of neural methods with the clarity and interpretability of rule-based reasoning.

**A major limitation in current rule-based neural networks is their reliance on fixed logical operators and non-learnable structural design.** These models enforce a predetermined architecture, with layers configured with static logical operations (either AND or OR) that cannot adapt during training (Wang et al., 2020; Beck et al., 2023). **The inability to dynamically select optimal logical operators** significantly compromises both model performance and rule quality. Furthermore, their

training methods are often heuristic-driven (Beck et al., 2023) or require repeated forward and backward passes (Wang et al., 2020; 2021; 2024), reducing efficiency. These constraints severely limit the models’ ability to discover the most appropriate logical structures for different datasets and generate sophisticated, high-quality rules for reliable interpretations.

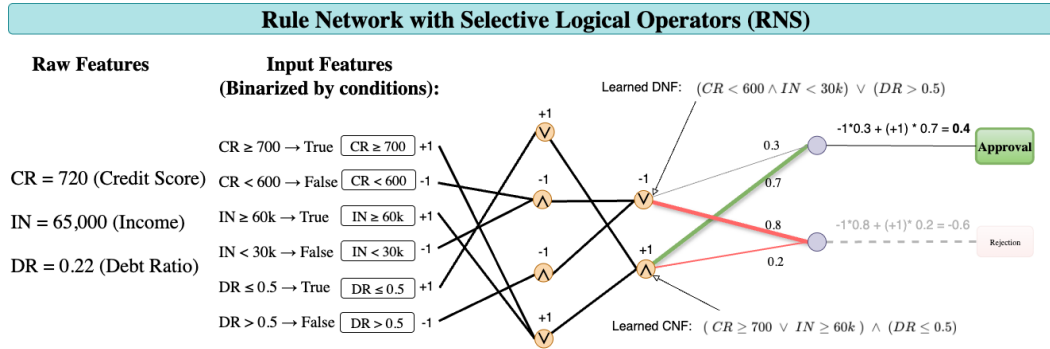


Figure 1: **Loan-application example of the RNS.** Raw attributes are binarized and passed through two learnable LSLs (AND/OR selected per neuron). The output layer learns class-specific weights from each clause to Approval/Rejection and sums them. Aggregating with learned weights gives Approval: 0.4 and Rejection:  $-0.6$ , yielding an *Approved* decision.

To address these fundamental limitations, we introduce the **Rule Network with Selective Logical Operators (RNS)**, which automatically learns neurons as AND or OR operators and simultaneously learns the connections between neurons, thereby forming CNF/DNF rules. RNS enables the network to uncover optimal logical structures while learning rules in both CNF and DNF formats, ensuring accurate classifications and interpretable insights through a transparent inference process. Figure 1 illustrates this process with a loan application example: raw features are binarized and passed through two learnable Logic Selection Layers (LSLs). The output layer assigns class-specific weights to each rule and aggregates them into final decisions. In this example, aggregation yields Approval 0.4 and Rejection  $-0.6$ , resulting in an *Approved* decision.

Key innovations of RNS include: two specialized Logic Selection Layers (LSLs), where weight parameters serve as hard selectors to determine AND or OR operations at each neuron; a Negation Layer for implementing logical negation; a Normal Form Constraint (NFC) to efficiently learn neuron connections; and the use of the Straight-Through Estimator (STE) for optimizing the discrete network. Through this design, RNS offers four unique advantages: (1) **Adaptive Logical Structure:** Neurons dynamically select logical operations, enabling flexible CNF or DNF rule learning; (2) **Complete Logic Operation:** Incorporates negation to achieve functional logic completeness; (3) **Efficient Rule Discovery:** Efficiently identifies optimal rules within the large search space of neuron connections; and (4) **Reliable Optimization:** Effectively optimizes discrete networks with non-differentiable components, mitigating the risk of gradient vanishing. Unlike existing approaches that rely on complex logical activation functions, RNS addresses the gradient vanishing problem through a lightweight design that combines STE with carefully designed logical activation functions.

Experiments on 13 datasets against 25 baselines demonstrate that RNS achieves superior performance in three critical aspects: (1) **prediction performance**, (2) **rule quality**, and (3) **training efficiency**. Notably, **RNS substantially outperforms state-of-the-art (SOTA) methods in rule quality, a factor even more critical than prediction accuracy for interpretable classification.** Its rule sets are not only more accurate and diverse but also considerably lower in complexity. Code is available at [https://anonymous.4open.science/r/RNS\\_-4A67/](https://anonymous.4open.science/r/RNS_-4A67/).

## 2 PRELIMINARIES

### 2.1 PROBLEM FORMULATION

A set of instances is represented as  $\mathcal{X}$ , where each instance  $\mathbf{x} \in \mathcal{X}$  is defined by a feature vector  $\mathbf{x} = [x_1, x_2, \dots, x_n]$ , comprising continuous or categorical features. Each instance is associated with a discrete class label  $y$ . The classification task aims to learn a function  $f : \mathbf{x} \rightarrow y$ . In this work, we design  $f$  as a rule-based model that automatically learns feature-based rules for prediction, with the learned rules providing inherent interpretability.

### 2.2 FEATURE BINARIZATION

Logical rules operate on Boolean truth values, so each feature must be represented in a True/False form. Since raw feature values cannot be directly evaluated as Boolean, we transform them into binary literals using conditions such as thresholds for continuous features or one-hot encodings for categorical features. For a categorical feature  $b_i$ , we use one-hot encoding to produce the corresponding binary vector  $\mathbf{b}_i$ . For a continuous feature  $c_j$ , we adopt the feature binning method from (Wang et al., 2021; 2024): a set of  $k$  upper bounds  $[\mathcal{H}_{j,1}, \dots, \mathcal{H}_{j,k}]$  and  $k$  lower bounds  $[\mathcal{L}_{j,1}, \dots, \mathcal{L}_{j,k}]$  is randomly sampled from the value range of  $c_j$ , and the binary representation of  $c_j$  is then derived as  $\tilde{\mathbf{c}} = [q(c_j - \mathcal{L}_{j,1}), \dots, q(c_j - \mathcal{L}_{j,k}), q(\mathcal{H}_{j,1} - c_j), \dots, q(\mathcal{H}_{j,k} - c_j)]$ , where  $q(x) = 1$  if  $x > 0$  and  $q(x) = -1$  otherwise. The full binarized input feature vector is  $\tilde{\mathbf{x}} = [\mathbf{b}_1, \dots, \mathbf{b}_p, \tilde{\mathbf{c}}_1, \dots, \tilde{\mathbf{c}}_t]$ . Beyond this random binning strategy, other heuristic- and learning-based methods can also be applied: AutoInt (Zhang et al., 2023) learns optimized bins jointly with model training, KInt (Dougherty et al., 1995) partitions values into clusters using K-means, and EntInt (Wang et al., 2020) selects bins that minimize label uncertainty. Designing an algorithm for how to bin features is beyond the scope of this work; instead, in Section I, we empirically compare different binning strategies with RNS.

### 2.3 NORMAL FORM RULES AS MODEL INTERPRETATION

Propositional logic is crucial to mathematical logic, focusing on propositions—statements with definite truth values – and logical connectives (e.g.,  $\wedge, \vee, \neg$ ) to form logical expressions. Among these, Conjunctive Normal Form (CNF) and Disjunctive Normal Form (DNF) are fundamental constructs: a formula  $z$  is in CNF if  $z = \bigwedge_i \bigvee_j l_{ij}$ , and in DNF if  $z = \bigvee_i \bigwedge_j l_{ij}$ , where each literal  $l_{ij}$  represents an atom or its negation. These standardized formats simplify logical deduction and analysis, facilitating efficient conversion of arbitrary logical expressions into forms suitable for both theoretical and practical applications. Building on this foundation, our work develops a logical rule-based classification model that predicts outcomes using automatically learned CNF and DNF expressions (rules). Specifically, we learn a set of logical rules  $\mathfrak{R} = \{z_1, \dots, z_m\}$ , where each rule  $z$  is constructed from binary features and their negations as literals. For each rule, we also learn a set of contribution scores  $\{s_{z,1}, \dots, s_{z,Y}\}$  that quantify its impact on each class  $y \in \{1, \dots, Y\}$ . Given an input  $\tilde{\mathbf{x}}$ , we evaluate the truth value of each rule in  $\mathfrak{R}$  and compute the logit for class  $k$  as  $\hat{y}_k = \sum_{i=1}^m z_i \times s_{z_i,k}$ . The learned rules thus provide a structured and mathematically transparent representation of the inference process, ensuring both accurate predictions and interpretability.

**Example: Loan Approval Decision.** Consider Figure 1, where an applicant with a credit score  $CR=720$ , income of  $IN=\$65,000$ , and debt ratio  $DR=0.22$  applies for a loan. The RNS first binarizes these features into logical conditions using learned thresholds, yielding the  $\pm 1$  encoded values:  $CR \geq 700 \rightarrow +1$ ,  $CR < 600 \rightarrow -1$ ,  $IN \geq 60k \rightarrow +1$ ,  $IN < 30k \rightarrow -1$ ,  $DR \leq 0.5 \rightarrow +1$ ,  $DR > 0.5 \rightarrow -1$ . Two LSL neurons then process these inputs: the first learns a DNF rule  $(CR < 600 \wedge IN < 30k) \vee (DR > 0.5)$  capturing rejection conditions, which evaluates to  $-1$  (false) since neither conjunction holds; the second learns a CNF rule  $(CR \geq 700 \vee IN \geq 60k) \wedge (DR \leq 0.5)$  capturing approval conditions, which evaluates to  $+1$  (true) as both clauses are satisfied. The output layer aggregates these logical decisions with learned weights, computing approval score  $(-1) \times 0.3 + (+1) \times 0.7 = 0.4$  and rejection score  $(-1) \times 0.8 + (+1) \times 0.2 = -0.6$ . Since the approval score exceeds the rejection score, the model outputs “Approved”—a decision fully traceable through interpretable logical rules rather than opaque neural computations.

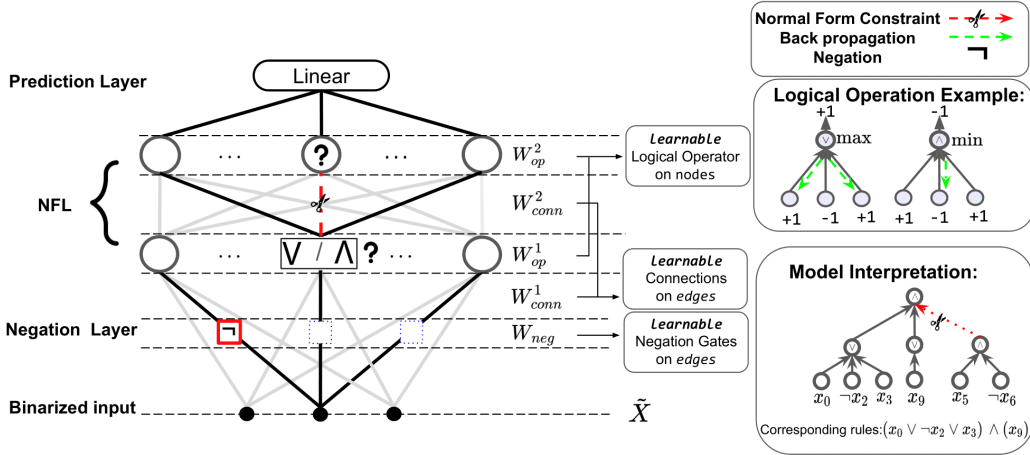


Figure 2: **RNS architecture and rule extraction.** From binarized inputs  $\tilde{X}$ , RNS learns interpretable CNF/DNF rules through a negation layer and two logic selection layers (NFL) constrained by normal form. The negation layer applies learnable gates ( $W_{neg}$ ) to select literals or their negations. NFL layers use learnable connection weights ( $W_{conn}^1, W_{conn}^2$ ) and operator weights ( $W_{op}^1, W_{op}^2$ ) to construct logical conjunctions and disjunctions. The normal form constraint enforces alternating operators between layers via masked back-propagation. A linear prediction layer aggregates rule activations for classification. The extracted rules (bottom-right) are human-readable logical formulas combining selected literals, e.g.,  $(x_0 \vee \neg x_2 \vee x_3) \wedge (x_9)$ .

### 3 METHOD

#### 3.1 OVERALL STRUCTURE

To achieve the goal in Section 2.3, we propose a rule network that supports logical rule computation and can be trained end-to-end. We identify four key challenges: (1) How to enable each neuron to dynamically select a logical operator? (2) How to incorporate negation, ensuring functional logic completeness? (3) How to efficiently learn logical connections within a large search space? (4) How to support effective optimization despite non-differentiable operations and gradient vanishing?

**Model Structure.** To address the first challenge, we design the Logic Selection Layer (LSL) with neurons that are learnable to select AND or OR logical operations, enabling connections to represent rules. To address the second challenge, a *Negation Layer* with learnable gates that flip the sign of input features when necessary. For the third challenge, we propose a *Normal Form Constraint* (NFC) that restricts valid connections across LSLs to reduce the search space. These components are trained using the proposed logical activation functions described in Section 3.4, which enables faithful logical computation with robust gradient flow. Non-differentiable steps such as  $\text{sign}(\cdot)$  are handled using the Straight-Through Estimator (STE). The overall architecture of RNS is shown in Figure 2.

Given an input sample, the input feature vector  $\mathbf{x}$  is binarized into  $\tilde{\mathbf{x}}$  and passed through the first LSL, where conjunction or disjunction operations are applied depending on the neuron’s operator selector. The Negation Layer operates on these connections to produce literals using original or negated features. The second LSL applies the next level of logic, generating a set of interpretable rules  $\mathfrak{R}$ , with each output neuron representing a learned logical clause. These rule outputs are then linearly combined with learned weights to produce class logits.

**Binary Neuron Values.** A key characteristic of RNS is that all neuron values (excluding predicted logits) are constrained to  $\pm 1$ , as is the binarized feature input described in Section 2.2. Unlike

existing works (Wang et al., 2020; 2021; 2024) that use 0/1, this design avoids the gradient vanishing problem, significantly improving the optimization process, detailed in Section 3.4 and Appendix K.

### 3.2 LOGIC SELECTION LAYER (LSL)

**Neuron Operator Selection.** To implement the nested structure of CNF and DNF, we stack two Logic Selection Layers (LSLs). Each LSL contains  $K$  neurons  $u_1, \dots, u_K$ , each as a learnable AND ( $\wedge$ ) or OR ( $\vee$ ) operator. The selection for the  $K$  neurons is parameterized by a weight vector  $\mathbf{w}_{op} \in \mathbb{R}^K$ . For a neuron  $u_i$ , the operator is determined by the sign of its weight:  $\tilde{\mathbf{w}}_{op}^i = \text{sign}(\mathbf{w}_{op}^i) \in \{+1, -1\}$ . The selection mechanism is defined as:

$$u_i = \begin{cases} \wedge, & \text{if } \tilde{w}_{op}^i = 1 \\ \vee, & \text{otherwise} \end{cases} \quad (1)$$

**Neuron Connection.** Similar to an MLP, we learn a weight matrix  $\mathbf{W}_{conn}$  to model connections between neurons in consecutive layers, where  $\mathbf{W}_{conn}^{i,j} \in \mathbb{R}$  denotes the weight from neuron  $u_j$  to  $u_i$ . RNS uses two such matrices: one for the (*input layer, 1st LSL*) and another for the (*1st LSL, 2nd LSL*). We binarize the weights via sign:  $\tilde{\mathbf{W}}_{conn}^{i,j} = \text{sign}(\mathbf{W}_{conn}^{i,j}) \in \{+1, -1\}$ . A connection exists if  $\tilde{\mathbf{W}}_{conn}^{i,j} = 1$ ; otherwise, it is inactive. The output  $v_i$  of a neuron  $u_i$  in the LSL is defined as:

$$v_i = \begin{cases} \bigwedge_{\tilde{\mathbf{W}}_{conn}^{i,j}=1} v_j, & \text{if } u_i = \wedge \\ \bigvee_{\tilde{\mathbf{W}}_{conn}^{i,j}=1} v_j, & \text{if } u_i = \vee \end{cases} \quad (2)$$

where  $v_j$  are the outputs of neurons from the previous layer connected to  $u_i$ .

**Logical Activation Functions.** In propositional logic with  $m$  inputs, the AND operation outputs +1 only if all connected neurons are +1, while the OR operation outputs +1 if at least one connected neuron is +1. Based on these principles, we use *min* and *max* functions as logical activation functions in RNS. The AND ( $\wedge$ ) operation is defined as the minimum of connected inputs, and the OR ( $\vee$ ) operation as the maximum:

$$v_i = \begin{cases} \min_{\tilde{\mathbf{W}}_{conn}^{i,j}=1} v_j, & \text{if } u_i = \wedge \\ \max_{\tilde{\mathbf{W}}_{conn}^{i,j}=1} v_j, & \text{if } u_i = \vee \end{cases} \quad (3)$$

Compared to the accumulative-multiplication-based activation functions proposed in prior works (Wang et al., 2020; 2021; 2024), this min-max-based method is more friendly for back-propagation and prevents gradient vanishing, detailed in Section 3.4.

### 3.3 NEGATION LAYER

Most existing works (Wang et al., 2020; 2021; 2024) support only {AND, OR}, which is not functionally complete for propositional logic (Mendelson, 1997; Enderton, 2001). To enable RNS to theoretically express any logical rule, we include the functionally complete set of operators {AND, OR, NEGATION} by introducing a negation layer to perform negation operations on features.

Since negation operates on input features, we apply the negation layer after the input layer. For each connection between  $v_j$  (from  $\tilde{\mathbf{x}}$ ) and a neuron  $u_i$  in the 1st LSL, we add a negation gate to decide whether to input the original value  $v_j$  or its negation  $-v_j$  to  $u_i$ . Each gate is parameterized by a weight  $\mathbf{W}_{neg}^{i,j} \in \mathbb{R}$ , and its binary version is obtained as  $\tilde{\mathbf{W}}_{neg}^{i,j} = \text{sign}(\mathbf{W}_{neg}^{i,j}) \in \{+1, -1\}$ . The negation operation is then defined as:  $Neg(v_j, \tilde{\mathbf{W}}_{neg}^{i,j}) = v_j \times \tilde{\mathbf{W}}_{neg}^{i,j}$ .

### 3.4 OPTIMIZATION

Learning a binary neural network is difficult because its discrete functions— $\text{sign}(\cdot)$ ,  $\text{max}(\cdot)$ , and  $\text{min}(\cdot)$ —are non-differentiable, making gradient computation challenging.

**Gradient of the Sign Function.** Inspired by the approach of searching for discrete solutions in continuous space (Courbariaux et al., 2015), we use the Straight-Through Estimator (STE) algorithm

to propagate gradients through non-differentiable operations during backpropagation. The STE assumes that the derivative of the sign function with respect to its input is 1, enabling the gradient to “pass through” the non-differentiable operation unchanged:  $\frac{\partial \text{sign}(x)}{\partial x} = 1$ .

**Gradient of Logical Activation Functions.** The max and min operations identify the maximum or minimum value across input neurons. While prior work (Lowe et al., 2022) primarily applies these functions in continuous pairwise scenarios, we extend them to handle discrete and multiple inputs. In pairwise cases, logical operators select one or two inputs, making gradient assignment and parameter updates straightforward. For multiple inputs, logical operators may involve several input neurons simultaneously, requiring careful gradient distribution to ensure proper updates. During the backward pass, uniform gradient updates are applied when multiple neurons share the same maximum or minimum value. This ensures fairness by equally distributing gradients among these neurons, promoting stability during optimization. Formally, for an input vector  $\mathbf{x}$  with elements  $x_1, x_2, \dots$ , the gradient of an input element  $x_i$  is:

$$\frac{\partial \max(\mathbf{x})}{\partial x_i} = \frac{1}{\sum_j \mathbb{I}(x_j = \max(\mathbf{x}))}, \quad \frac{\partial \min(\mathbf{x})}{\partial x_i} = \frac{1}{\sum_j \mathbb{I}(x_j = \min(\mathbf{x}))}. \quad (4)$$

**Gradient Vanishing.** Our design addresses the gradient vanishing problem that affects existing methods (Wang et al., 2020; 2021), caused by two main factors. First, these methods use binary states  $\{0, 1\}$ , resulting in many neurons outputting 0. These zero outputs propagate during the forward pass, nullifying gradients in the backward pass. Second, these methods rely on cumulative multiplications for logical activation functions like AND and OR. For instance, the AND operation is defined as:  $AND(\mathbf{x}) = \prod_i x_i$  when  $x_i \in \{0, 1\}$ . The gradient of the AND operation with respect to an input  $x_j$  is given by:  $\partial AND(\mathbf{x})/\partial x_j = \prod_{i \neq j} x_i$ . If any  $x_i = 0$ , then the gradient becomes 0.

In contrast, we adopt  $\pm 1$  for defining binary states in the neural network, preventing the generation of 0 outputs and 0 gradients. Furthermore, instead of designing complex continuous activation functions, our simple yet effective max/min functions avoid the use of multiplication. The straightforward uniform gradient update introduced in Equation 4 alleviates the problem of gradient vanishing. A more detailed analysis is provided in Appendix K.

### 3.5 NORMAL FORM CONSTRAINT (NFC)

Next, we turn our attention to a critical challenge in RNS – learning the connections between the two LSLs is computationally intensive. With  $K$  neurons in each LSL, learning and determining  $K^2$  potential connections significantly reduces the model’s efficiency and efficacy. To mitigate this, we design a normal form constraint that reduces the search space for learning the connections between two LSLs. Specifically, since CNF and DNF have a nested structure – where operations at the two levels must differ (CNF is a conjunction of disjunction rules, while DNF is a disjunction of conjunction rules) – we enforce a constraint that only neurons of different types from the two LSLs can be connected. For neurons  $u_i$  and  $u_j$  from two LSLs, we define a mask parameter  $M^{i,j}$  as:

$$M^{i,j} = \tilde{\mathbf{w}}_{op}^i \oplus \tilde{\mathbf{w}}_{op}^j = -\tilde{\mathbf{w}}_{op}^i \tilde{\mathbf{w}}_{op}^j, \quad (5)$$

Where  $\oplus$  denotes the XOR operation. A connection exists between  $u_i$  and  $u_j$  if  $M^{i,j} \times \tilde{\mathbf{W}}_{conn}^{i,j} = 1$ , otherwise, no connection is formed. During optimization, we update  $\tilde{\mathbf{W}}_{conn}^{i,j}$  only when  $M^{i,j} = 1$ , ensuring that connections are limited to valid neuron pairs and reducing the computational complexity of the learning process.

Assuming there are  $C_1$  and  $C_2$  conjunction neurons in two LSLs respectively, and  $D_1$  and  $D_2$  disjunction neurons respectively ( $C_1 + D_1 = C_2 + D_2 = K$ ), then potential connections under NFC are  $C_1 D_2 + C_2 D_1 \leq K^2$ . The empirical study in Section I demonstrates that NFC benefits both the efficiency and efficacy of the model.

## 4 EXPERIMENT

We conduct comprehensive experiments to evaluate RNS and answer the following research questions: **RQ1:** Learning high-quality rules is the core goal of an interpretable model. Can RNS learn high-quality rules? **RQ2:** How does RNS perform w.r.t. classification accuracy compared to SOTA

324 baselines? **RQ3**: How efficient is RNS in terms of model complexity and training time? **RQ4**:  
 325 How do the proposed Negation Layer, NFC, and different binning functions impact RNS? Note that  
 326 detailed ablation studies for RQ4 are provided in Appendix I. **RQ5**: What are the impacts of different  
 327 hyperparameters in RNS?  
 328

329 **4.1 EXPERIMENTAL SETTINGS**

330  
 331 **Datasets.** Extensive experiments are conducted on nine small datasets (adult, bank-marketing, ban-  
 332 knote, chess, c-4, letRecog, magic04, wine, tic-tac-toe) and four large datasets (activity, dota2,  
 333 facebook, fashion-mnist). These datasets are widely used for evaluating classification perfor-  
 334 mance (Letham et al., 2015; Wang et al., 2017; Yang et al., 2017). Details are provided in Table 5.  
 335 Datasets are categorized as ‘Discrete’ or ‘Continuous’ based on whether their features are exclusively  
 336 of one type, while datasets containing both feature types are classified as ‘Mixed’.

337 **Performance Evaluation.** We assess classification performance using the F1 score (Macro-average  
 338 for multi-class cases) with 5-fold cross-validation, reporting the average over five iterations. To pro-  
 339 vide a comprehensive evaluation, we compare models across all datasets using average rank (Demšar,  
 340 2006) (lower is better) and normalized mean (Marton et al., 2023) (higher is better) as metrics.

341  
 342 **Performance Evaluation.** We assess classification performance using the F1 score (Macro-average  
 343 for multi-class cases) with 5-fold cross-validation, reporting the average over five iterations. To pro-  
 344 vide a comprehensive evaluation, we compare models across all datasets using average rank (Demšar,  
 345 2006) (lower is better) and normalized mean (Marton et al., 2023) (higher is better) as metrics.  
 346 We evaluate RNS against two categories of models: interpretable and non-interpretable complex  
 347 methods. The interpretable models include rule-based approaches such as RRL (Wang et al., 2021;  
 348 2024), RIPPER (Cohen, 1995), CRS (Wang et al., 2020), PTURS (Yang & van Leeuwen, 2024),  
 349 DRNET (Qiao et al., 2021), LORD (Beck et al., 2023), C4.5 (Quinlan, 2014), CART (Breiman, 2017),  
 350 SBRL (Yang et al., 2017), CORELS (Angelino et al., 2018), Logistic Regression (LR) (Kleinbaum  
 351 et al., 2008), and KNN (Peterson, 2009). The non-interpretable models include BNN (Courbariaux  
 352 et al., 2016b), PLNN (Chu et al., 2018), SVM (Schölkopf & Smola, 2002), RF (Breiman, 2001),  
 353 LightGBM (LGBM) (Ke et al., 2017), XGBoost (XGB) (Chen & Guestrin, 2016), FT (Gorishniy  
 354 et al., 2021), SAINT (Somepalli et al., 2021), NODE (Popov et al., 2019), STG (Yamada et al., 2020),  
 355 TabNet (Arik & Pfister, 2021), TabTransformer (Huang et al., 2020), and VIME (Yoon et al., 2020).

Model	adult	bank	banknote	chess	c-4	letRecog	magic	tic-tac-toe	wine	activity	dota2	fb	fashion	N-Mean ↑	AvgRank ↓
<i>Non-interpretable Models</i>															
BNN	77.26	72.49	99.64	78.55	61.94	81.06	79.50	98.92	95.77	97.86	54.76	85.94	85.33	0.800(14)	15.31(15)
FT	79.01	77.04	99.93	80.64	72.45	97.17	85.95	97.84	94.63	98.56	59.70	86.52	89.23	0.843(7)	8.62(8)
LGBM	80.36	75.28	99.48	80.58	70.53	96.51	86.67	99.40	98.44	99.41	58.81	85.87	89.91	0.845(6)	7.46(6)
NODE	80.55	77.16	99.93	80.64	71.90	<b>97.20</b>	86.20	<b>100.0</b>	97.78	97.70	60.01	87.93	89.63	0.850(3)	5.77(3)
PLNN	73.55	72.40	<b>100.0</b>	77.85	64.55	92.34	83.07	<b>100.0</b>	76.07	98.27	59.46	89.43	89.36	0.806(13)	11.77(11)
RF	79.22	72.67	99.40	75.00	62.72	96.59	86.48	<b>100.0</b>	98.31	97.80	57.39	87.49	88.35	0.828(10)	10.69(10)
SAINT	79.31	75.60	99.04	79.37	72.85	96.72	85.46	96.95	95.50	98.94	59.58	88.79	89.69	0.842(8)	8.77(9)
STG	76.38	60.24	90.49	60.33	62.17	78.18	68.38	93.43	94.39	84.89	41.02	60.18	81.15	0.697(24)	23.15(26)
SVM	63.63	66.78	<b>100.0</b>	79.58	69.85	95.57	79.43	<b>100.0</b>	96.05	98.67	57.76	87.20	84.46	0.808(11)	12.54(12)
TabNet	80.94	77.54	96.34	80.78	72.13	93.44	85.08	<b>100.0</b>	98.37	96.07	59.16	86.53	88.98	0.840(9)	8.00(7)
TabTransformer	79.35	75.67	93.78	77.64	71.24	81.75	81.03	95.77	95.16	93.07	59.12	86.86	87.65	0.807(12)	14.85(14)
VIME	78.24	76.59	98.91	76.27	52.27	80.05	83.19	93.48	92.18	92.37	57.01	88.52	81.12	0.783(20)	17.08(18)
XGB	80.64	74.71	99.55	80.66	70.65	96.38	<b>86.69</b>	99.48	97.78	99.38	58.53	88.90	89.82	0.847(4)	6.38(4)
<i>Interpretable Models</i>															
C4.5	77.77	71.24	98.45	79.90	61.66	88.20	82.44	98.45	95.48	94.24	52.08	80.71	80.49	0.793(16)	17.31(19)
CART	77.06	71.38	97.85	79.15	61.24	87.62	81.20	97.85	94.39	93.35	51.91	81.50	79.61	0.787(18)	18.85(21)
CORELS	70.56	66.86	98.49	24.86	51.72	61.13	77.37	98.49	97.43	51.61	46.21	34.93	38.06	0.582(25)	22.08(25)
CRS	80.95	73.34	94.93	80.21	65.88	84.96	80.87	94.93	97.78	95.05	56.31	91.38	66.92	0.795(15)	13.85(13)
DRNET	<b>76.83</b>	<b>68.45</b>	<b>97.92</b>	<b>76.34</b>	<b>62.18</b>	<b>88.73</b>	<b>79.26</b>	<b>95.87</b>	<b>92.41</b>	<b>93.52</b>	<b>52.67</b>	<b>81.93</b>	<b>80.37</b>	<b>0.780(21)</b>	<b>19.54(22)</b>
KNN	77.57	75.61	<b>100.0</b>	75.21	65.18	91.92	77.92	92.33	96.24	97.09	51.61	68.61	82.35	0.785(19)	16.00(16)
LORD	80.72	74.90	99.51	80.61	70.77	96.32	86.52	99.45	97.85	99.32	58.61	88.75	89.76	0.847(5)	7.00(5)
LR	78.43	69.81	98.82	33.06	49.87	72.05	75.72	98.82	95.16	98.47	59.34	88.62	84.53	0.742(23)	16.54(17)
PTURS	77.92	71.84	96.73	77.58	65.43	89.37	76.18	95.42	90.85	94.76	53.82	84.27	82.19	0.789(17)	18.38(20)
RIPPER	74.69	69.76	96.00	70.95	64.78	92.94	77.92	97.79	89.16	88.08	55.67	64.18	78.66	0.758(22)	20.38(23)
RRL	80.42	77.18	<b>100.0</b>	79.66	72.01	96.14	86.24	<b>100.0</b>	98.37	98.96	60.08	90.11	89.64	0.851(2)	4.62(2)
SBRL	79.88	72.67	94.44	26.44	48.54	64.32	82.52	94.44	95.84	11.34	34.83	31.16	47.38	0.552(26)	21.15(24)
RNS (Ours)	<b>81.24</b>	<b>77.62</b>	<b>100.0</b>	<b>81.19</b>	<b>72.93</b>	95.82	86.68	<b>100.0</b>	<b>98.80</b>	<b>99.49</b>	<b>60.17</b>	<b>90.93</b>	<b>90.04</b>	<b>0.857(1)</b>	<b>1.77(1)</b>

374 Table 1: F1 scores (%) across 13 datasets. The top rows show non-interpretable models; the bottom  
 375 rows show interpretable models. Bold denotes best performance. N-Mean and AvgRank are computed  
 376 across all methods. Blue-highlighted rows indicate our additional results.  
 377

378  
379  
380  
381  
382  
383  
384  
385  
386  
387  
388  
389  
390  
391  
392  
393  
394  
395  
396  
397  
398  
399  
400  
401  
402  
403  
404  
405  
406  
407  
408  
409  
410  
411  
412  
413  
414  
415  
416  
417  
418  
419  
420  
421  
422  
423  
424  
425  
426  
427  
428  
429  
430  
431

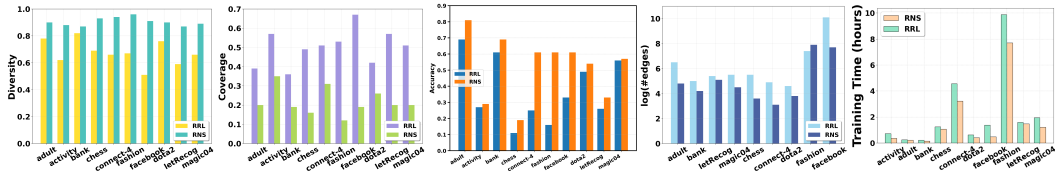


Figure 3: (a) Diversity. (b) Coverage. (c) Accuracy. (d) Avg Rule Length. (e) Training Time.

## 4.2 RULE QUALITY (RQ1)

The primary advantage of rule learning models over black-box methods is their ability to provide transparent and interpretable insights. However, this advantage is entirely dependent on the quality of the rules they generate. Learned rules must be of high quality, possessing characteristics like simplicity, accuracy, and generalizability. Such rules are crucial for fostering user trust, facilitating debugging, and discovering meaningful knowledge from data. Poorly constructed or overly complex rules defeat the purpose of using an interpretable model in the first place.

In this section, we evaluate the quality of the learned rules. Specifically, we benchmark RNS against RRL, the SOTA rule-based neural network, as a comparative baseline. The evaluation is conducted in two parts: (1) a comparison of the rules generated by these models using various existing **rule quality metrics**, and (2) a simulation experiment to assess the models’ ability to recover the underlying rules behind the data.

### 4.2.1 RULE QUALITY METRICS

Following prior work Yu et al. (2023); Lakkaraju et al. (2016), we evaluate rule quality using three key metrics: diversity, coverage, and single-rule accuracy. Extensive experiments are conducted across 10 datasets. **Accuracy** measures the prediction accuracy of a single rule for the instances it covers. As shown in Figure 3 and Figure 13, RNS achieves higher or comparable single-rule accuracy on nearly all datasets. **Coverage** quantifies the proportion of data instances covered by a rule. Lower coverage indicates that rules are more specific and easier for human experts to understand. As depicted in Figure 3 and Figure 12, RNS consistently yields rules with lower average coverage across most datasets, indicating that its rules focus on more localized, less redundant patterns. **Diversity** measures the overlap ratio between pairs of rules, with higher diversity reflecting that rules capture distinct, non-redundant logic. In Figure 3 and Figure 11, RNS consistently achieves higher diversity scores across all tested datasets. The diversity gap is especially prominent for more complex datasets such as *facebook* and *dota2*, demonstrating RNS’s ability to extract a set of rules that cover a wider variety of patterns with minimal redundancy.

In summary, RNS produces rule sets that are not only more accurate but also more diverse and interpretable due to their reduced coverage and increased diversity. This improvement in rule quality—across accuracy, coverage, and diversity—demonstrates the effectiveness of RNS in generating more meaningful and distinct rules compared to previous approaches.

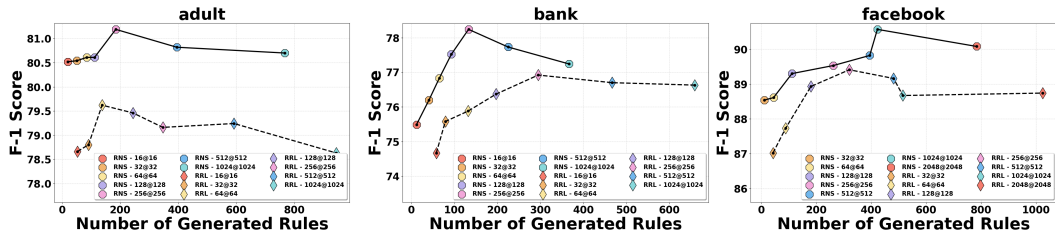


Figure 4: Comparison of F1 Scores for RNS and RRL with Different Numbers of Generated Rules.

Ground-Truth Rules	RNS	RRL	RRL w naive negation
$(x_1 \vee x_2) \wedge \neg x_3$	$x_1 \wedge \neg x_3, x_2 \wedge \neg x_3$	$x_1 \wedge x_2 \wedge x_3, x_1, x_2 \wedge x_3$	$x_1 \wedge \neg x_2 \wedge \neg x_3, x_1, \neg x_3$
$x_1 \vee (\neg x_2 \wedge \neg x_3)$	$\neg x_2 \wedge \neg x_3, x_1$	$x_3 \vee x_2 \vee (x_2 \vee x_1)$	$x_1, x_2, x_1 \wedge \neg x_2 \wedge x_3$
$x_1 \wedge \neg x_2 \wedge x_3$	$x_1 \wedge \neg x_2 \wedge x_3$	$x_1 \wedge x_3, x_2, (x_1 \wedge x_3) \wedge x_1$	$x_1, x_2, x_1 \wedge \neg x_2 \wedge x_3$
$x_1 \vee \neg x_2 \vee \neg x_3$	$x_1 \vee \neg x_2 \vee \neg x_3$	$x_2 \wedge x_3, (x_2 \wedge x_3) \wedge x_2$	$x_1 \wedge \neg x_3, \neg x_2 \wedge \neg x_3$

Table 2: Logical rules learned by RNS, RRL, and RRL with naive negation on synthetic data.

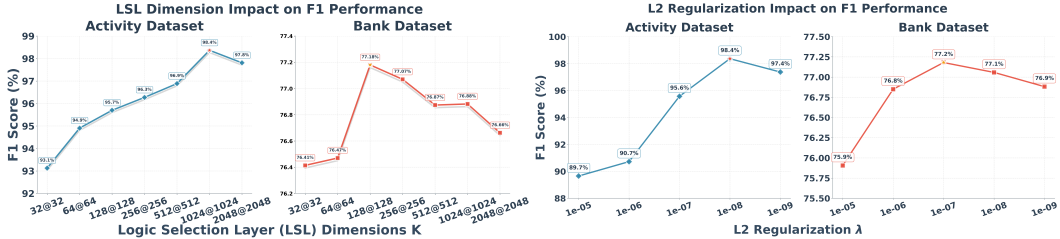


Figure 5: Impact of LSL dimension K and L2 regularization for both *activity* and *bank*.

### 4.3 CLASSIFICATION PERFORMANCE (RQ2)

Table 3 demonstrates the superior performance of RNS compared to both interpretable and non-interpretable baseline models across 13 datasets. RNS achieves the highest Normalized Mean (N-Mean) score of 0.997, outperforming all baselines, and its average rank of 1.69 further highlights its dominance, as it consistently ranks first or near the top across most datasets. Among interpretable baselines, RRL is the strongest competitor, yet RNS consistently outperforms it across the majority of datasets—for example, on the *bank* dataset, RNS achieves 77.62%, a 0.44% improvement over RRL, and on *fashion*, RNS attains 90.04%, compared to RRL’s 89.64%. Although complex models such as LightGBM and XGBoost are recognized for their strong performance, RNS surpasses them on several benchmarks, achieving 72.93% on the *c-4* dataset (exceeding LightGBM and XGBoost by 2.40% and 2.28%, respectively) and 99.49% on the *activity* dataset (outperforming LightGBM at 99.41% and XGBoost at 99.38%). These results highlight that RNS is not only interpretable but also matches or exceeds the performance of the strongest black-box baselines across diverse tabular datasets, with some baseline results drawn from Wang et al. (2024) to ensure fair comparison with a consistent experimental setup.

#### 4.3.1 SIMULATION EXPERIMENT

We conduct a controlled simulation experiment to assess RNS’s ability to learn the exact logical rules used to generate synthetic data.

**Setup.** We generate synthetic data using predefined logical rules to evaluate rule reconstruction ability. Three probability parameters  $\mathbf{p} = (p_1, p_2, p_3)$  correspond to feature variables  $\{x_1, x_2, x_3\}$ , each drawn from  $U(0, 1)$ . Using these as Bernoulli parameters, we generate  $X_{gen} \in \mathbb{R}^{3 \times 50000}$  through Bernoulli sampling, resulting in  $n = 50000$  binary vectors where each element is sampled from  $Bernoulli(p_i)$ . Labels are assigned based on logical rules (e.g.,  $x_1 \wedge x_2 \wedge x_3 \rightarrow 1$  assigns label 1 if all features are 1). We define four rule types presented in Table 2 and split the data evenly into training/test sets. Both RNS and RRL use logical layer dimensions of 64@64.

**Results.** Table 2 shows RNS achieves near 100% accuracy while recovering exact ground-truth rule structures. RRL produces overly simplified, redundant, or logically incomplete rules despite using naive negation settings. Detailed analysis is provided in Appendix H.

### 4.4 EFFICIENCY (RQ3)

We evaluate learning efficiency based on **the number of learned rules**, **the length of the learned rules**, and **computational time**. Large rule sets with lengthy conditions are difficult to interpret, so smaller rule sets with concise rules are preferred. Figure 4 and Figure 9 illustrate the relationship between the number of learned rules and the F1 score. RNS consistently outperforms RRL, achieving

Rules	Contribution to Negative Class	Contribution to Positive Class	Coverage
$(\text{duration} > 146.719 \wedge \neg \text{job} \in \{\text{housemaid}, \text{entrepreneur}\} \wedge \neg \text{month} \in \{\text{nov}, \text{mar}, \text{feb}\} \wedge \text{balance} < 1757.276) \vee$ $(\text{duration} < 850.975 \wedge \neg \text{job} = \text{technician} \wedge \text{month} = \text{mar}) \vee (\text{duration} > 146.719 \wedge \neg \text{month} = \text{jan} \wedge \neg \text{contact} = \text{unknown}) \vee$ $(\text{duration} > 19.057 \wedge \neg \text{age} > 67.87 \wedge \text{month} = \text{oct} \wedge \neg \text{day} < 5.569) \vee$ $(\neg \text{duration} < 850.975 \wedge \text{age} < 66.97 \wedge \neg \text{job} = \text{housemaid} \wedge \neg \text{month} = \text{jan} \wedge \neg \text{poutcome} = \text{success}) \vee$ $(\text{duration} > 59.781 \wedge \neg \text{job} = \text{unknown}) \vee (\text{duration} > 19.057 \wedge \neg \text{age} > 28.949) \vee$ $(\text{duration} > 59.781 \wedge \neg \text{age} > 67.87 \wedge \neg \text{job} \in \{\text{self-employed}, \text{housemaid}\} \wedge \neg \text{month} \in \{\text{oct}, \text{nov}\} \wedge \neg \text{loan} = \text{yes} \wedge \neg \text{default} = \text{yes})$	-0.4159	0.3929	0.4236
$(\text{duration} > 59.781 \wedge \neg \text{month} \in \{\text{jul}, \text{aug}\} \wedge \text{balance} < 1537.899 \wedge \text{previous} > 1.028 \wedge \text{pdays} < 229.326) \vee$ $(\neg \text{month} = \text{mar} \wedge \neg \text{poutcome} = \text{success} \wedge 102.083 < \text{pdays} < 304.124 \wedge \text{previous} < 6.754)$	0.3314	-0.3539	0.3617
$(\neg \text{month} = \text{aug} \wedge \neg \text{age} < 12.397 \wedge \text{campaign} > 0.066 \wedge \text{day} < 33.738 \wedge \text{pdays} < 244.292)$	-0.5581	0.5437	0.3319

Figure 6: Logical rules from RNS on the bank-marketing dataset.

higher performance with fewer rules. Additionally, Figure 3 and Figure 10 show that RNS produces shorter rules, averaging fewer than eight conditions, improving comprehensibility. In terms of computational time, as shown in Figure 3, RNS converges faster than the best baseline across all datasets. In general, RNS is more efficient than the baseline as it consistently produces fewer rules of a more concise size, using shorter training time.

#### 4.5 HYPERPARAMETER STUDY (RQ5)

We evaluate LSL dimension  $K$  and L2 regularization in Figure 5. For the *activity* dataset, F1 improves steadily with larger  $K$ , peaking at 1024@1024, while the *bank* dataset achieves its best performance earlier at 128@128. Similarly, moderate L2 values (around  $10^{-7}$ – $10^{-8}$ ) yield the highest scores, with *activity* reaching 98.4% and *bank* stabilizing near 77.2%. These results suggest that *activity* dataset benefits from larger architectures and balanced regularization, whereas *bank* saturates with smaller models and is less sensitive to  $\lambda$ .

#### 4.6 MODEL INTERPRETATION

RNS provides interpretable rules that capture meaningful patterns in the bank-marketing dataset, shown in Figure 14, such as seasonal effects, the link between longer calls and successful deposits, and demographic or age-related tendencies. These insights illustrate how RNS reflects plausible customer behaviors while remaining non-causal, as detailed in Appendix J.

## 5 ADVANTAGES AND LIMITATIONS OF RNS

**Advantages.** RNS advances interpretable rule learning via (i) learnable AND/OR operator selection per neuron, (ii) functional completeness with NOT via a Negation Layer, and (iii) robust optimization using max/min logical activations,  $\pm 1$  states, and STE. These yield superior prediction, rule quality, and efficiency over prior work (e.g., RRL).

**Limitations.** Although more efficient than prior neural rule learners, RNS is heavier than heuristic algorithms (e.g., C4.5, RIPPER). Future work includes improved scalability for extremely high-dimensional or streaming data and integrating domain priors.

## 6 CONCLUSION

We proposed RNS, a selective discrete neural network that learns CNF/DNF rules via two Logic Selection Layers with learnable AND/OR neurons, a Negation Layer, and an NFC for valid, efficient connections. With STE-enabled training and logical max/min activations, RNS achieves strong performance, rule quality, and efficiency across diverse datasets. Future directions include recommendations and text applications.

## REFERENCES

- 540  
541  
542 Elaine Angelino, Nicholas Larus-Stone, Daniel Alabi, Margo Seltzer, and Cynthia Rudin. Learning  
543 certifiably optimal rule lists for categorical data. *Journal of Machine Learning Research*, 18(234):  
544 1–78, 2018.
- 545 Sercan Ö Arik and Tomas Pfister. Tabnet: Attentive interpretable tabular learning. In *Proceedings of*  
546 *the AAAI Conference on Artificial Intelligence*, volume 35, pp. 6679–6687, 2021.
- 547 Florian Beck, Johannes Fürnkranz, and Van Quoc Phuong Huynh. Layerwise learning of mixed  
548 conjunctive and disjunctive rule sets. In *International Joint Conference on Rules and Reasoning*,  
549 pp. 95–109. Springer, 2023.
- 550 Leo Breiman. Random forests. *Machine learning*, 45:5–32, 2001.
- 551 Leo Breiman. *Classification and regression trees*. Routledge, 2017.
- 552  
553  
554 Adrian Bulat and Georgios Tzimiropoulos. Xnor-net++: Improved binary neural networks, 2019.
- 555  
556 Tianqi Chen and Carlos Guestrin. Xgboost: A scalable tree boosting system. In *Proceedings of the*  
557 *22nd acm sigkdd international conference on knowledge discovery and data mining*, pp. 785–794,  
558 2016.
- 559 Pengyu Cheng, Chang Liu, Chunyuan Li, Dinghan Shen, Ricardo Henao, and Lawrence  
560 Carin. Straight-through estimator as projected Wasserstein gradient flow. In *arXiv preprint*  
561 *arXiv:1910.02176*, 2019a.
- 562 Pengyu Cheng, Chang Liu, Chunyuan Li, Dinghan Shen, Ricardo Henao, and Lawrence Carin.  
563 Straight-through estimator as projected wasserstein gradient flow, 2019b.
- 564  
565 Lingyang Chu, Xia Hu, Juhua Hu, Lanjun Wang, and Jian Pei. Exact and consistent interpretation  
566 for piecewise linear neural networks: A closed form solution. In *Proceedings of the 24th ACM*  
567 *SIGKDD International Conference on Knowledge Discovery & Data Mining*, pp. 1244–1253,  
568 2018.
- 569 Peter Clark and Tim Niblett. The cn2 induction algorithm. *Mach. Learn.*, 3(4):261–283, mar 1989.  
570 ISSN 0885-6125. doi: 10.1023/A:1022641700528. URL [https://doi.org/10.1023/A:](https://doi.org/10.1023/A:1022641700528)  
571 [1022641700528](https://doi.org/10.1023/A:1022641700528).
- 572  
573 William W Cohen. Fast effective rule induction. In *Machine learning proceedings 1995*, pp. 115–123.  
574 Elsevier, 1995.
- 575  
576 Matthieu Courbariaux, Yoshua Bengio, and Jean-Pierre David. Binaryconnect: Training deep neural  
577 networks with binary weights during propagations. *Advances in neural information processing*  
578 *systems*, 28, 2015.
- 579 Matthieu Courbariaux, Yoshua Bengio, and Jean-Pierre David. Binaryconnect: Training deep neural  
580 networks with binary weights during propagations, 2016a.
- 581 Matthieu Courbariaux, Itay Hubara, Daniel Soudry, Ran El-Yaniv, and Yoshua Bengio. Binarized  
582 neural networks: Training deep neural networks with weights and activations constrained to +1 or  
583 -1, 2016b.
- 584  
585 Janez Demšar. Statistical comparisons of classifiers over multiple data sets. *The Journal of Machine*  
586 *learning research*, 7:1–30, 2006.
- 587 James Dougherty, Ron Kohavi, and Mehran Sahami. Supervised and unsupervised discretization of  
588 continuous features. In *Machine learning proceedings 1995*, pp. 194–202. Elsevier, 1995.
- 589 Anton Dries, Luc De Raedt, and Siegfried Nijssen. Mining predictive k-cnf expressions. *IEEE*  
590 *Transactions on Knowledge and Data Engineering*, 22(5):743–748, 2009.
- 591  
592 Herbert B Enderton. *A mathematical introduction to logic*. Elsevier, 2001.
- 593  
Eibe Frank and Ian H Witten. Generating accurate rule sets without global optimization. 1998.

- 594 Nicholas Frosst and Geoffrey Hinton. Distilling a neural network into a soft decision tree. *arXiv*  
595 *preprint arXiv:1711.09784*, 2017.
- 596
- 597 Yury Gorishniy, Ivan Rubachev, Valentin Khruikov, and Artem Babenko. Revisiting deep learning  
598 models for tabular data. *Advances in Neural Information Processing Systems*, 34:18932–18943,  
599 2021.
- 600 Tzung-Pai Hong and Shian-Shyong Tsang. A generalized version space learning algorithm for noisy  
601 and uncertain data. *IEEE Transactions on Knowledge and Data Engineering*, 9(2):336–340, 1997.
- 602
- 603 Xin Huang, Ashish Khetan, Milan Cvitkovic, and Zohar Karnin. Tabtransformer: Tabular data  
604 modeling using contextual embeddings. In *arXiv preprint arXiv:2012.06678*, 2020.
- 605 Arcchit Jain, Clément Gautrais, Angelika Kimmig, and Luc De Raedt. Learning cnf theories using  
606 mdl and predicate invention. In *IJCAI*, pp. 2599–2605, 2021.
- 607
- 608 Soochan Jeong, Rayadurgam Srikant, and Lei Ying. Finite sample analysis of the ODE method for  
609 stochastic approximation and reinforcement learning. In *International Conference on Machine*  
610 *Learning*, pp. 14846–14886. PMLR, 2023.
- 611 Guolin Ke, Qi Meng, Thomas Finley, Taifeng Wang, Wei Chen, Weidong Ma, Qiwei Ye, and Tie-Yan  
612 Liu. Lightgbm: A highly efficient gradient boosting decision tree. *Advances in neural information*  
613 *processing systems*, 30, 2017.
- 614
- 615 Minje Kim and Paris Smaragdīs. Bitwise neural networks, 2016.
- 616 Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint*  
617 *arXiv:1412.6980*, 2014.
- 618
- 619 David G Kleinbaum, K Dietz, M Gail, M Klein, and M Klein. Logistic regression. *A Self-Learning*  
620 *Tekst*, 2008.
- 621 Fayez Lahoud, Radhakrishna Achanta, Pablo Márquez-Neila, and Sabine Sūsstrunk. Self-binarizing  
622 networks, 2019.
- 623
- 624 Himabindu Lakkaraju, Stephen H Bach, and Jure Leskovec. Interpretable decision sets: A joint  
625 framework for description and prediction. In *Proceedings of the 22nd ACM SIGKDD international*  
626 *conference on knowledge discovery and data mining*, pp. 1675–1684, 2016.
- 627 Benjamin Letham, Cynthia Rudin, Tyler H. McCormick, and David Madigan. Interpretable classifiers  
628 using rules and bayesian analysis: Building a better stroke prediction model. *The Annals of*  
629 *Applied Statistics*, 9(3), September 2015. ISSN 1932-6157. doi: 10.1214/15-aos848. URL  
630 <http://dx.doi.org/10.1214/15-AOAS848>.
- 631
- 632 Wenmin Li, Jiawei Han, and Jian Pei. Cmar: Accurate and efficient classification based on multiple  
633 class-association rules. In *Proceedings 2001 IEEE international conference on data mining*, pp.  
634 369–376. IEEE, 2001.
- 635 Bing Liu, Wynne Hsu, and Yiming Ma. Integrating classification and association rule mining. In  
636 *Proceedings of the fourth international conference on knowledge discovery and data mining*, pp.  
637 80–86, 1998.
- 638
- 639 Zechun Liu, Baoyuan Wu, Wenhan Luo, Xin Yang, Wei Liu, and Kwang-Ting Cheng. Bi-real net:  
640 Enhancing the performance of 1-bit cnns with improved representational capability and advanced  
641 training algorithm, 2018.
- 642 Wei-Yin Loh. Classification and regression trees. *Wiley interdisciplinary reviews: data mining and*  
643 *knowledge discovery*, 1(1):14–23, 2011.
- 644
- 645 Scott C. Lowe, Robert Earle, Jason d’Eon, Thomas Trappenberg, and Sageev Oore. Logical activation  
646 functions: Logit-space equivalents of probabilistic boolean operators, 2022.
- 647
- Sascha Marton, Stefan Lüdtkke, Christian Bartelt, and Heiner Stuckenschmidt. Grande: Gradient-based  
decision tree ensembles. *arXiv preprint arXiv:2309.17130*, 2023.

- 648 Elliott Mendelson. *Schaum's outline of theory and problems of beginning calculus*. McGraw-Hill,  
649 1997.
- 650 Ryszard S. Michalski. Aqval/1—computer implementation of a variable-valued logic system v11  
651 and examples of its application to pattern recognition. In *International Joint Conference on*  
652 *Artificial Intelligence*, 1973. URL [https://api.semanticscholar.org/CorpusID:](https://api.semanticscholar.org/CorpusID:60492559)  
653 60492559.
- 654
- 655 Christoph Molnar. *Interpretable machine learning*. Lulu. com, 2020.
- 656
- 657 Raymond J Mooney. Encouraging experimental results on learning cnf. *Machine Learning*, 19:79–92,  
658 1995.
- 659 Giulia Pagallo and David Haussler. Boolean feature discovery in empirical learning. *Machine*  
660 *Learning*, 5:71–99, 1990. URL [https://api.semanticscholar.org/CorpusID:](https://api.semanticscholar.org/CorpusID:5661437)  
661 5661437.
- 662
- 663 Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor  
664 Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. Pytorch: An imperative style,  
665 high-performance deep learning library. *Advances in neural information processing systems*, 32,  
666 2019.
- 667 Ali Payani and Faramarz Fekri. Learning algorithms via neural logic networks. *arXiv preprint*  
668 *arXiv:1904.01554*, 2019.
- 669
- 670 Felix Petersen, Christian Borgelt, Hilde Kuehne, and Oliver Deussen. Deep differentiable logic  
671 gate networks. In *Advances in Neural Information Processing Systems (NeurIPS)*, volume 35, pp.  
672 2006–2019, 2022.
- 673
- 674 Leif E Peterson. K-nearest neighbor. *Scholarpedia*, 4(2):1883, 2009.
- 675
- 676 Sergei Popov, Stanislav Morozov, and Artem Babenko. Neural oblivious decision trees for deep  
677 learning on tabular data. In *Advances in Neural Information Processing Systems*, pp. 6285–6295,  
678 2019.
- 679
- 680 Litao Qiao, Weijia Wang, and Bill Lin. Learning accurate and interpretable decision rule sets from  
681 neural networks, 2021. URL <https://arxiv.org/abs/2103.02826>.
- 682
- 683 J. Ross Quinlan. Learning logical definitions from relations. *Machine Learning*, 5:239–266, 1990.  
684 URL <https://api.semanticscholar.org/CorpusID:6746439>.
- 685
- 686 J Ross Quinlan. *C4. 5: programs for machine learning*. Elsevier, 2014.
- 687
- 688 Mohammad Rastegari, Vicente Ordonez, Joseph Redmon, and Ali Farhadi. Xnor-net: Imagenet  
689 classification using binary convolutional neural networks, 2016.
- 690
- 691 Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin. ” why should i trust you?” explaining the  
692 predictions of any classifier. In *Proceedings of the 22nd ACM SIGKDD international conference*  
693 *on knowledge discovery and data mining*, pp. 1135–1144, 2016.
- 694
- 695 Cynthia Rudin. Stop explaining black box machine learning models for high stakes decisions and use  
696 interpretable models instead, 2019. URL <https://arxiv.org/abs/1811.10154>.
- 697
- 698 MICHALSKI R. S. On the quasi-minimal solution of the general covering problem. *Proceedings*  
699 *of the 5th International Symposium on Information Processing*, 3:125–128, 1969. URL [https:](https://cir.nii.ac.jp/crid/1571135649817801472)  
700 [://cir.nii.ac.jp/crid/1571135649817801472](https://cir.nii.ac.jp/crid/1571135649817801472).
- 701
- 702 Charbel Sakr, Jungwook Choi, Zhuo Wang, Kailash Gopalakrishnan, and Naresh Shanbhag. True  
703 gradient-based training of deep binary activated neural networks via continuous binarization. In  
704 *2018 IEEE international conference on acoustics, speech and signal processing (ICASSP)*, pp.  
705 2346–2350. IEEE, 2018.
- 706
- 707 Bernhard Schölkopf and Alexander J Smola. *Learning with kernels: support vector machines,*  
708 *regularization, optimization, and beyond*. MIT press, 2002.

- 702 Gowthami Somepalli, Micah Goldblum, Avi Schwarzschild, C Bayan Bruss, and Tom Goldstein.  
703 Saint: Improved neural networks for tabular data via row attention and contrastive pre-training.  
704 *arXiv preprint arXiv:2106.01342*, 2021.
- 705 William Sverdluk. *Dynamic version spaces in machine learning*. Wayne State University, 1992.
- 706  
707 Tong Wang, Cynthia Rudin, Finale Doshi-Velez, Yimin Liu, Erica Klampfl, and Perry MacNeille. A  
708 bayesian framework for learning rule sets for interpretable classification. *The Journal of Machine*  
709 *Learning Research*, 18(1):2357–2393, 2017.
- 710  
711 Zhuo Wang, Wei Zhang, Ning Liu, and Jianyong Wang. Transparent classification with multilayer  
712 logical perceptrons and random binarization, 2020.
- 713  
714 Zhuo Wang, Wei Zhang, Ning Liu, and Jianyong Wang. Scalable rule-based representation learning  
715 for interpretable classification, 2021.
- 716  
717 Zhuo Wang, Wei Zhang, Ning Liu, and Jianyong Wang. Learning interpretable rules for scalable data  
718 representation and classification. *IEEE Transactions on Pattern Analysis and Machine Intelligence*,  
719 46(2):1121–1133, February 2024. ISSN 1939-3539. doi: 10.1109/tpami.2023.3328881. URL  
<http://dx.doi.org/10.1109/TPAMI.2023.3328881>.
- 720  
721 Yutaro Yamada, Ofir Lindenbaum, Sahand Negahban, and Yuval Kluger. Feature selection using  
722 stochastic gates. In *International Conference on Machine Learning*, pp. 10648–10659. PMLR,  
723 2020.
- 724  
725 Hongyu Yang, Cynthia Rudin, and Margo Seltzer. Scalable bayesian rule lists. In *International*  
*conference on machine learning*, pp. 3921–3930. PMLR, 2017.
- 726  
727 Lincen Yang and Matthijs van Leeuwen. Probabilistic truly unordered rule sets, 2024. URL  
<https://arxiv.org/abs/2401.09918>.
- 728  
729 Penghang Yin, Jiancheng Lyu, Shuai Zhang, Stanley Osher, Yingyong Qi, and Jack Xin. Understand-  
730 ing straight-through estimator in training activation quantized neural nets, 2019.
- 731  
732 Xiaoxin Yin and Jiawei Han. Cpar: Classification based on predictive association rules. In *Proceed-*  
*ings of the 2003 SIAM international conference on data mining*, pp. 331–335. SIAM, 2003.
- 733  
734 Jinsung Yoon, Yao Zhang, James Jordon, and Mihaela van der Schaar. Vime: Extending the success of  
735 self-and semi-supervised learning to tabular domain. In *Advances in Neural Information Processing*  
736 *Systems*, volume 33, pp. 11033–11043, 2020.
- 737  
738 Lu Yu, Meng Li, Ya-Lin Zhang, Longfei Li, and Jun Zhou. Finrule: Feature interactive neural  
739 rule learning. In *Proceedings of the 32nd ACM International Conference on Information and*  
*Knowledge Management*, pp. 3020–3029, 2023.
- 740  
741 Wei Zhang, Yongxiang Liu, Zhuo Wang, and Jianyong Wang. Learning to binarize continuous features  
742 for neuro-rule networks. In *Proceedings of the Thirty-Second International Joint Conference on*  
743 *Artificial Intelligence*, pp. 4584–4592, 2023.
- 744  
745  
746  
747  
748  
749  
750  
751  
752  
753  
754  
755

## 7 ETHICS STATEMENT

This work presents RNS, a rule-based neural network for interpretable classification. We acknowledge and adhere to the ICLR Code of Ethics in all aspects of this research. Our work focuses on developing interpretable machine learning methods that can provide transparent decision-making processes, which aligns with ethical AI principles of explainability and accountability.

**Data and Experimental Ethics:** All datasets used in our experiments are publicly available benchmark datasets commonly used in machine learning research. We did not collect new data involving human subjects, and our use of existing datasets follows standard academic practices with proper attribution. The datasets include UCI Machine Learning Repository datasets and other established benchmarks that have been widely used in prior interpretable machine learning research.

**Potential Applications and Societal Impact:** While our method improves interpretability in machine learning, we acknowledge that rule-based models, like all AI systems, can potentially encode or amplify biases present in training data. We encourage practitioners to carefully evaluate fairness and bias when applying RNS to sensitive domains such as healthcare, finance, or criminal justice. The interpretable nature of our model can actually aid in identifying and addressing such biases by making decision logic transparent and auditable.

**Research Integrity:** We have conducted all experiments with scientific rigor, reported results honestly, and made efforts to ensure reproducibility by providing implementation details and planning to release code. We have appropriately cited prior work and clearly distinguished our contributions from existing methods. Our experimental comparisons use established baselines and evaluation metrics to ensure fair assessment.

## 8 USE OF LARGE LANGUAGE MODELS

We used large language models (LLMs) solely for text polishing and improving the clarity of our writing. All technical contributions, experimental design, implementation, analysis, and scientific insights are entirely our own work. The LLMs were not used for generating ideas, conducting experiments, writing code, or producing any substantive content. Specifically, LLMs assisted only with grammar corrections, sentence structure improvements, and enhancing the readability of our manuscript. The mathematical formulations, algorithmic innovations, experimental methodology, and all results interpretation remain exclusively the product of the authors' research efforts.

## 9 REPRODUCIBILITY STATEMENT

We have made extensive efforts to ensure the reproducibility of our results. **Implementation Details:** Section "Reproducibility" in the Appendix provides comprehensive hyperparameter settings, training procedures, and architectural specifications for RNS. We specify the number of Logic Selection Layer neurons (grid-searched from 32 to 4096), regularization coefficients ( $10^{-5}$  to  $10^{-9}$ ), binning strategies, batch sizes, learning rates, and training schedules for both small and large datasets.

**Code and Data Availability:** All datasets used are publicly available from the UCI Machine Learning Repository and other standard sources, with complete statistics provided in Table 2 of the Appendix. We will release our PyTorch implementation of RNS along with experimental scripts upon paper acceptance. The anonymous code repository is currently available at the URL provided in the paper.

**Experimental Reproducibility:** Our experiments use 5-fold cross-validation with multiple random seeds to ensure statistical reliability. We provide detailed baseline configurations following prior work citations, use established evaluation metrics (F1 scores, rule quality measures), and specify all preprocessing steps, including feature binarization methods. The simulation experiments use controlled synthetic data generation procedures that are fully specified to enable exact replication.

**Baseline Comparisons:** We use implementations and settings from established prior work, particularly following the high-quality codebase from Wang et al. (2024) for fair comparison with existing rule-based neural networks. All hardware specifications (NVIDIA A100 80GB GPU, Linux server) and software dependencies (PyTorch) are documented to facilitate reproduction of computational results.

## A REBUTTAL

This appendix addresses reviewer feedback and presents additional experimental results conducted during the rebuttal period.

### A.1 REBUTTAL - CONVERGENCE ANALYSIS (W3, REVIEWER MKDJ)

We analyze convergence for a simplified RNS module with a single Logic Selection Layer (LSL) operating on binary inputs  $\mathbf{x} \in \{\pm 1\}^d$ . Each neuron  $u_j$  selects inputs  $S_j \subset \{1, \dots, d\}$  and an operator via  $\sigma_j \in \mathbb{R}$ : if  $\sigma_j < 0$ ,  $u_j$  computes  $z_j(\mathbf{x}) = \min_{i \in S_j} x_i$  (AND); if  $\sigma_j > 0$ ,  $z_j(\mathbf{x}) = \max_{i \in S_j} x_i$  (OR). The network minimizes population loss  $L(\theta) = \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}}[\ell(f_\theta(\mathbf{x}), y)]$ .

**STE for min/max operations.** We use straight-through estimators for both min/max operations and the sign of  $\sigma_j$ . For  $z = \max(x_{i_1}, \dots, x_{i_k})$ , we set

$$\frac{\partial z}{\partial x_{i_m}} = \begin{cases} \frac{1}{|\mathcal{M}|}, & x_{i_m} = \max_{\ell} x_{i_\ell} \\ 0, & \text{otherwise,} \end{cases} \quad (6)$$

uniformly distributing gradients among all maximal inputs. For the selector, we treat  $\text{sign}(\sigma_j)$  as having unit derivative near zero (identity STE). Let  $\tilde{\nabla}L(\theta)$  denote the resulting coarse gradient and consider updates  $\theta_{t+1} = \theta_t - \eta \tilde{\nabla}L(\theta_t)$ .

Following convergence analyses for activation-quantized networks (Yin et al., 2019) and STE-based quantization-aware training (Jeong et al., 2023), we assume: (1) bounded data support with  $\pm 1$  inputs; (2) compact parameter set; (3)  $L$ -Lipschitz population gradient; (4) proper STE alignment:  $\langle \mathbb{E}[\tilde{\nabla}L(\theta)], \nabla L(\theta) \rangle \geq c \|\nabla L(\theta)\|^2$  for  $c > 0$ ; (5) step size  $0 < \eta < 2/L$ .

Under these conditions, coarse-gradient descent satisfies:  $L(\theta_{t+1}) \leq L(\theta_t)$  (monotone descent);  $\sum_{t=0}^{\infty} \|\nabla L(\theta_t)\|^2 < \infty$ , hence  $\|\nabla L(\theta_t)\| \rightarrow 0$ ; and every limit point is a critical point of  $L$ . When the target is exactly an AND/OR over the same inputs, convergence reaches the global minimum.

**Empirical Convergence Results.** Figure 15 shows training loss trajectories across eight diverse datasets, validating our theoretical guarantees. All datasets exhibit: (1) rapid initial descent indicating effective gradient flow; (2) near-monotonic decrease consistent with  $L(\theta_{t+1}) \leq L(\theta_t)$ ; (3) stable convergence to low loss values.

Convergence speeds vary by dataset characteristics: `wine` converges within 50 epochs due to its small size and clear separation; `dota2` stabilizes in 50-75 epochs; larger datasets like `fashion` (784 features, 10 classes) and `letRecog` require 200-300 epochs. The `bank` dataset shows minor oscillations at epochs 300-400 from mini-batch stochasticity on imbalanced data, but all datasets achieve stable convergence with  $\|\nabla L(\theta_t)\| \rightarrow 0$ .

These results confirm that proper STE design enables reliable gradient-based optimization of discrete logical networks across diverse problem scales and class distributions, avoiding the vanishing gradient issues of product-based activations.

### A.2 REBUTTAL - ADDITIONAL BASELINE COMPARISONS (REVIEWER XJK9)

We include two supplementary interpretable baselines, PTURS and DRNET, highlighted in blue in Table 3.

Table 3 demonstrates the superior performance of RNS compared to both interpretable and non-interpretable baseline models across 13 datasets. RNS achieves the highest Normalized Mean (N-Mean) score of 0.857, outperforming all baselines, and its average rank of 1.77 further highlights its dominance, as it consistently ranks first or near the top across most datasets. Among interpretable baselines, RRL is the strongest competitor with an N-Mean of 0.851 and average rank of 4.62, yet RNS consistently outperforms it across the majority of datasets—for example, on the `bank` dataset, RNS achieves 77.62%, a 0.44% improvement over RRL’s 77.18%, and on `fashion`, RNS attains 90.04%, compared to RRL’s 89.64%. Although complex non-interpretable models such as NODE (N-Mean: 0.850, AvgRank: 5.77), XGBoost (N-Mean: 0.847, AvgRank: 6.38), and LightGBM (N-Mean: 0.845, AvgRank: 7.46) are recognized for their strong performance, RNS surpasses them on several

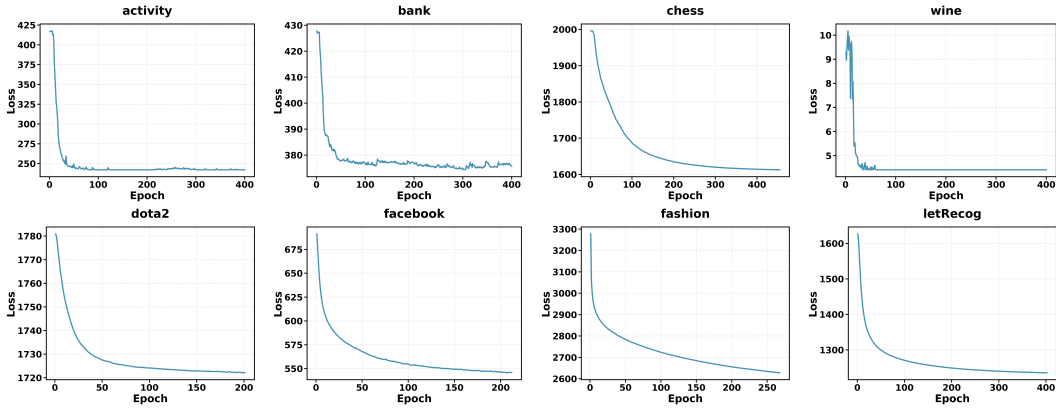


Figure 7: **Training loss convergence across diverse datasets.** Loss trajectories for RNS on eight datasets with varying sizes and complexities. All curves show rapid initial descent and stable convergence, validating our theoretical analysis.

benchmarks, achieving 72.93% on the c-4 dataset (exceeding LightGBM and XGBoost by 2.40% and 2.28%, respectively) and 99.49% on the activity dataset (outperforming LightGBM at 99.41% and XGBoost at 99.38%). Notably, RNS also outperforms LORD, another strong interpretable rule-based baseline (N-Mean: 0.847, AvgRank: 7.00), across most datasets, including achieving superior performance on adult (81.24% vs. 80.72%), bank (77.62% vs. 74.90%), chess (81.19% vs. 80.61%), c-4 (72.93% vs. 70.77%), and dota2 (60.17% vs. 58.61%). These results highlight that RNS is not only interpretable but also matches or exceeds the performance of the strongest black-box baselines across diverse tabular datasets, with some baseline results drawn from Wang et al. (2024) to ensure fair comparison with a consistent experimental setup.

Model	adult	bank	banknote	chess	c-4	letRecog	magic	tic-tac-toe	wine	activity	dota2	fb	fashion	N-Mean $\uparrow$	AvgRank $\downarrow$
<i>Non-interpretable Models</i>															
BNN	77.26	72.49	99.64	78.55	61.94	81.06	79.50	98.92	95.77	97.86	54.76	85.94	85.33	0.800(14)	15.31(15)
FT	79.01	77.04	99.93	80.64	72.45	97.17	85.95	97.84	94.63	98.56	59.70	86.52	89.23	0.843(7)	8.62(8)
LGBM	80.36	75.28	99.48	80.58	70.53	96.51	86.67	99.40	98.44	99.41	58.81	85.87	89.91	0.845(6)	7.46(6)
NODE	80.55	77.16	99.93	80.64	71.90	<b>97.20</b>	86.20	<b>100.0</b>	97.78	97.70	60.01	87.93	89.63	0.850(3)	5.77(3)
PLNN	73.55	72.40	<b>100.0</b>	77.85	64.55	92.34	83.07	<b>100.0</b>	76.07	98.27	59.46	89.43	89.36	0.806(13)	11.77(11)
RF	79.22	72.67	99.40	75.00	62.72	96.59	86.48	<b>100.0</b>	98.31	97.80	57.39	87.49	88.35	0.828(10)	10.69(10)
SAINT	79.31	75.60	99.04	79.37	72.85	96.72	85.46	96.95	95.50	98.94	59.58	88.79	89.69	0.842(8)	8.77(9)
STG	76.38	60.24	90.49	60.33	62.17	78.18	68.38	<b>93.43</b>	94.39	84.89	41.02	60.18	81.15	0.697(24)	23.15(26)
SVM	63.63	66.78	<b>100.0</b>	79.58	69.85	95.57	79.43	<b>100.0</b>	96.05	98.67	57.76	87.20	84.46	0.808(11)	12.54(12)
TabNet	80.94	77.54	96.34	80.78	72.13	93.44	85.08	<b>100.0</b>	98.37	96.07	59.16	86.53	88.98	0.840(9)	8.00(7)
TabTransformer	79.35	75.67	93.78	77.64	71.24	81.75	81.03	95.77	95.16	93.07	59.12	86.86	87.65	0.807(12)	14.85(14)
VIME	78.24	76.59	98.91	76.27	52.27	80.05	83.19	93.48	92.18	92.37	57.01	88.52	81.12	0.783(20)	17.08(18)
XGB	80.64	74.71	99.55	80.66	70.65	96.38	<b>86.69</b>	99.48	97.78	99.38	58.53	88.90	89.82	0.847(4)	6.38(4)
<i>Interpretable Models</i>															
C4.5	77.77	71.24	98.45	79.90	61.66	88.20	82.44	98.45	95.48	94.24	52.08	80.71	80.49	0.793(16)	17.31(19)
CART	77.06	71.38	97.85	79.15	61.24	87.62	81.20	97.85	94.39	93.35	51.91	81.50	79.61	0.787(18)	18.85(21)
CORELS	70.56	66.86	98.49	24.86	51.72	61.13	77.37	98.49	97.43	51.61	46.21	34.93	38.06	0.582(25)	22.08(25)
CRS	80.95	73.34	94.93	80.21	65.88	84.96	80.87	94.93	97.78	95.05	56.31	91.38	66.92	0.795(15)	13.85(13)
DRNET	76.83	68.45	97.92	76.34	62.18	88.73	79.26	95.87	92.41	93.52	52.67	81.93	80.37	0.780(21)	19.54(22)
KNN	77.57	75.61	<b>100.0</b>	75.21	65.18	91.92	77.92	92.33	96.24	97.09	51.61	68.61	82.35	0.785(19)	16.00(16)
LORD	80.72	74.90	99.51	80.61	70.77	96.32	86.52	99.45	97.85	99.32	58.61	88.75	89.76	0.847(5)	7.00(5)
LR	78.43	69.81	98.82	33.06	49.87	72.05	75.72	98.82	95.16	98.47	59.34	88.62	84.53	0.742(23)	16.54(17)
PTURS	77.92	71.84	96.73	77.58	65.43	89.37	76.18	95.42	90.85	94.76	53.82	84.27	82.19	0.789(17)	18.38(20)
RIPPER	74.69	69.76	96.00	70.95	64.78	92.94	77.92	97.79	89.16	88.08	55.67	64.18	78.66	0.758(22)	20.38(23)
RRL	80.42	77.18	<b>100.0</b>	79.66	72.01	96.14	86.24	<b>100.0</b>	98.37	98.96	60.08	90.11	89.64	0.851(2)	4.62(2)
SBRL	79.88	72.67	94.44	26.44	48.54	64.32	82.52	94.44	95.84	11.34	34.83	31.16	47.38	0.552(26)	21.15(24)
RNS (Ours)	<b>81.24</b>	<b>77.62</b>	<b>100.0</b>	<b>81.19</b>	<b>72.93</b>	95.82	86.68	<b>100.0</b>	<b>98.80</b>	<b>99.49</b>	<b>60.17</b>	<b>90.93</b>	<b>90.04</b>	<b>0.857(1)</b>	<b>1.77(1)</b>

Table 3: F1 scores (%) across 13 datasets. The top rows show non-interpretable models; the bottom rows show interpretable models. Bold denotes best performance. N-Mean and AvgRank are computed across all methods. Blue-highlighted rows indicate our additional results.

Ground-Truth Rules	RNS	RRL	RRL w/ naive negation
<i>Medium Complexity</i>			
$(x_1 \vee x_2) \wedge \neg x_3 \wedge x_4$	$x_1 \wedge \neg x_3 \wedge x_4, x_2 \wedge \neg x_3 \wedge x_4$	$x_1 \wedge x_2 \wedge x_4, x_1 \wedge x_4, x_2$	$x_1 \wedge x_4, x_2 \wedge \neg x_3, x_4$
$x_1 \vee (\neg x_2 \wedge x_3 \wedge \neg x_5)$	$x_1, \neg x_2 \wedge x_3 \wedge \neg x_5$	$x_1 \vee x_3 \vee (x_3 \vee x_5)$	$x_1, x_3, \neg x_2 \wedge x_3$
$(x_1 \wedge x_2) \vee (x_3 \wedge \neg x_4 \wedge x_5)$	$x_1 \wedge x_2, x_3 \wedge \neg x_4 \wedge x_5$	$x_1 \wedge x_2, x_3 \wedge x_5, x_1 \wedge x_2 \wedge x_3$	$x_1 \wedge x_2, x_3 \wedge x_5, \neg x_4$
<i>High Complexity</i>			
$(x_1 \wedge x_2 \wedge \neg x_3) \vee (x_4 \wedge x_5 \wedge x_7) \vee (\neg x_8 \wedge x_9)$	$x_1 \wedge x_2 \wedge \neg x_3, x_4 \wedge x_5 \wedge x_7, \neg x_8 \wedge x_9$	$x_1 \wedge x_2, x_4 \wedge x_5 \wedge x_7, x_7 \wedge x_9$	$x_1 \wedge x_2, x_4 \wedge x_5 \wedge x_7, x_9$

Table 4: Logical rules learned by RNS, RRL, and RRL with naive negation on extended synthetic data with more features. RNS consistently recovers the exact ground-truth structure, while RRL produces incomplete or overly simplified rules.

### A.3 REBUTTAL - ADDITIONAL SIMULATION EXPERIMENTS (W2, REVIEWER MKDJ; W5, REVIEWER J98E)

To address concerns about the scalability of rule recovery to more complex logical structures, we conduct additional simulation experiments beyond the 3-variable setting in Table 2. These experiments probe whether RNS maintains its advantage in exact rule recovery as problem dimensionality and rule complexity increase.

**Setup.** We extend the simulation framework from Section 4.3.1 to higher-dimensional settings with more features. For each complexity level, we generate synthetic datasets following the same procedure: probability parameters  $\mathbf{p} = (p_1, \dots, p_d)$  are sampled uniformly from  $U(0, 1)$ , and binary feature vectors are generated via Bernoulli sampling to produce  $n = 150,000$  samples. Labels are assigned based on predefined logical rules of increasing structural complexity, including nested disjunctions, conjunctions with multiple negations, and rules with longer conjunctive clauses. All models are configured with logical layer dimensions of 2048@2048, and datasets are split evenly into training and test sets.

**Results.** As shown in Table 4, RNS consistently recovers the exact ground-truth logical structure across both medium-complexity and high-complexity rules, demonstrating robust scalability. In contrast, both RRL variants struggle to capture complete logical structures as complexity increases.

For the medium complexity rules, RNS precisely recovers all ground-truth patterns, including complex structures like  $(x_1 \vee x_2) \wedge \neg x_3 \wedge x_4$  (decomposed into its two disjunctive components) and  $x_1 \vee (\neg x_2 \wedge x_3 \wedge \neg x_5)$ . RRL produces incomplete rules that miss critical negations or introduce spurious disjunctions (e.g.,  $x_1 \vee x_3 \vee (x_3 \vee x_5)$  instead of the correct structure). RRL with naive negation performs slightly better but still produces overly simplified rules that fail to capture all logical constraints.

For the high-complexity rules, RNS successfully decomposes the complex three-way disjunction  $(x_1 \wedge x_2 \wedge \neg x_3) \vee (x_4 \wedge x_5 \wedge x_7) \vee (\neg x_8 \wedge x_9)$  into its exact constituent clauses. In contrast, RRL systematically drops negation constraints (e.g., missing  $\neg x_3$ ) and introduces spurious conjunctions (e.g.,  $x_7 \wedge x_9$ ), while RRL with naive negation produces overly general rules that lose structural precision (e.g., reducing  $\neg x_8 \wedge x_9$  to just  $x_9$ ).

These results demonstrate that RNS’s advantage in exact rule recovery is not limited to simple, low-dimensional settings but extends to substantially more complex logical structures, validating its utility for interpretable learning in realistic scenarios.

### A.4 REBUTTAL - EXTENDED RULE QUALITY ANALYSIS (W3, REVIEWER MVR6; Q3, REVIEWER J98E)

In response to reviewer feedback, we extend our rule quality analysis to include additional interpretable baselines: C4.5 (decision tree), RIPPER (decision list), and LORD (layered rule learning). Figure 8 presents a comprehensive comparison on Adult and Bank datasets, revealing a fundamental incompatibility between our rule quality metrics and classical methods. Classical baselines exhibit metric degeneracy: uniformly high diversity and uniformly low coverage across both datasets. This degeneracy arises from structural constraints rather than rule quality. Decision trees partition the input space into disjoint regions where each leaf rule applies to a mutually exclusive subset, forcing diversity near 1 and coverage near  $1/\#\text{leaves}$ . Decision lists use first-match semantics that create similar effective partitions. Critically, these values are architecturally determined and provide

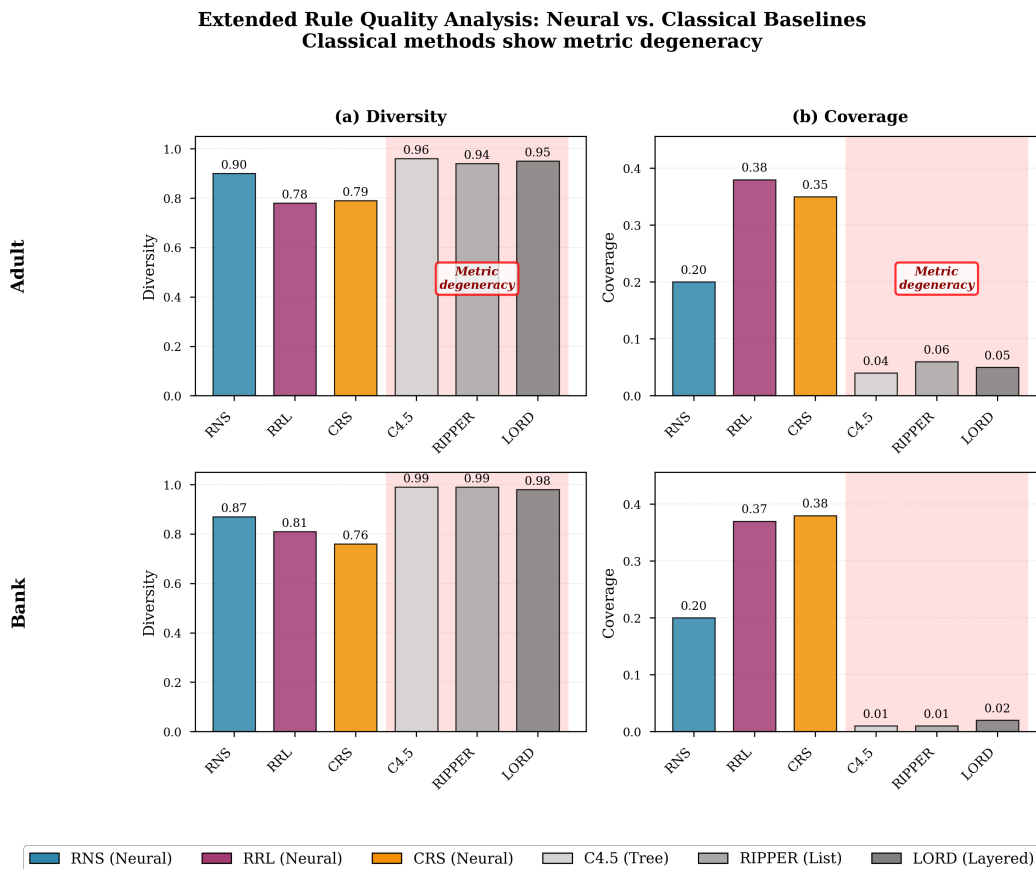


Figure 8: **Extended rule quality analysis comparing neural and classical baselines.** Classical methods (C4.5, RIPPER, LORD) exhibit metric degeneracy: uniformly high diversity and uniformly low coverage due to their disjoint rule semantics. These values provide no discriminative information about rule quality, making coverage and diversity meaningful only for parallel rule-set models (RNS, RRL, CRS) where rules can overlap.

no discriminative information, making these metrics uninformative for comparing rule quality in partition-based models.

In contrast, neural rule-set methods show meaningful variation: RNS achieves diversity of 0.87-0.90 and coverage of 0.20, while RRL shows 0.78-0.81 and 0.37-0.38, respectively. These differences reflect genuine architectural choices regarding rule specialization and coordination. Classical methods also achieve lower F1 scores (Table 3): on Adult, RNS achieves 81.24% compared to 77.77% (C4.5), 74.69% (RIPPER), and 80.72% (LORD); on Bank, 77.62% compared to 71.24%, 69.76%, and 74.90%. This 3-7% performance gap indicates quality issues beyond metric design. We conclude that coverage and diversity metrics are only meaningful for parallel rule-set models (RNS, RRL, CRS) where rules can fire simultaneously, explaining why the main paper focuses rule quality comparison on these methods. Including classical baselines would not strengthen the comparison—it would simply demonstrate that partition-based architectures yield uninformative metric values by design while also exhibiting lower predictive performance.

#### A.5 REBUTTAL - ADDITIONAL QUALITATIVE EXAMPLES (W1, REVIEWER MKDJ)

**Learned Rule from Adult Dataset.** We examine a learned rule from the Adult income dataset to demonstrate that RNS naturally discovers interpretable CNF (Conjunctive Normal Form) structures. The following rule predicts low income ( $\leq 50K$ ) through a conjunction of disjunctive clauses:

1026  
 1027  $(\text{capital-gain} < 17,939 \vee \text{education-num} < 12.8) \wedge$   
 1028  $(\neg \text{Exec-managerial} \vee \text{capital-gain} < 15,318) \wedge$   
 1029  $(\neg \text{Prof-specialty} \vee \text{education-num} < 12.8)$   
 1030

1031 Each clause constrains a wealth-building mechanism: (1) requiring either low capital gains or limited  
 1032 education, (2) excluding executives with substantial investments ( $> \$15,000$ ), and (3) excluding  
 1033 professional specialists with college degrees. The conjunctive structure encodes that low income  
 1034 results from lacking multiple advantages simultaneously—neither education nor capital gains nor  
 1035 high-status occupation alone suffices to escape this prediction. This demonstrates RNS’s ability to  
 1036 capture interpretable domain structure where compound constraints reflect economic reality.

1037 **Learned Rule from Wine Quality Dataset.** We examine a learned rule from the Wine Quality  
 1038 dataset that demonstrates CNF structure. The following Class-2 rule predicts high-quality wine and  
 1039 achieves 94.37% support, covering nearly all instances of this quality class:

1041  $(\text{total-SO}_2 < 3.729 \vee \text{free-SO}_2 < 3.490) \wedge$   
 1042  $(\text{alcohol} > 0.301 \wedge \text{alcohol} < 0.391) \wedge$   
 1043  $(\text{citric-acid} > 1.638 \wedge \text{citric-acid} < 1.700) \wedge$   
 1044  $(\text{density} > 0.066 \wedge \text{density} < 0.132) \wedge$   
 1045  $(\text{volatile-acidity} < 4.797) \wedge (\text{pH} > 0.120) \wedge$   
 1046  $(\text{total-SO}_2 > -0.094 \vee \text{free-SO}_2 > 0.938)$   
 1047

1048 This rule exhibits CNF structure where wine quality requires multiple chemical properties to simulta-  
 1049 neously satisfy constraints. Most clauses are unit literals establishing upper or lower bounds (e.g.,  
 1050  $1.638 < \text{citric-acid} < 1.700$  creates a narrow acceptable range of width 0.062), while two clauses  
 1051 are disjunctions offering alternative sulfur dioxide conditions. Unlike the Adult dataset where high  
 1052 income admits multiple alternative paths (DNF), wine quality demands that all conjunctive clauses be  
 1053 satisfied—reflecting the domain knowledge that wine quality depends on balanced chemistry rather  
 1054 than any single exceptional property. The exceptionally high support (94.37%) indicates this rule  
 1055 captures the essential characteristics defining this quality class.  
 1056

## 1057 B RELATED WORK

### 1058 B.1 TRADITIONAL RULE LEARNING METHODS

1059 Historically, example-based rule learning algorithms (S., 1969; Michalski, 1973; Pagallo & Haussler,  
 1060 1990) were initially proposed by selecting a random example and finding the best rule to cover it.  
 1061 However, due to their computational inefficiency, CN2 (Clark & Niblett, 1989) explicitly changed  
 1062 the strategy to finding the best rule that covers as many examples as possible. Building on these  
 1063 heuristic algorithms, FOIL (Quinlan, 1990), a system for learning first-order Horn clauses, was  
 1064 subsequently developed. While some algorithms learn rule sets directly, such as RIPPER (Cohen,  
 1065 1995), PART (Frank & Witten, 1998), and CPAR (Yin & Han, 2003), others post-process a decision  
 1066 tree (Quinlan, 2014) or construct sets of rules by post-processing association rules, like CBA (Liu  
 1067 et al., 1998) and CMAR (Li et al., 2001). All these algorithms use different strategies to find and use  
 1068 sets of rules for classification.  
 1069

1070 Most of these algorithms are based on Disjunctive Normal Form (DNF) (S., 1969; Michalski, 1973;  
 1071 Pagallo & Haussler, 1990; Clark & Niblett, 1989; Liu et al., 1998; Li et al., 2001; Frank & Witten,  
 1072 1998; Yin & Han, 2003; Cohen, 1995) expressions. CNF learners have been shown to perform  
 1073 competitively with DNF learners (Mooney, 1995), inspiring a line of CNF learning algorithms (Dries  
 1074 et al., 2009; Jain et al., 2021; Beck et al., 2023; Sverdlik, 1992; Hong & Tsang, 1997). Traditional  
 1075 rule-based models are valued for their interpretability but struggle to find the global optimum due to  
 1076 their discrete, non-differentiable nature. Extensive exploration of heuristic methods (Quinlan, 2014;  
 1077 Loh, 2011; Cohen, 1995) has not consistently yielded optimal solutions.  
 1078

1079 In response, recent research has turned to Bayesian frameworks to enhance model structure (Letham  
 et al., 2015; Wang et al., 2017; Yang et al., 2017), employing strategies such as if-then rules (Lakkaraju

et al., 2016) and advanced data structures for quicker training (Angelino et al., 2018). Despite these advancements, extended search times, scalability challenges, and performance issues limit the practicality of rule-based models compared to ensemble methods like Random Forest (Breiman, 2001) and Gradient Boosted Decision Trees (Chen & Guestrin, 2016; Ke et al., 2017), which trade off interpretability for improved performance.

## B.2 RULE LEARNING NEURAL NETWORKS

Neural rule learning-based methods integrate rule learning with advanced optimization techniques, enabling the discovery of complex and nuanced rules that combine the interpretability of symbolic models with the generalization power of neural networks. Unlike tree-based models, which explicitly follow feature-condition rules, neural approaches rely on weight parameters to control the rule learning process, offering improved robustness and scalability through data-driven training. However, existing approaches such as neural decision trees and rule extraction from neural networks (Frosst & Hinton, 2017; Ribeiro et al., 2016; Wang et al., 2020; 2021; Zhang et al., 2023) face challenges in fidelity and scalability. In particular, RRL (Wang et al., 2021; 2024), a state-of-the-art rule-based neural network, requires a predefined structure of CNF and DNF layers, which limits flexibility, leads to inefficient rule discovery, and exacerbates optimization issues such as gradient vanishing (Wang et al., 2020).

Recent work on differentiable logic gate networks (Petersen et al., 2022) learns distributions over all 16 binary Boolean operators with fixed sparse connectivity for fast hardware inference ( $> 10^6$  inferences/s). However, DLGN uses continuous relaxations during training and discretizes only at inference, whereas RNS maintains discrete operations throughout via Straight-Through Estimator. Moreover, DLGN lacks explicit rule extraction mechanisms, while RNS evaluates rule coverage, diversity, and recovery for interpretable CNF/DNF learning.

Our proposed Rule Network with Selective Logical Operators (RNS) addresses these challenges through its novel architecture and optimization strategies, improving scalability, rule quality, and training stability, as detailed in Section 3.

## B.3 BINARIZED NEURAL NETWORK

A related topic to this work is Binarized Neural Networks (BNNs), which optimize deep neural networks by employing binary weights. The deployment of deep neural networks typically requires substantial memory storage and computing resources. To achieve significant memory savings and energy efficiency during inference, recent efforts have focused on learning binary model weights while maintaining the performance levels of their floating-point counterparts Courbariaux et al. (2015; 2016a); Rastegari et al. (2016); Bulat & Tzimiropoulos (2019); Liu et al. (2018). Innovations such as bit logical operations Kim & Smaragdīs (2016) and novel training strategies for self-binarizing networks Lahoud et al. (2019), along with integrating scaling factors for weights and activations Sakr et al. (2018), have advanced BNNs. However, due to the binary nature of their weights, BNNs face optimization challenges. The Straight-Through Estimator (STE) method Courbariaux et al. (2015; 2016a); Cheng et al. (2019b) allows gradients to "pass through" non-differentiable functions, making it particularly effective for discrete optimization.

Despite both using binarized model weights and employing STE for optimization, our work diverges significantly from BNNs. First, RNS adopts specialized logical activation functions to perform logical operations on features, whereas BNNs typically use the Sign function to produce binary outputs. Second, BNNs are fully connected neural networks, while RNS features a learning mechanism for its connections. Most importantly, these distinctions enable RNS to learn logical rules for both prediction and interpretability, setting it apart from BNNs, which are primarily designed to enhance model efficiency.

## C DATASET STATISTICS

In Table 5, the first nine data sets are small, while the last four are large. The "Discrete" or "Continuous" feature type indicates that all features in the data set are either discrete or continuous,

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  
1175  
1176  
1177  
1178  
1179  
1180  
1181  
1182  
1183  
1184  
1185  
1186  
1187

Dataset	#instances	#classes	#features	feature type
adult	32561	2	14	mixed
bank	45211	2	16	mixed
banknote	1372	2	4	continuous
chess	28056	18	6	discrete
connect-4	67557	3	42	discrete
letRecog	20000	26	16	continuous
magic04	19020	2	10	continuous
tic-tac-toe	958	2	9	discrete
wine	178	3	13	continuous
activity	10299	6	561	continuous
dota2	102944	2	116	discrete
facebook	22470	4	4714	discrete
fashion	70000	10	784	continuous

Table 5: Datasets statistics.

respectively. The "Mixed" feature type indicates that the corresponding data set contains both discrete and continuous features.

## D REPRODUCIBILITY

**Data Preprocessing.** All datasets are converted into tabular feature vectors to ensure consistent treatment across different data modalities. For standard tabular benchmarks (Adult, Bank, Chess, Connect-4, etc.), we use the original attributes provided by the datasets. Categorical features are one-hot encoded, while continuous features are treated as real-valued inputs (as summarized in Table 5) and then binarized using the scheme described in Section 2.2. For Fashion-MNIST, each 28×28 grayscale image is flattened into a 784-dimensional continuous feature vector, and pixel values are passed through the same feature binarization layer as other continuous features. The binarization method uses random binning in main experiments, with ablations over EntInt/KInt/AutoInt in Appendix H. No additional feature engineering specific to images is applied; Fashion-MNIST is treated purely as a high-dimensional tabular dataset. All methods operate on the same base preprocessed features to ensure fair comparison: interpretable neuro-rule baselines (RRL, CRS) use the same binarized representation as in their original code following Wang et al. (Wang et al., 2021; 2024), while non-interpretable tabular models (FT-Transformer, SAINT, NODE, TabNet, LGBM, XGB) are trained on the same continuous/tabular features before binarization.

**Model Configuration.** RNS includes two Logic Selection Layers (LSLs), with the number of logical neurons grid-searched from 32 to 4096 based on dataset complexity. For training, we use cross-entropy loss with L2 regularization to control model complexity, with the regularization coefficient searched in the range  $10^{-5}$  to  $10^{-9}$ . The number of bins in the feature binarization layer is selected from {15, 30, 50}. The model is trained using the Adam optimizer (Kingma & Ba, 2014) with a batch size of 32. For small datasets, training runs for 400 epochs, with the learning rate reduced by 10% every 100 epochs. For large datasets, training is conducted for 100 epochs with a similar learning rate schedule, decreasing every 20 epochs. Baseline settings follow those described in (Wang et al., 2021; 2024). RNS is implemented in PyTorch (Paszke et al., 2019), and we use the high-quality code base from (Wang et al., 2021; 2024) for baseline comparisons. Experiments are performed on a Linux server equipped with an NVIDIA A100 80GB GPU. All code and data are available at [https://anonymous.4open.science/r/RNS\\_4A67/](https://anonymous.4open.science/r/RNS_4A67/).

## E HYPERPARAMETER STUDY

**LSL Dimension K.** We analyze the effect of layer dimension K (number of neurons) in the two LSLs. A larger K increases model complexity, potentially improving performance but risking overfitting. We test dimensions K in {32@32, 64@64, 128@128, 256@256, 512@512, 1024@1024, 2048@2048} on a small dataset (bank) and a large dataset (activity) shown in Figure 5. The F1 score initially rises and then falls as K increases, peaking at 1024@1024 for the activity dataset and 128@128 for the

bank dataset. This indicates that the activity dataset benefits from larger models, while the bank dataset performs best with smaller models, suggesting that optimal model size depends on dataset complexity.

**L2 Regularization  $\lambda$ .** We examine the impact of the L2 regularization weight  $\lambda$ , which controls RNS complexity. A larger  $\lambda$  reduces model complexity, resulting in fewer and simpler logical rules. We vary  $\lambda$  from  $10^{-5}$  to  $10^{-9}$  and show the model’s performance in Figure 5. The F1 score initially rises and then falls as  $\lambda$  increases, indicating that performance improves with an appropriately chosen  $\lambda$ . Balancing  $\lambda$  is essential to avoid overfitting or underfitting and achieve optimal model complexity.

## F EFFICIENCY

We evaluate learning efficiency by the number and length of learned rules across all datasets, as shown in Figure 9 and Figure 10, respectively. Figure 9 presents scatter plots of F1 score against  $\log(\#\text{edges})$  across 10 datasets. The boundary connecting the results of different RNS architectures separates the upper left corner from the best baseline methods. This indicates that RNS consistently learns fewer rules while achieving high prediction accuracy across various scenarios. The average length of rules in RNS trained on different datasets is shown in Figure 10. The average rule length is less than 6 for all datasets except fashion and facebook, which are unstructured datasets and have more complex features. These results indicate that the rules learned by RNS are generally easy to understand across different scenarios.

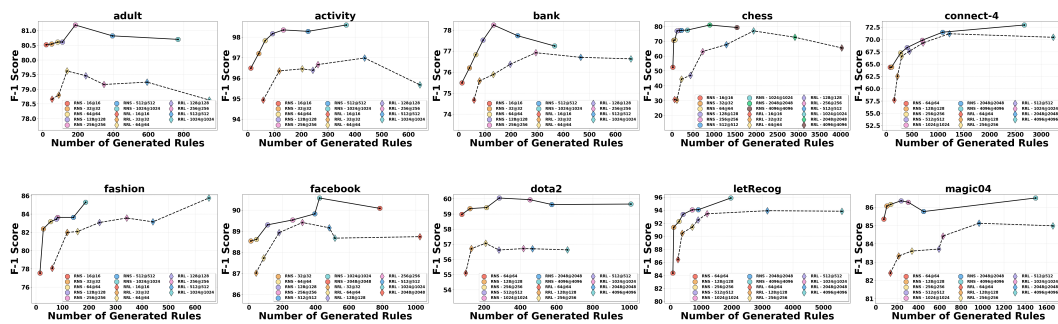


Figure 9: Rule Number comparison across different architectures.

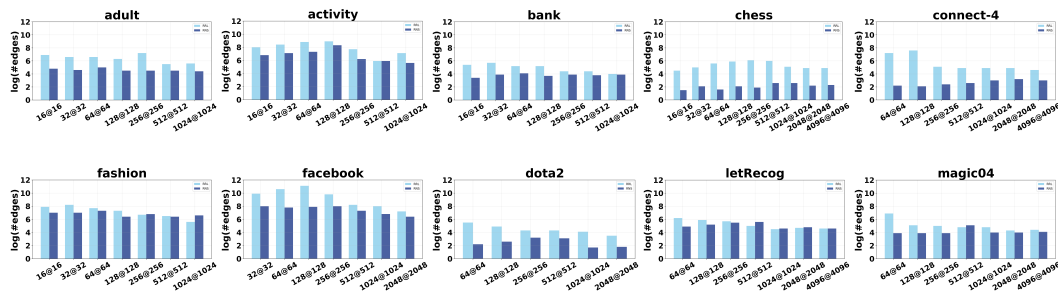


Figure 10: Rule Length comparison across different architectures.

## G RULE QUALITY

**Rule Quality Metrics.** Following prior work Yu et al. (2023); Lakkaraju et al. (2016), we evaluate rule quality using three key metrics: diversity, coverage, and single-rule accuracy. Let  $I_r$  denote the set of data instances covered by rule  $r$ , and  $D$  denote the complete dataset. The metrics are defined as:

**Single-rule Accuracy** measures the prediction accuracy of a single rule for the instances it covers:

$$\text{Single-rule Accuracy} = \frac{|\{i \in I_r : \text{rule prediction matches } y_i\}|}{|I_r|} \quad (7)$$

This metric evaluates how accurately a single rule classifies the data instances it covers when used independently for prediction. It measures the proportion of correctly classified instances among all instances that satisfy the rule’s conditions. Higher single-rule accuracy indicates that the rule is more reliable for making predictions on its own, without relying on other rules in the rule set. **Coverage** quantifies the proportion of data instances covered by a rule:

$$\text{Coverage} = \frac{|I_r|}{|D|} \quad (8)$$

This metric measures the scope or generality of a rule. Lower coverage indicates that rules are more specific and easier for human experts to understand, as they apply to a smaller, more focused subset of the data. Very high coverage may suggest overly general rules that lack specificity.

**Diversity** measures the overlap ratio between pairs of rules, with higher diversity reflecting that rules capture distinct, non-redundant logic:

$$\text{Diversity} = 1 - \frac{|I_i \cap I_j|}{|I_i \cup I_j|} \quad (9)$$

This metric quantifies how different rules are from each other by measuring their overlap. Higher diversity values indicate that rules capture different patterns in the data with minimal redundancy, leading to a more comprehensive and non-redundant rule set. Low diversity suggests that multiple rules are covering similar data instances, which reduces the interpretability and efficiency of the rule set.

**Results.** We conduct extensive experiments across 10 datasets, as shown in Figure 11, Figure 12, and Figure 13. RNS consistently produces rules with superior results: higher accuracy, greater diversity, and lower coverage deviation. This indicates that the rules learned by RNS are easier to distinguish and exhibit better prediction and generalization power.

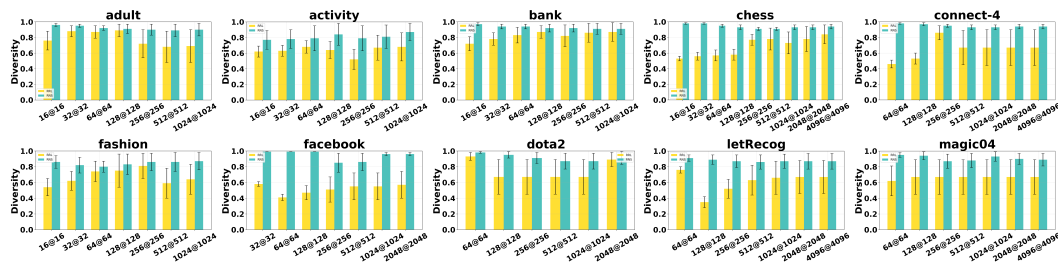


Figure 11: Rule Diversity comparison across different architectures.

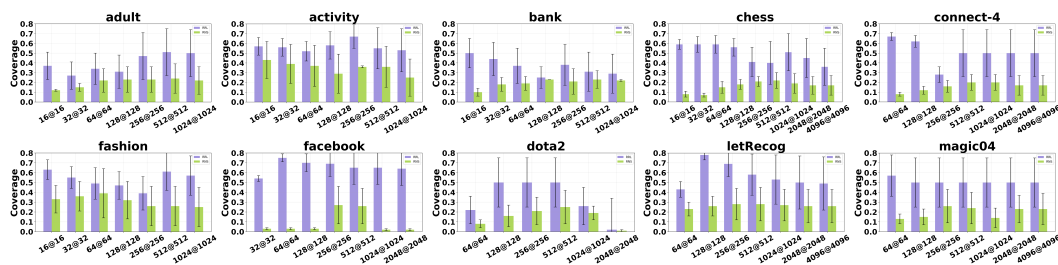


Figure 12: Rule Coverage comparison across different architectures.

## H SIMULATION EXPERIMENT

We conduct a controlled simulation experiment to assess RNS’s ability to learn the exact logical rules used to generate synthetic data.

1296  
1297  
1298  
1299  
1300  
1301  
1302  
1303  
1304  
1305  
1306  
1307  
1308  
1309  
1310  
1311  
1312  
1313  
1314  
1315  
1316  
1317  
1318  
1319  
1320  
1321  
1322  
1323  
1324  
1325  
1326  
1327  
1328  
1329  
1330  
1331  
1332  
1333  
1334  
1335  
1336  
1337  
1338  
1339  
1340  
1341  
1342  
1343  
1344  
1345  
1346  
1347  
1348  
1349

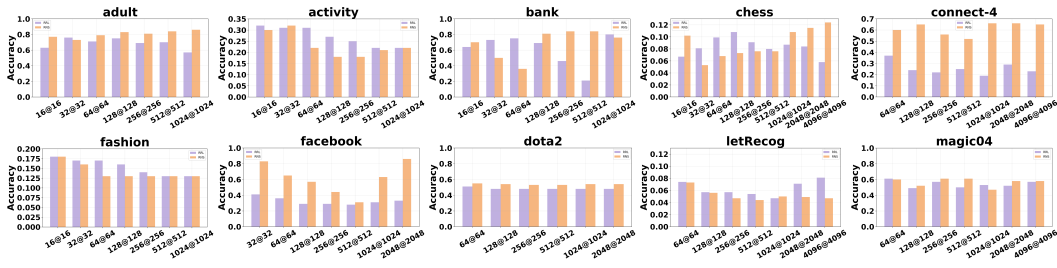


Figure 13: Rule Accuracy comparison across different architectures.

**Setup.** We generate synthetic data based on predefined logical rules and train rule-based models to evaluate their ability to reconstruct these ground-truth rules. The dataset is synthesized by defining three probability parameters,  $\mathbf{p} = (p_1, p_2, p_3)$ , corresponding to feature variables  $\{x_1, x_2, x_3\}$ , each drawn independently from a uniform distribution  $U(0, 1)$ . Using these as Bernoulli parameters, we generate  $X_{gen} \in \mathbb{R}^{3 \times 50000}$  by repeating Bernoulli sampling, resulting in  $n = 50000$  binary vectors of length 3, where each element is sampled from  $Bernoulli(p_i)$ . Labels are assigned to these binary vectors based on specific logical rules. For example, a rule  $x_1 \wedge x_2 \wedge x_3 \rightarrow 1$  assigns a label of 1 if all three features are 1. We define four different types of rules, presented as "Ground-Truth Rules" in Table 2. The dataset is split evenly into training and test sets. Both RNS and RRL are configured with logical layer dimensions of 64@64.

**Results.** As shown in Table 2, RNS not only achieves near 100% accuracy but also recovers the exact structure of the ground-truth rules. In contrast, RRL, even with a naive negation setting, tends to produce overly simplified, redundant, or logically incomplete rules.

For example, RNS is able to precisely recover rules such as  $x_1 \wedge \neg x_2 \wedge x_3$  and  $x_2 \wedge \neg x_3$ , faithfully mirroring the underlying logic used to generate the data. RRL, on the other hand, often produces rules that are either too general (e.g.,  $x_1, x_2$ ) or combine terms in a way that does not fully reflect the intended logic (e.g.,  $(x_1 \wedge x_3), x_2$ ). This pattern holds even when RRL is augmented with the naive negation setting: RNS consistently recovers all ground-truth rules, whereas RRL fails to identify several rules in their correct logical form and often outputs multiple trivial or repeated variants for a single ground-truth pattern.

## I ABLATION STUDY (RQ4)

**Negation Layer and NFC.** To evaluate the effectiveness of the Negation Layer and Normal Form Constraint (NFC), we compare RNS trained with and without these components. The Negation Layer enables functional completeness by supporting negation operations in learned rules, while NFC narrows the search space for learning connections between logical layers, improving performance. Table 6 shows F1 scores for RNS and its variants on the bank (small) and activity (large) datasets. The F1 score decreases without these components, highlighting their importance. NFC also improves learning efficiency. With all other factors constant, training time decreases by 11.22% and 25% for the activity and bank datasets, respectively, as shown in Table 7, emphasizing its utility.

**Binning Function.** We evaluate RNS’s performance with different binning functions introduced in Section 2.2. F1 scores for RNS with these methods are shown in Table 8. RNS exhibits strong flexibility, with RanInt performing the best due to its stochastic diversity, which helps prevent overfitting. Additionally, its computational efficiency and simplicity are advantageous. KInt, in contrast, clusters based solely on feature values, neglecting target variables, which can reduce its effectiveness. EntInt incorporates label information to minimize entropy within bins, potentially improving accuracy. AutoInt, while adaptable, incurs significant computational overhead due to its parameter optimization, posing challenges in practice.

Model Variant	Activity	Bank
RNS w/o Negation Layer	97.68	76.92
RNS w/o NFC	97.81	76.36
<b>RNS (Full)</b>	<b>98.37</b>	<b>77.18</b>

Table 6: Ablation study: F1 Score (%) of RNS and its variants on two datasets.

Configuration	Activity	Bank
RNS w/o NFC	1h 38m 9s	1h 24m 50s
<b>RNS (with NFC)</b>	<b>1h 27m 3s</b>	<b>1h 3m 56s</b>

Table 7: Training efficiency: Impact of Normal Form Constraint (NFC) on training time.

Binning Method	Activity	Bank
KInt (K-means)	92.71	74.37
EntInt (Entropy)	97.09	77.02
AutoInt (Auto-interval)	95.62	75.59
<b>RanInt (Random)</b>	<b>98.80</b>	<b>77.18</b>

Table 8: Feature preprocessing: F1 Score (%) comparison of different binning methods.

## J MODEL INTERPRETATION

Figure 14 presents the three most discriminative rules extracted from RNS, trained on the bank-marketing dataset to identify customer profiles likely to subscribe to a term deposit through telesales. These rules provide detailed profiles by highlighting specific conditions that increase subscription likelihood. ‘Coverage’ denotes the proportion of training samples satisfying a rule.

The transparent interpretations reveal several insights. Subscriptions are less likely during colder seasons, possibly due to reluctance to make financial decisions during holidays. The model identifies a correlation between longer call durations and successful marketing, suggesting that extended dialogues often indicate interest in deposits. It also highlights that housemaids, entrepreneurs, and self-employed individuals are less inclined towards term deposits, potentially due to financial constraints or the need for liquidity. Additionally, the analysis shows age-related trends in call duration, with younger individuals having shorter calls and middle-aged clients engaging in longer discussions. These findings demonstrate the model’s ability to reflect plausible real-world behaviors while avoiding direct causal assertions.

## K GRADIENT VANISHING

Despite the high interpretability of discrete logical layers, training them is challenging due to their discrete parameters and non-differentiable structures. This challenge is addressed by drawing inspiration from the training methods used in binary neural networks, which involve searching for discrete solutions within a continuous space. Wang et al. (2021) leverages the logical activation functions proposed by Payani & Fekri (2019) in RRL:

$$AND(h, W_{conn}^i) = \prod_{j=1}^n F_c(h_j, W_{i,j}) \quad (10)$$

$$OR(h, W_{conn}^i) = 1 - \prod_{j=1}^n (1 - F_d(h_j, W_{i,j})) \quad (11)$$

where  $F_c(h, w) = 1 - w(1 - h)$  and  $F_d(h, w) = h \cdot w$ .

Rules	Contribution to Negative Class	Contribution to Positive Class	Coverage
(duration > 146.719 $\wedge$ $\neg$ job $\in$ {housemaid, entrepreneur} $\wedge$ $\neg$ month $\in$ {nov, mar, feb} $\wedge$ balance < 1757.276) $\vee$ (duration < 850.975 $\wedge$ $\neg$ job = technician $\wedge$ month = mar) $\vee$ (duration > 146.719 $\wedge$ $\neg$ month = jan $\wedge$ $\neg$ contact = unknown) $\vee$ (duration > 19.057 $\wedge$ $\neg$ age > 67.87 $\wedge$ month = oct $\wedge$ $\neg$ day < 5.569) $\vee$ ( $\neg$ duration < 850.975 $\wedge$ age < 66.97 $\wedge$ $\neg$ job = housemaid $\wedge$ $\neg$ month = jan $\wedge$ $\neg$ payout = success) $\vee$ (duration > 59.781 $\wedge$ $\neg$ job = unknown) $\vee$ (duration > 19.057 $\wedge$ age > 28.949) $\vee$ (duration > 59.781 $\wedge$ $\neg$ age > 67.87 $\wedge$ $\neg$ job $\in$ {self-employed, housemaid} $\wedge$ $\neg$ month $\in$ {oct, nov} $\wedge$ $\neg$ loan = yes $\wedge$ $\neg$ default = yes)	-0.4159	0.3929	0.4236
(duration > 59.781 $\wedge$ $\neg$ month $\in$ {jul, aug} $\wedge$ balance < 1537.899 $\wedge$ previous > 1.028 $\wedge$ pdays < 229.326) $\vee$ ( $\neg$ month = mar $\wedge$ $\neg$ payout = success $\wedge$ 102.083 < pdays < 304.124 $\wedge$ previous < 6.754)	0.3314	-0.3539	0.3617
( $\neg$ month = aug $\wedge$ $\neg$ age < 12.397 $\wedge$ campaign > 0.066 $\wedge$ day < 33.738 $\wedge$ pdays < 244.292)	-0.5581	0.5437	0.3319

Figure 14: Logical rules from RNS on the bank-marketing dataset.

If  $\mathbf{h}$  and  $W_i$  are both binary vectors, then  $\text{Conj}(\mathbf{h}, W_i) = \bigwedge_{W_{i,j}=1} h_j$  and  $\text{Disj}(\mathbf{h}, W_i) = \bigvee_{W_{i,j}=1} h_j$ .  $F_c(\mathbf{h}, W)$  and  $F_d(\mathbf{h}, W)$  decide how much  $h_j$  would affect the operation according to  $W_{i,j}$ . After using continuous weights and logical activation functions, the AND and OR operators, denoted by  $R$  and  $S$ , are defined as follows:

$$r_i^{(l)} = \text{AND}(u^{(l-1)}, W_i^{(l,0)}) \quad (12)$$

$$s_i^{(l)} = \text{OR}(u^{(l-1)}, W_i^{(l,1)}) \quad (13)$$

Although the whole logical layer becomes differentiable by applying this continuous relaxation, the above functions are not compatible with RNS. In this setting, the output of the node  $h \in [0, 1]$  conflicts with our binarized setting, where all logical parameters are either  $-1$  or  $+1$ . This not only breaks the inherently discrete nature of RNS but also suffers from the serious vanishing gradient problem. The reason can be found by analyzing the partial derivative of each node with respect to its directly connected weights and with respect to its directly connected nodes as follows:

$$\frac{\partial r}{\partial W_{i,j}} = (u_j^{(l-1)} - 1) \cdot \prod_{k \neq j} F_c(u_k^{(l-1)}, W_{i,k}^{(l,0)}) \quad (14)$$

$$\frac{\partial r}{\partial u_j^{(l-1)}} = W_{i,j}^{(l,0)} \cdot \prod_{k \neq j} F_c(u_k^{(l-1)}, W_{i,k}^{(l,0)}) \quad (15)$$

Since  $u_j^{(l-1)}$  and  $W_{i,k}^{(l,0)}$  fall within the range of 0 to 1, the values of  $F_c(\cdot)$  also lie within this range. If the number of inputs is large and most  $F_c(\cdot)$  are not 1, the gradient tends toward 0 due to multiplications. Additionally, in the discrete setting, only when  $u_j^{(l-1)}$  and  $1 - F_d(\cdot)$  are all 1, can the gradient be non-zero, which is quite rare in practice.

## K.1 DETAILED ANALYSIS

Training discrete logical layers in neural networks is challenging due to their binary parameters and non-differentiable operations. A common approach to make such models trainable is to relax the logical functions to a continuous, differentiable form Wang et al. (2024). These relaxations enable gradient-based training, but they often suffer from severe **vanishing gradient** problems due to the multiplicative structure of the logical functions. Moreover, the RRL activations output values in  $[0, 1]$ , which conflicts with a fully binarized  $\pm 1$  logic setting (such as our RNS approach) and breaks the discrete semantics of the network. We detail below why vanishing gradients occur in such formulations and describe the solutions proposed in RRL as well as the different strategy taken by RNS to overcome these issues.

## K.2 VANISHING GRADIENTS IN MULTIPLICATIVE LOGICAL ACTIVATIONS (RRL)

In a conventional formulation of a logical AND gate Payani & Fekri (2019), the output is the product of its binary inputs. For binary  $x_i \in \{0, 1\}$ , one can write

$$\text{AND}(x_1, \dots, x_n) = \prod_{i=1}^n x_i, \quad (16)$$

and similarly, a logical OR can be written (using De Morgan’s law) as

$$\text{OR}(x_1, \dots, x_n) = 1 - \prod_{i=1}^n (1 - x_i). \quad (17)$$

While these expressions are correct for binary values, their direct use in a neural network leads to zero gradients almost everywhere – a phenomenon known as *gradient vanishing*. The gradient of the AND with respect to one input is

$$\frac{\partial \text{AND}(x)}{\partial x_j} = \prod_{i \neq j} x_i, \quad (18)$$

which is zero whenever any other input  $x_i$  is zero. More generally, if  $x_i$  are in  $[0, 1]$ , this derivative equals the product of the other  $(n - 1)$  inputs. As  $n$  grows or if the inputs are often fractional values  $< 1$ , this product becomes extremely small – decaying exponentially with  $n$  – thus severely diminishing the gradient signal. In the case of OR (product of  $(1 - x_i)$  terms), the same issue arises by symmetry (the gradient will contain a product of many factors in  $[0, 1]$ ). The consequence is that standard product-based logic units have large regions of the input space where the gradient is essentially zero, hampering learning.

In a fully discrete setting, if a conjunction has more than one false literal, changing any single input from 0 to 1 does not alter the output (since at least one false remains to make the AND false). Thus, the gradient at that point is exactly zero in all those input directions – creating broad “dead zones” where the network cannot learn.

To enable gradient-based training, RRL replaces the hard logical operations with differentiable approximations that produce continuous outputs. Let  $h = (h_1, \dots, h_n)$  be the input vector to a logical neuron (with  $h_j \in [0, 1]$  as relaxations of Boolean values) and let  $W_i = (W_{i,1}, \dots, W_{i,n})$  be a set of weights indicating which inputs are involved in the  $i$ -th clause (for example,  $W_{i,j} \in \{0, 1\}$  or  $[0, 1]$ , with 1 meaning the  $j$ th input is included in the clause). RRL defines smoothed conjunction and disjunction functions as:

$$\text{Conj}(h, W_i) = \prod_{j=1}^n F_c(h_j, W_{i,j}), \quad F_c(h, w) = 1 - w(1 - h), \quad (19)$$

$$\text{Disj}(h, W_i) = 1 - \prod_{j=1}^n (1 - F_d(h_j, W_{i,j})), \quad F_d(h, w) = h \cdot w, \quad (20)$$

where  $F_c(h, w)$  and  $F_d(h, w)$  are specific smooth blending functions for inputs  $h$  and weights  $w$ . If  $h \in \{0, 1\}^n$  and  $W_i$  is binary,  $\text{Conj}(h, W_i)$  reduces to the logical AND of the selected inputs (those with  $W_{i,j} = 1$ ), and  $\text{Disj}(h, W_i)$  becomes the logical OR. For continuous  $h$  and  $W_i$ , these give differentiable outputs in  $[0, 1]$ .

However, the gradient remains problematic. Consider  $r_i = \text{Conj}(h, W_i)$ . Using the chain rule:

$$\frac{\partial r_i}{\partial W_{i,j}} = \left( \prod_{k \neq j} F_c(h_k, W_{i,k}) \right) \frac{\partial F_c(h_j, W_{i,j})}{\partial W_{i,j}} = (h_j - 1) \prod_{k \neq j} F_c(h_k, W_{i,k}), \quad (21)$$

$$\frac{\partial r_i}{\partial h_j} = \left( \prod_{k \neq j} F_c(h_k, W_{i,k}) \right) \frac{\partial F_c(h_j, W_{i,j})}{\partial h_j} = W_{i,j} \prod_{k \neq j} F_c(h_k, W_{i,k}). \quad (22)$$

Each partial derivative is proportional to a product of  $(n - 1)$  terms  $F_c(h_k, W_{i,k})$ , each of which lies in  $[0, 1]$ . Unless all those terms are very close to 1, the product will be small; if any term is 0, the product (and thus the gradient) is zero. In other words, the magnitude of the gradient *shrinks multiplicatively* with every additional input in the clause.

A similar calculation for  $s_i = \text{Disj}(h, W_i)$  yields, by symmetry (with  $G_k := 1 - F_d(h_k, W_{i,k})$ ):

$$\frac{\partial s_i}{\partial W_{i,j}} = h_j \prod_{k \neq j} G_k, \quad \frac{\partial s_i}{\partial h_j} = W_{i,j} \prod_{k \neq j} G_k, \quad (23)$$

1512 which again contains a product of  $(n - 1)$  factors  $G_k \in [0, 1]$ . Thus, for large  $n$  or for typical  
 1513 fractional values of  $h_j$  and  $W_{i,j}$ , these gradients **vanish** to near-zero. In practice, when many inputs  
 1514 of an AND are not extremely close to 1, the gradient through that AND node becomes negligibly  
 1515 small. And in the discrete limit ( $h, W \in \{0, 1\}$ ),  $\frac{\partial r_i}{\partial W_{i,j}}$  and  $\frac{\partial r_i}{\partial h_j}$  are zero in almost all cases: only if  
 1516 *all* other inputs of the AND are 1 (so that the output is sensitive to the remaining input/literal) will  
 1517 a change in that input make a difference. Such a situation—e.g., a clause where all but one literal  
 1518 is true—is rare, meaning the network spends most of its time in regimes where the loss gradient is  
 1519 zero with respect to each individual parameter. This underscores the fundamental incompatibility of  
 1520 naive multiplicative logical units with gradient-based learning: the more literals in a rule, the more  
 1521 pronounced the vanishing gradient problem becomes.

1522  
 1523 **Incompatibility with  $\pm 1$  Binarization (RNS).** An additional issue is that the RRL activations  
 1524  $\text{Conj}(h, W)$  and  $\text{Disj}(h, W)$  produce outputs in the range  $[0, 1]$  (since they are essentially probabilities  
 1525 or fractional truth values). This is incompatible with the design of RNS, where all internal logical  
 1526 signals are meant to be binary  $\{-1, +1\}$  values at runtime. In RNS, we require a hard True or  
 1527 False, encoded as  $+1$  or  $-1$ . Using RRL’s continuous outputs inside an RNS network would break  
 1528 the discrete semantics and require additional mechanisms to re-binarize the outputs at each layer.  
 1529 Moreover, as discussed above, those continuous outputs suffer from vanishing gradients when used  
 1530 in deep clauses. Thus, while RRL’s relaxation makes a logical layer differentiable, it does so at the  
 1531 cost of deviating from binary  $\pm 1$  representation and encountering extremely small gradients for large  
 1532 rules. These drawbacks motivate alternative approaches that maintain binary representations and  
 1533 avoid multiplicative shrinkage.

### 1534 K.3 RRL’S MITIGATIONS: LOG-DOMAIN SMOOTHING AND GRADIENT GRAFTING

1535  
 1536 RRL and related frameworks have proposed a couple of techniques to alleviate the vanishing gradient  
 1537 and training difficulties while still using product-based logic. We briefly summarize two such  
 1538 strategies: a log-domain scaled activation function, and a hybrid training method known as gradient  
 1539 grafting.

1540  
 1541 **Log-space smoothing of logical activations.** One idea introduced with RRL is to modify the prod-  
 1542 uct formulation by amplifying small values in the product through a log transformation. Specifically,  
 1543 define a projection function  $P(v)$  that boosts a small product  $v$ :

$$1544 P(v) = \frac{1}{1 - \log v}, \quad \frac{dP}{dv} = \frac{P(v)^2}{v}. \quad (24)$$

1545  
 1546 Here  $v > 0$  is a value (in practice,  $v$  will be a small positive number representing the product of  
 1547 several terms). The function  $P(v)$  is chosen such that when  $v$  is small,  $-\log v$  is large, so  $P(v)$  will  
 1548 significantly exceed  $v$  (for example, if  $v = 10^{-3}$ , then  $P(v) \approx 1/(1 - (-6.9)) \approx 1/7.9 \approx 0.127$ ,  
 1549 whereas  $v$  itself is 0.001). In this way,  $P(v)$  stretches the lower end of the range upward.

1550  
 1551 RRL incorporates this into the logical units by first computing the product of inputs in log-space and  
 1552 then projecting back. Using a small  $\varepsilon > 0$  to avoid  $\log 0$ , the *improved* conjunction and disjunction  
 1553 are defined as:

$$1554 \text{Conj}^+(h, W_i) = P\left(\prod_{j=1}^n (F_c(h_j, W_{i,j}) + \varepsilon)\right), \quad (25)$$

$$1555 \text{Disj}^+(h, W_i) = 1 - P\left(\prod_{j=1}^n (1 - F_d(h_j, W_{i,j}) + \varepsilon)\right). \quad (26)$$

1556  
 1557 When  $\varepsilon \rightarrow 0$  and  $h, W$  are binary,  $\text{Conj}^+$  and  $\text{Disj}^+$  recover the exact AND/OR as before. However,  
 1558 for intermediate values, these use  $P(\cdot)$  to prevent the product from becoming too small. If we let

$$1559 v = \prod_{j=1}^n (F_c(h_j, W_{i,j}) + \varepsilon)$$

1566 be the raw product inside  $P$ , then by the chain rule the gradient of  $\text{Conj}^+$  with respect to a parameter  
 1567 becomes:

$$1569 \frac{\partial \text{Conj}^+}{\partial W_{i,j}} = \frac{P(v)^2}{v} \left( \prod_{k \neq j} (F_c(h_k, W_{i,k}) + \varepsilon) \right) \frac{\partial F_c(h_j, W_{i,j})}{\partial W_{i,j}}, \quad (27)$$

$$1572 \frac{\partial \text{Conj}^+}{\partial h_j} = \frac{P(v)^2}{v} \left( \prod_{k \neq j} (F_c(h_k, W_{i,k}) + \varepsilon) \right) \frac{\partial F_c(h_j, W_{i,j})}{\partial h_j}. \quad (28)$$

1576 Comparing these to the original gradients equation 21–equation 22, we see that the pure product  
 1577 term  $\prod_{k \neq j} F_c(h_k, W_{i,k})$  is now multiplied by  $\frac{P(v)^2}{v}$ . For moderately small  $v$ , this factor can be  
 1578 significantly larger than 1, thereby attenuating the vanishing effect. Intuitively, instead of the gradient  
 1579 shrinking proportional to  $v$  (the product of many small terms), it shrinks proportional to  $P(v)^2$ , and  
 1580 since  $P(v) > v$  when  $v$  is small, the decay is slower. This log-space trick can appreciably increase  
 1581 gradient magnitudes when the clause length  $n$  is not too large or the inputs are not too extreme.

1582 Comparing these to the original gradients equation 21–equation 22, we see that the pure product  
 1583 term  $\prod_{k \neq j} F_c(h_k, W_{i,k})$  is now multiplied by  $\frac{P(v)^2}{v}$ . For moderately small  $v$ , this factor can be  
 1584 significantly larger than 1, thereby attenuating the vanishing effect.

1586 To understand how the gradient shrinks, let’s analyze this mathematically. Consider the product term:

$$1588 v = \prod_{j=1}^n F_c(h_j, W_{i,j}) = \prod_{j=1}^n (1 - W_{i,j}(1 - h_j)) \quad (29)$$

1591 For typical intermediate values where  $h_j \approx 0.5$  and  $W_{i,j} \approx 0.5$ , we have:

$$1593 F_c(h_j, W_{i,j}) \approx 1 - 0.5(1 - 0.5) = 0.75 \quad (30)$$

1595 Therefore, the product becomes:

$$1597 v \approx (0.75)^n \quad (31)$$

1598 As  $n$  increases, this decays exponentially:

$$1600 n = 10 : v \approx 0.75^{10} \approx 0.056 \quad (32)$$

$$1602 n = 20 : v \approx 0.75^{20} \approx 0.003 \quad (33)$$

$$1603 n = 50 : v \approx 0.75^{50} \approx 5.7 \times 10^{-7} \quad (34)$$

$$1604 n = 100 : v \approx 0.75^{100} \approx 3.2 \times 10^{-13} \quad (35)$$

1606 The gradient magnitude without log-smoothing is proportional to  $v$ :

$$1608 \left| \frac{\partial r_i}{\partial W_{i,j}} \right| \propto v \approx \alpha^n \quad \text{where } \alpha < 1 \quad (36)$$

1611 With the log-space projection  $P(v)$ , the gradient becomes:

$$1613 \left| \frac{\partial \text{Conj}^+}{\partial W_{i,j}} \right| \propto \frac{P(v)^2}{v} = \frac{1}{v(1 - \log v)^2} \quad (37)$$

1616 For small  $v$ , we have  $P(v) \approx \frac{1}{-\log v}$ , so:

$$1618 \frac{P(v)^2}{v} \approx \frac{1}{v(\log v)^2} \quad (38)$$

Let’s evaluate this for different clause sizes:

$$n = 10 : \quad v \approx 0.056, \quad \frac{P(v)^2}{v} \approx \frac{1}{0.056 \times (-2.88)^2} \approx 2.15 \quad (39)$$

$$n = 20 : \quad v \approx 0.003, \quad \frac{P(v)^2}{v} \approx \frac{1}{0.003 \times (-5.81)^2} \approx 9.87 \quad (40)$$

$$n = 50 : \quad v \approx 5.7 \times 10^{-7}, \quad \frac{P(v)^2}{v} \approx \frac{1}{5.7 \times 10^{-7} \times (-14.4)^2} \approx 8.5 \times 10^3 \quad (41)$$

While  $P(v)^2/v$  grows as  $v$  shrinks, for very small  $v$  (large  $n$ ), the growth is only polynomial in  $\log(1/v) \approx n \log(1/\alpha)$ , not enough to fully compensate for the exponential decay. Eventually, for very large  $n$ :

$$\lim_{n \rightarrow \infty} \frac{P(v)^2}{v} \approx \frac{1}{\alpha^n \cdot n^2 \cdot (\log \alpha)^2} \rightarrow 0 \quad (42)$$

While  $P(v)^2/v$  grows as  $v$  shrinks, for very small  $v$  (large  $n$ ), the growth is only polynomial in  $n$ , not enough to fully compensate for the exponential decay.

To see this clearly, recall that  $v \approx \alpha^n$  where  $\alpha < 1$  (e.g.,  $\alpha = 0.75$ ). Then:

$$P(v) = \frac{1}{1 - \log v} = \frac{1}{1 - \log(\alpha^n)} = \frac{1}{1 - n \log \alpha} \quad (43)$$

Since  $\alpha < 1$ , we have  $\log \alpha < 0$ , so  $-\log \alpha > 0$ . For large  $n$ :

$$P(v) \approx \frac{1}{n|\log \alpha|} \quad (44)$$

Therefore:

$$\frac{P(v)^2}{v} \approx \frac{1/n^2(\log \alpha)^2}{\alpha^n} = \frac{1}{\alpha^n \cdot n^2 \cdot (\log \alpha)^2} \quad (45)$$

The key observation is:

- The numerator grows as  $\mathcal{O}(1/n^2)$  (polynomial decay)
- The denominator decays as  $\mathcal{O}(\alpha^n)$  (exponential decay)

Since exponential decay dominates polynomial growth:

$$\lim_{n \rightarrow \infty} \frac{1}{\alpha^n \cdot n^2 \cdot (\log \alpha)^2} = \lim_{n \rightarrow \infty} \frac{1}{n^2(\log \alpha)^2} \cdot \frac{1}{\alpha^n} \rightarrow 0 \quad (46)$$

because  $\frac{1}{\alpha^n}$  approaches 0 much faster than  $\frac{1}{n^2}$  approaches infinity.

Thus, while the log-space trick delays the vanishing, it cannot prevent it for large  $n$ . Intuitively, instead of the gradient shrinking proportional to  $v$  (the product of many small terms), it shrinks proportional to  $P(v)^2$ , and since  $P(v) > v$  when  $v$  is small, the decay is slower. This log-space trick can appreciably increase gradient magnitudes when the clause length  $n$  is not too large or the inputs are not too extreme.

However, this modification does **not fully eliminate** the vanishing gradient problem. When  $n$  is very large or many inputs are significantly below 1, the initial product  $v$  becomes extremely tiny (e.g.  $10^{-10}$  or smaller), making  $\log v$  a large negative number. In such extreme cases,  $P(v)$  itself approaches 0 (since  $1 - \log v$  is huge), and consequently  $P(v)^2/v$  can also become very small. In the worst case, if any factor in the product is 0,  $v = 0$  and no finite smoothing can help ( $P(0)$  is undefined without  $\varepsilon$  and effectively  $P(v)^2/v$  remains near 0 for very small  $v$ ). Thus, for very complex clauses or highly non-saturated inputs, the gradients may still collapse to nearly zero. Moreover, once we ultimately project the network to discrete weights and inputs ( $h, W \in \{0, 1\}$ ), the gradient at those exact binary points is again zero in most directions (as discussed earlier). In summary,  $P(v)$ -based smoothing improves gradient flow for moderately small signals but does not fundamentally overcome vanishing gradients when training very large logical expressions. Additional strategies are required to train purely discrete models.

**Gradient Grafting for discrete training.** Another technique used in RRL to handle training with binary decisions is **Gradient Grafting**. The idea is to maintain two parallel models during training: a continuous one that is used for backpropagation, and a discrete one that defines the actual loss. Let  $\theta$  denote the set of trainable continuous parameters (weights) and  $q(\theta)$  be a binarization function that maps these to discrete values (for instance, thresholding each weight to 0 or 1). We denote by  $\hat{Y} = F(\theta, X)$  the output of the continuous model on input  $X$ , and by  $\bar{Y} = F(q(\theta), X)$  the output of the corresponding binarized model (i.e., the actual logical network with hard decisions). Training proceeds by computing the loss  $L(\bar{Y})$  on the discrete model’s output, but then updating  $\theta$  using gradients from the continuous model. In formula:

$$\theta_{t+1} = \theta_t - \eta \left[ \frac{\partial L(\bar{Y})}{\partial \bar{Y}} \right] \left[ \frac{\partial \hat{Y}}{\partial \theta_t} \right], \quad (47)$$

where  $\eta$  is the learning rate. In this update,  $\frac{\partial L(\bar{Y})}{\partial \bar{Y}}$  is the gradient of the loss with respect to the discrete model’s output (this measures how the final loss would change if the discrete output changed), and  $\frac{\partial \hat{Y}}{\partial \theta_t}$  is the Jacobian of the continuous model’s output with respect to its parameters (which is well-defined and non-zero because  $\hat{Y}$  is produced by smooth activations). The chain of these two terms provides an effective surrogate gradient for  $\theta$  that steers the discrete model’s loss  $L(\bar{Y})$  downward, even though  $\bar{Y}$  itself has zero or undefined gradients w.r.t.  $\theta$ . In essence, the discrete network’s error signal is “grafted” onto the continuous network’s sensitivity.

This approach circumvents the vanishing gradient at the discrete points by always following the continuous proxy’s gradients, and can successfully optimize a logical network where direct backpropagation would fail. Gradient grafting, however, comes at the cost of increased training complexity. The optimization is no longer a simple gradient descent on a single well-defined objective; instead, it couples two models (one binary, one continuous) and relies on their interplay. The update in Eq. equation 47 is not the true gradient of any single loss function with respect to  $\theta$ , since  $L(\bar{Y})$  is evaluated on a different forward path than  $\hat{Y}$ . Thus, careful tuning and heuristics may be needed to make this training scheme converge reliably. Nonetheless, this method has been shown to help train discrete logical networks that would otherwise be stuck due to vanishing gradients.

In summary, RRL’s approach to training discrete logical rules involves smoothing the logical operations (to keep gradients alive) and using a hybrid training procedure to inject discrete loss information, partially mitigating the vanishing gradient issue.

#### K.4 RNS: $\pm 1$ ENCODING AND min/max LOGICAL ACTIVATIONS

RNS takes a fundamentally different route to avoid vanishing gradients: it redesigns the logical neuron computations and encoding so that gradients do not collapse in the first place, even at discrete points. There are two key aspects to this strategy:

**(1)  $\pm 1$  Binary Encoding (No Zeros).** Instead of representing False as 0 and True as 1, RNS encodes Boolean values as  $-1$  (false) and  $+1$  (true). All intermediate logical signals in the network use this  $\{-1, +1\}$  domain. The advantage of this encoding is that multiplying by  $-1$  flips a signal’s truth value while multiplying by  $+1$  leaves it unchanged – and importantly, **zero is never an output of a logical unit**. This means we never encounter the situation of a gradient being multiplied by 0 (which was a major issue in the 0/1 encoding). By design, removing 0 from the state space ensures that no single input can annihilate the gradient by being zero.

**(2) Min/Max-Based Logical Operations (Non-multiplicative).** Instead of using products to represent AND/OR, RNS uses *piecewise-linear extremum functions* that exactly mimic logic under the  $\pm 1$  encoding. Specifically, if a set of inputs  $\{x_1, x_2, \dots, x_n\}$  are all either  $-1$  or  $+1$ , we define:

$$\text{AND}(x_1, \dots, x_n) = \min\{x_1, \dots, x_n\}, \quad \text{OR}(x_1, \dots, x_n) = \max\{x_1, \dots, x_n\}.$$

For example, if any  $x_j = -1$  (False), the minimum will be  $-1$ , correctly giving AND = False; only if all  $x_j = +1$  will the minimum be  $+1$  (True). Similarly, the maximum returns  $+1$  if any input is true. These operators perfectly realize the Boolean logic of AND/OR for  $\pm 1$  inputs, but unlike products, they are not multiplicative and do not cause exponential shrinkage of gradients. Instead,

they behave as selectors: the AND output is whichever input is the “most false” (most negative), and the OR output is the “most true” (most positive).

Because min and max are piecewise linear functions, we can define well-behaved subgradients for them. Suppose  $y = \max(x_1, \dots, x_n)$ . In the case of no ties, exactly one input attains this maximum value; say  $x_k$  is the largest. A small change in  $x_k$  will change  $y$  equally (a 1-to-1 slope), whereas changes in any other  $x_j$  (that are below the max) have no effect on  $y$  (as long as they don’t exceed  $x_k$ ). One convenient choice of subgradient is to assign  $\frac{\partial y}{\partial x_k} = 1$  for one of the maximal indices  $k$ , and  $\frac{\partial y}{\partial x_j} = 0$  for all other  $j \neq k$ . More generally, when there are ties (multiple inputs share the max value), the gradient can be split among them. A common subgradient is:

$$\frac{\partial y}{\partial x_i} = \begin{cases} \frac{1}{|M|} & \text{if } i \in M, \\ 0 & \text{if } i \notin M, \end{cases} \quad \text{where } M = \{j : x_j = \max_k x_k\}. \quad (48)$$

An analogous definition applies for  $y = \min(x_1, \dots, x_n)$ : the subgradient is  $1/|m|$  for inputs  $i$  that attain the minimum and 0 for others (where  $m = \{j : x_j = \min_k x_k\}$ ). In words, the gradient of a max gate is distributed equally to the input(s) that are currently “winning” (i.e., the True inputs in an OR, or the least False inputs in an AND when all are true). Crucially, this gradient *does not vanish with  $n$* : at least one input of an AND/OR receives a substantial gradient (of order 1), indicating it is responsible for the output. Even in the worst case of a tie among all  $n$  inputs (e.g., all inputs are  $-1$  for an AND, or all  $+1$  for an OR), each input would get a gradient of  $1/n$  – which decays only linearly with  $n$ , not exponentially. In most cases, only one or a few inputs determine the extremum, and they get a full magnitude gradient. There is no multiplicative chaining of many factors as in RRL’s  $F_c$  or  $F_d$  products. Additionally, negation in RNS is handled simply by a sign-flip: each literal may have a trainable indicator that either uses the variable as  $+1$  (positive literal) or negates it ( $-1$ ). Negation is just multiplication by  $-1$ , which is trivially differentiable (its subgradient is  $-1$  when active, or essentially treated as a constant factor).

**Training with Straight-Through Estimators (STE).** Both the  $\pm 1$  encoding and the use of min / max activations are still inherently non-differentiable as functions (the max has a kink where two inputs are equal, and the sign function that produces  $\pm 1$  outputs is discontinuous). However, RNS is amenable to standard techniques for training binarized networks, particularly the *straight-through estimator* (STE) for gradients. In backpropagation, whenever a gradient hits a non-differentiable threshold (such as the sign function that produces  $\pm 1$  outputs), RNS simply treats that operation as an identity mapping for the sake of gradient computation – meaning the gradient is “straight-through” passed as if the threshold were not there ?. This is a common approach in binary neural networks to approximate gradients through quantization. In the case of the min / max gates, we use the subgradient defined in Eq. equation 48 during backpropagation. The combination of STE and subgradient for min / max ensures that at *discrete* operating points, the network can still propagate meaningful gradients.

For example, consider an AND implemented as  $y = \min(x_1, \dots, x_n)$  with all inputs either  $+1$  or  $-1$ . Suppose  $y = -1$  (False) because at least one  $x_j = -1$ . In the forward pass, this yields  $y = -1$ . In the backward pass, the subgradient will assign a non-zero derivative to each input that was equal to the minimum (i.e., to each  $x_j$  that is  $-1$ ). This correctly identifies that increasing any of those  $x_j$  from  $-1$  to  $+1$  (i.e., flipping a false literal to true) would change the AND output to a higher value (potentially making the whole conjunction true if that was the only false). Thus, even though the forward function is flat for changes in any single input (since you need all falses to flip for the AND to turn true), the chosen subgradient provides a direction for learning: it tells the network to try flipping those false literals. By contrast, in a multiplicative AND, if more than one input is false, the true gradient is strictly zero for a single-input change – there is no signal at all indicating which inputs are candidates to flip. An STE cannot magically invent a meaningful gradient in that case without risking divergence from the actual function’s behavior. In RNS, however, the STE is effectively aligning with the inherent combinatorial structure of the logic: it distributes blame to all currently-false conditions for an AND (or to all currently-true conditions for an OR that outputs true, in case of ties). This yields a robust training signal even in fully discrete regimes.

## 1782 K.5 WHY RNS AVOIDS VANISHING GRADIENTS

1783  
1784 In summary, the design choices in RNS eliminate the root causes of vanishing gradients that plague  
1785 RRL’s product-based logic:

- 1786 • **Activation Form:** RRL uses multiplicative activations (product for AND, complement-  
1787 product for OR), which cause gradient magnitudes to contract multiplicatively with clause  
1788 size. Even the improved RRL with the  $P(v)$  log-smoothing still relies on a product (mod-  
1789 ulated by a corrective factor) and can suffer when many terms are far from 1. In contrast,  
1790 RNS uses min / max activations, which are piecewise linear selectors. There is no long  
1791 product over many inputs; the gradient comes from identifying the extremal input(s). This  
1792 fundamental difference means RNS does not inherently squeeze the gradient as the number  
1793 of inputs grows.
- 1794 • **Binary Encoding:** RRL’s 0/1 encoding introduces an actual zero output (false = 0) that can  
1795 outright nullify gradients (any factor of 0 in a product zeroes out the whole gradient). RNS’s  
1796  $\pm 1$  encoding avoids this so that a false value is  $-1$  instead of 0, so it never multiplicatively  
1797 annihilates an entire gradient. Every input always has the potential to influence the output  
1798 by changing sign, and the gradient methods used in RNS take advantage of that.
- 1799 • **Gradient Scaling with Clause Size:** In RRL, the magnitude of a gradient component  
1800 is on the order of  $\prod_{k \neq j} \alpha_k$  for some  $\alpha_k \in [0, 1]$  related to each of the other inputs (e.g.  
1801  $F_c(h_k, W_{i,k})$ ). This leads to an *exponential decay* in gradient as  $n$  increases, unless all  $\alpha_k$   
1802 are extremely close to 1. The log-domain trick rescales this by a factor of  $P(v)^2/v$ , which  
1803 slows the decay but does not stop it for very large  $n$ . In RNS, by contrast, a gradient to  
1804 a decisive input is  $\mathcal{O}(1)$  – it does not diminish with  $n$  at all (one input typically gets full  
1805 gradient 1 if it alone determines the output). In worst-case tie scenarios, the gradient might  
1806 be split among  $n$  inputs, giving each about  $1/n$ , i.e., decaying linearly with  $n$ , which is  
1807 far milder than exponential decay. Thus, RNS scales to clauses with many literals without  
1808 facing an overwhelming vanishing gradient issue.
- 1809 • **Discrete Training Behavior:** RRL must resort to special training techniques like gradient  
1810 grafting to handle discrete parameters, because directly backpropagating through a binarized  
1811 product-form network yields zero gradients in most places. RNS, on the other hand, can be  
1812 trained with standard backpropagation augmented with STE, which is a simpler and more  
1813 direct approach. The reason STE works well for RNS is that its surrogate gradient (identity  
1814 for the sign function, plus the min / max subgradient) aligns with actual changes that would  
1815 affect the output. In RRL’s case, an STE would have to assign gradients to inputs that in  
1816 reality do not affect the output unless combined with others – a fundamentally ambiguous  
1817 credit assignment. Therefore, RNS avoids the need for a two-model training setup; one  
1818 can optimize the  $\pm 1$  network in one go by using surrogate gradients, without the gradients  
1819 vanishing.

1820 Overall, by using  $\pm 1$  encoding and min/max logic, RNS preserves discrete interpretability while  
1821 ensuring that gradient signals remain strong and informative. The network is able to learn large  
1822 logical formulas (many-input clauses) because it never multiplies a long chain of fractional terms  
1823 during backpropagation. Each logical neuron in RNS passes a gradient to the input(s) that currently  
1824 determine its output, providing a clear learning direction. These properties enable RNS to train  
1825 effectively, where a product-based logical network would struggle or stall due to vanishing gradients.

## 1826 L CONVERGENCE OF STE WITH MIN/MAX LOGICAL ACTIVATIONS

1827  
1828 We consider a simplified RNS module consisting of a single Logic Selection Layer (LSL) with binary  
1829 inputs  $\mathbf{x} \in \{\pm 1\}^d$ . Each neuron  $u_j$  selects a subset  $S_j \subset \{1, \dots, d\}$  and an operator type via a  
1830 selector parameter  $\sigma_j \in \mathbb{R}$ : if  $\sigma_j < 0$ ,  $u_j$  is an AND gate,  $z_j(\mathbf{x}) = \min_{i \in S_j} x_i$ ; if  $\sigma_j > 0$ ,  $u_j$  is an  
1831 OR gate,  $z_j(\mathbf{x}) = \max_{i \in S_j} x_i$ . All internal signals are constrained to  $\pm 1$ , as in the main text and  
1832 Appendix J. The network output  $f_\theta(\mathbf{x})$  (including the linear prediction head) is trained on population  
1833 loss  $L(\theta) = \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}}[\ell(f_\theta(\mathbf{x}), y)]$ .

1834 **STE for min/max and selector parameters.** During backpropagation we use a straight-through  
1835 estimator (STE) for both the min/max operations and the sign of  $\sigma_j$ . For a neuron output  $z =$

1836  $\max(x_{i_1}, \dots, x_{i_k})$ , we assign

$$1837 \frac{\partial z}{\partial x_{i_m}} = \begin{cases} 1 & x_{i_m} = \max_{\ell} x_{i_\ell} \\ \frac{1}{|\mathcal{M}|} & \text{otherwise,} \end{cases} \quad (49)$$

1841 and analogously for the min operation, where  $\mathcal{M}$  is the set of inputs achieving the extremum. This  
1842 uniformly distributes the gradient among all “winning” inputs. For the selector, we treat  $\text{sign}(\sigma_j)$  as  
1843 having derivative 1 in a small band around 0 (identity STE), so that mis-specified operators receive a  
1844 non-zero gradient signal. Let  $\tilde{\nabla}L(\theta)$  denote the resulting *coarse gradient* given by STE, and consider  
1845 the full-batch update  $\theta_{t+1} = \theta_t - \eta \tilde{\nabla}L(\theta_t)$ .  
1846

1847 Yin et al. (2019) analyze coarse-gradient descent for activation-quantized networks and show that for  
1848 a “proper” STE, the expected coarse gradient is positively correlated with the true population gradient,  
1849 and its negative direction is a descent direction for the population loss; with a Lipschitz population  
1850 gradient and sufficiently small step size, the resulting algorithm converges to a critical point of the  
1851 population loss. Their framework covers one-hidden-layer architectures with non-differentiable  
1852 activations and biased gradients very similar in spirit to our min/max layer. More recently, Jeong  
1853 et al. (2023) provide finite-sample convergence guarantees for STE-based quantization-aware training  
1854 (QAT), showing that with enough samples the STE gradient iterates converge to the true binary  
1855 weights. Separately, Cheng et al. (2019a) interpret the STE as the discretization of a projected  
1856 Wasserstein gradient flow and establish convergence of the induced flow to stationary points of the  
1857 objective over discrete measures.

1858 Specializing these ideas to our min/max LSL, we assume: (1) *Population loss*:  $L(\theta)$  is defined as an  
1859 expectation over a fixed data distribution  $\mathcal{D}$  with bounded support; all inputs are  $\pm 1$ . (2) *Bounded*  
1860 *parameters*:  $\theta$  (including  $\sigma_j$ ) remains in a compact set (enforced by regularization or clipping).  
1861 (3) *Lipschitz expected gradient*: The population loss has an  $L$ -Lipschitz gradient on this set, after  
1862 smoothing via STE as in Yin et al. (2019). (4) *Proper STE alignment*: The expected coarse gradient  
1863 satisfies  $\langle \mathbb{E}[\tilde{\nabla}L(\theta)], \nabla L(\theta) \rangle \geq c \|\nabla L(\theta)\|^2$  for some  $c > 0$  whenever  $\nabla L(\theta) \neq 0$ . (5) *Small step*  
1864 *size*: The learning rate satisfies  $0 < \eta < 2/L$ .

1865 Under these assumptions, the coarse-gradient descent iterates satisfy:  $L(\theta_{t+1}) \leq L(\theta_t)$  for all  $t$   
1866 (monotone descent up to  $O(\eta^2)$  terms);  $\sum_{t=0}^{\infty} \|\nabla L(\theta_t)\|^2 < \infty$ , hence  $\|\nabla L(\theta_t)\| \rightarrow 0$ ; and every  
1867 limit point of  $\{\theta_t\}$  is a critical point of the population loss  $L$ . This is a direct adaptation of the  
1868 convergence proofs for coarse-gradient descent in activation-quantized networks Yin et al. (2019) and  
1869 for STE-based QAT Jeong et al. (2023) to our min/max logical layer. In the special case where the  
1870 target function is exactly an AND or OR over the same inputs, the only critical points with minimal  
1871 loss correspond to the correct operator selection (up to symmetries), so the method converges to a  
1872 global minimum in the simplified model.

## 1873 L.1 EMPIRICAL CONVERGENCE RESULTS

1874 To validate the theoretical convergence guarantees, we present empirical loss curves from training  
1875 RNS on eight diverse datasets. Figure 15 shows the training loss trajectories across representative  
1876 datasets spanning different domains and complexities. These results demonstrate the practical  
1877 convergence behavior of our STE-based optimization approach.  
1878

1879 Across all datasets, we observe consistent convergence patterns that align with our theoretical analysis:  
1880 (1) *Rapid initial descent*: The loss decreases sharply in early epochs, indicating effective gradient  
1881 flow through the discrete min/max operations. (2) *Monotonic decrease*: After the initial phase,  
1882 the loss exhibits monotonic or near-monotonic descent, consistent with the theoretical guarantee  
1883  $L(\theta_{t+1}) \leq L(\theta_t)$ . (3) *Stable convergence*: The loss curves plateau at low values, suggesting  
1884 convergence to critical points of the empirical risk.

1885 Different datasets exhibit varying convergence speeds and characteristics: The `wine` dataset demon-  
1886 strates extremely rapid convergence within just 50 epochs, achieving stable low loss values quickly  
1887 due to its relatively small size and clear class separation. The `dota2` dataset shows the fastest  
1888 convergence among medium-sized datasets, stabilizing within 50-75 epochs. In contrast, larger and  
1889 more complex datasets like `fashion` and `letRecog` require 200-300 epochs to reach convergence,  
with `fashion` displaying the slowest but most consistent descent due to its high dimensionality (784

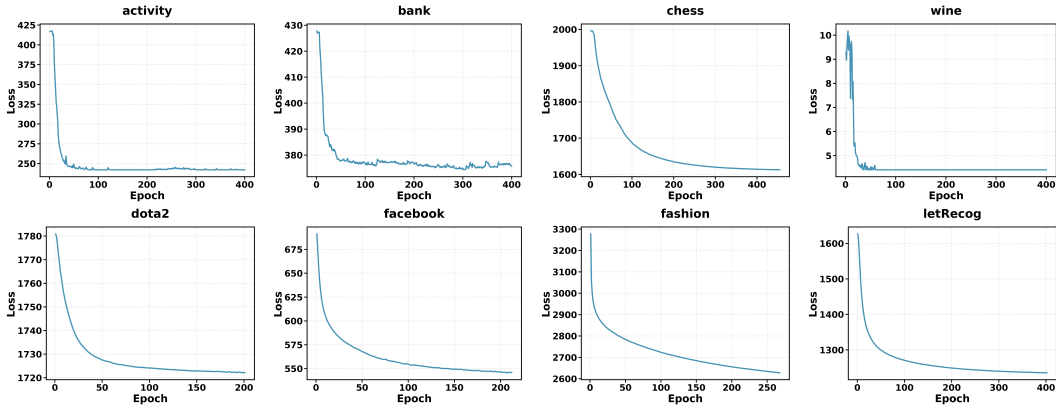


Figure 15: **Training loss convergence across diverse datasets.** Each subplot shows the training loss trajectory for RNS across datasets with varying sizes, dimensionalities, and complexities. Top row (left to right): activity, bank, chess, wine. Bottom row: dota2, facebook, fashion, letRecog. All curves exhibit rapid initial descent followed by stable convergence, validating our theoretical analysis that STE-based gradient descent with min/max logical activations converges to critical points of the loss.

features) and 10-class structure. The bank dataset exhibits slight oscillations in later epochs around epoch 300-400, likely reflecting the stochastic nature of mini-batch training on imbalanced data. Nevertheless, all datasets achieve stable convergence with diminishing gradient norms  $\|\nabla L(\theta_t)\| \rightarrow 0$ .

These empirical results strongly support our theoretical analysis. The consistent convergence across diverse datasets—ranging from small-scale multiclass classification (wine) to large-scale problems (fashion, letRecog) and from balanced (chess) to imbalanced (bank) class distributions—demonstrates that the proper STE design and min/max activations successfully enable gradient-based optimization of discrete logical networks. Importantly, the absence of training instabilities, gradient explosions, or divergence across all nine datasets confirms that our approach avoids the vanishing gradient problems that plague product-based logical activations. The monotonic descent observed in most datasets validates the theoretical guarantee that each update reduces or maintains the loss when using a properly designed STE with appropriate step size.