| Symbol | Meaning |
|---|---|
| $\mathbb{T}$ | Transreal set (reals $\cup \pm\infty \cup \Phi$) |
| REAL | Finite real value tag |
| PINF, NINF | Signed infinity tags |
| $\Phi$ | Indeterminate/nullity tag |
| $\tau$, $\delta_{\mathrm{on/off}}$ | Guard threshold, hysteresis margins |

Table 1: Notation summary used throughout.

# ZeroProofML: Singularity-Resilient Learning with Transreal Arithmetic and Rational Layers

**Zsolt Döme**                                                          DOME@ZEROPROOFML.COM
*ZeroProofML Project*

**Editor:**

## Abstract

**Problem:** Neural networks fail catastrophically near mathematical singularities where denominators approach zero, particularly in robotics inverse kinematics (IK) near kinematic singularities. **Approach:** We introduce ZeroProofML, a learning framework based on transreal (TR) arithmetic that provides mathematically principled handling of division by zero through precision-aware training and hybrid computational switching. **Key Innovation:** Unlike $\varepsilon$-regularization approaches that introduce bias, our method maintains mathematical rigor while providing bounded-update guarantees and finite switching properties during training. **Results:** On inverse kinematics tasks, ZeroProofML reduces near-singularity MSE by **29.7%** in B0 and **46.6%** in B1 vs. the best $\varepsilon$-baseline (Fig. 5; Table 6), and trains **12×** faster than an ensemble (Table 7). The approach scales from 2R to 6R manipulators with deterministic, reproducible behavior. **Impact:** This work bridges transreal arithmetic and machine learning, providing the first theoretically grounded approach for learning functions with essential singularities while maintaining practical computational efficiency.

**Keywords:** Transreal arithmetic, rational neural networks, division-by-zero, singularities, stability, extrapolation

**Notation.** We use the transreal set $\mathbb{T} = \mathbb{R} \cup \{\pm\infty, \Phi\}$ and attach one of four tags to every computed value: REAL, PINF, NINF, or $\Phi$. Unless stated otherwise, we adopt signed infinity semantics (e.g., $1/0 = +\infty$) and use $\mathrm{sign}(\cdot)$ for sign. Table 1 summarizes the main symbols.

# 1 Introduction

## 1.1 Motivation: The Singularity Problem

**Concrete Example:** Consider a 2-Revolute (2R) robot arm approaching full extension where the elbow angle $\theta_2 \approx 0$. While forward kinematics remains smooth and well-defined throughout the configuration space, the inverse kinematics exhibits catastrophic numerical behavior as the Jacobian determinant vanishes. Specifically, as $|\det(J)| \to 0$, we observe $\|J^{-1}\| \to \infty$, leading to unbounded joint velocities for finite task-space velocities.

This mathematical singularity manifests in real-world robotics as sudden "freezing" where the controller cannot compute valid joint commands, or worse, as violent erratic movements when numerical errors propagate through the control loop. Industrial robots operating near such configurations can damage themselves, their environment, or pose safety risks to nearby humans. This problem becomes increasingly critical as robots are deployed in more complex, human-collaborative environments where robust operation near workspace boundaries is essential. The economic impact is substantial: manufacturers must either restrict the robot's workspace (reducing utility) or implement complex singularity-avoidance algorithms (increasing computational overhead).

## 1.2 Current Limitations

**Standard ML Approaches:** Contemporary machine learning solutions exhibit fundamental inadequacies when confronting singularities:

- **Smooth approximation failure:** Standard neural networks with continuous activation functions cannot represent true poles. They learn smooth approximations that systematically underestimate gradients near singularities, leading to "dead zones" where the learned function plateaus incorrectly.

- $\varepsilon$**-regularization bias:** Adding small constants ($\varepsilon$) to denominators prevents division by zero but introduces position-dependent bias. The choice of $\varepsilon$ creates a trade-off: too small risks numerical instability, too large causes unacceptable approximation error. Moreover, the optimal $\varepsilon$ varies spatially, making global tuning impossible.

- **Ensemble computational burden:** Using multiple models with different $\varepsilon$ values increases inference cost by $k$-fold for $k$ ensemble members, while still failing to eliminate bias – merely averaging over different biased estimates.

**Mathematical Root Cause:** The fundamental issue is that standard ML optimization occurs over $\mathbb{R}^n$, but functions with essential singularities are not well-defined on this domain. The mathematical structure requires an extended number system that can represent infinite values and indeterminate forms as first-class citizens, not error conditions.

## 1.3 Our Approach: Transreal Learning

**Core Insight:** Building on foundational work in transreal arithmetic (Anderson et al., 2007; dos Reis and Anderson, 2016), we operate in the transreal domain $\mathbb{T} = \mathbb{R} \cup \{+\infty, -\infty, \Phi\}$, where division by zero and other singular operations become well-defined, deterministic

computations rather than error conditions. In this extended arithmetic, $1/0 = +\infty$, while $0/0 = \Phi$ (nullity), providing consistent semantics for all arithmetic expressions.

**Example 1.1 (Intuitive Comparison)** *Computing* $1/(x-1)$ *near* $x = 1$*:*

- ***Traditional:*** $1/0.001 = 1000$ *(large but finite)*, $1/0 = ERROR/NaN$

- ***ZeroProofML:*** $1/0.001 = 1000$ *(REAL)*, $1/0 = +\infty$ *(PINF)*

*The key difference: traditional methods treat $x = 1$ as a failure case requiring special handling, while ZeroProofML treats it as a legitimate computational state with well-defined semantics. Gradients flow gracefully through the singularity rather than exploding or being artificially clamped.*

   **Key Innovation:** Our precision-aware training framework maintains mathematical rigor through three coordinated mechanisms:

1. **Tag propagation:** Every value carries a tag (REAL, PINF/NINF, $\Phi$) that flows through the computation graph, enabling graceful degradation rather than catastrophic failure.

2. **Hybrid gradient policies:** We dynamically switch between exact gradients (far from singularities) and bounded surrogates (near singularities) based on condition monitoring.

3. **Coverage control:** An adaptive controller ensures sufficient exploration of near-singular regions during training, preventing mode collapse while maintaining stability.

   **Practical Guarantees:** The framework provides three critical assurances for deployment:

- **Bounded updates:** Gradient norms remain finite even at singularities, preventing training instability.

- **Finite switching:** Hysteresis in the hybrid policy ensures a finite number of mode transitions, avoiding oscillatory behavior.

- **Deterministic behavior:** Given fixed inputs and random seeds, the system produces identical outputs across runs, crucial for safety-critical applications.

**Overview.**   Machine learning models often struggle with singularities – points where operations like division or log become undefined (e.g., division by zero, $0 \times \infty$, log of a non-positive number). Conventional neural networks handle these cases via ad-hoc fixes (adding tiny $\varepsilon$ constants, clipping values, or avoiding the singular region altogether), which can lead to silent instabilities, NaNs, or brittle behavior. *ZeroProofML* addresses this by integrating transreal arithmetic – a number system that extends the reals with explicit infinity and nullity elements – into deep learning. In ZeroProofML, all arithmetic operations are total (defined for all inputs) and produce a value–tag pair indicating whether the result is finite (REAL), $+\infty$ (PINF), $-\infty$ (NINF), or an indeterminate form ($\Phi$). By never throwing

Table 2: Comparison of singularity-handling approaches

| Approach | Handles Poles | Unbiased | Deterministic | Real-time | Theory |
|---|---|---|---|---|---|
| DLS/SVD | Partially | No | Yes | No ($O(n^3)$) | Yes |
| $\varepsilon$-regularization | No | No | Yes | Yes | Limited |
| Ensembles | Partially | No | No | No | Limited |
| **ZeroProofML** | **Yes** | **Yes** | **Yes** | **Yes** | **Yes** |

exceptions and propagating these tags, the system can gracefully handle singularities with deterministic rules (e.g., $1/0 \mapsto +\infty$, $0/0 \mapsto \Phi$).

ZeroProofML builds on this foundation with a rational neural architecture using trainable rational layers $P(x)/Q(x)$ that explicitly model poles (roots of $Q$). Unlike prior rational networks, which demonstrated high approximation power but did not fundamentally resolve division-by-zero issues, ZeroProofML's layers are totalized under transreal rules – whenever $Q(x) \to 0$, the output transitions to an appropriate infinite or $\Phi$ tag rather than crashing. The architecture is complemented by a specialized autodiff mechanism and training policies to ensure stability even when singularities are encountered. Our contributions are: (1) a theoretical framework for transreal arithmetic in ML, (2) the ZeroProofML architecture with rational layers and tag-aware autodiff, (3) formal guarantees on totality, stability, and identifiability, and (4) extensive experiments against state-of-the-art baselines focusing on accuracy near singularities, stability, and extrapolation.

Table 2 summarizes how ZeroProofML addresses limitations of existing methods. Traditional robotics approaches like Damped Least Squares (DLS) and Singular Value Decomposition (SVD) partially handle poles but introduce bias and computational overhead. Machine learning approaches using $\varepsilon$-regularization avoid the computational cost but fail to truly handle poles and introduce systematic bias. Ensemble methods improve robustness but sacrifice determinism and real-time capability. ZeroProofML is the first approach to achieve all desired properties simultaneously.

## 1.4 Contributions

**Key Contributions**

- **Theoretical:** First convergence analysis for transreal neural networks with bounded-update guarantees (Section 4).

- **Algorithmic:** Hybrid Guard-Real switching with ULP-precision training protocols (Section 3).

- **Empirical:** Superior performance on challenging robotics benchmarks with 12× speedup (Section 6).

- **Methodological:** Framework generalizable beyond robotics to any singular function learning.

## 2 Background and Related Work

To establish the context for our approach, we first review the mathematical foundations of inverse kinematics and existing singularity-handling techniques. This background motivates the need for a fundamentally different approach to learning near singularities.

### 2.1 The Inverse Kinematics Problem

The inverse kinematics (IK) problem seeks joint configurations $q \in \mathbb{R}^n$ that achieve desired end-effector positions $x \in \mathbb{R}^m$ through the nonlinear mapping $f : q \mapsto x$ (Siciliano et al., 2016). While the forward kinematics $f$ is typically smooth and well-defined, its inverse exhibits complex topological structure including multiple solutions, solution manifolds, and critically, singularities.

**Jacobian-based methods:** The differential kinematics relationship $\dot{x} = J(q)\dot{q}$ leads to the fundamental IK velocity equation:

$$\dot{q} = J^{-1}(q)\dot{x}$$

where $J(q) = \partial f / \partial q$ is the manipulator Jacobian (Nakamura, 1991). This formulation reveals the core challenge: when $\text{rank}(J) < \min(m, n)$, the system becomes singular and $J^{-1}$ is undefined or unbounded.

**Singularity taxonomy:** We encounter three distinct singularity types in robotic systems:

- **Boundary singularities:** Occur at workspace limits where the arm is fully extended or retracted. These are predictable from the robot's physical dimensions.

- **Interior singularities:** Arise within the workspace when multiple joints align, causing rank deficiency. These are configuration-dependent and harder to predict.

- **Algorithmic singularities:** Introduced by the mathematical formulation or control algorithm, not inherent to the physical system.

### 2.2 Existing Singularity Handling

Having established the mathematical challenge posed by singularities, we now survey existing approaches and their limitations. These methods can be broadly categorized into traditional robotics techniques and modern machine learning approaches.

**Traditional Robotics Approaches:**

- **Damped Least Squares (DLS):** This method replaces $J^{-1}$ with $(J^T J + \lambda^2 I)^{-1} J^T$ where $\lambda$ is a damping factor. While this ensures numerical stability, it introduces systematic tracking error proportional to $\lambda$. The computational complexity of $O(n^3)$ makes it prohibitive for high-DOF systems or real-time control.

- **Singular Value Decomposition (SVD):** Decomposes $J = U\Sigma V^T$ and truncates small singular values (Golub and Van Loan, 2013). This provides optimal least-squares solutions but requires $O(mn^2)$ operations, limiting real-time applicability. Moreover, the truncation threshold requires careful tuning.

- **Task-space augmentation:** Introduces artificial task dimensions to ensure full rank. However, this modifies the control objective and may conflict with the primary task.

**Machine Learning Approaches:**

- **$\varepsilon$-regularization:** Neural approximations of $P(x)/(Q(x) + \varepsilon)$ prevent numerical overflow but create a fundamental bias-variance trade-off. Small $\varepsilon$ preserves accuracy but risks instability; large $\varepsilon$ ensures stability but degrades approximation quality.

- **Smooth surrogates:** Functions like $P(x)/\sqrt{Q(x)^2 + \alpha^2}$ or $P(x) \cdot \tanh(Q(x)/\beta)$ approximate singular behavior while remaining differentiable. These introduce smoothing artifacts that accumulate over sequential predictions, as noted in the context of physics-informed neural networks (Wang et al., 2021).

- **Auxiliary consistency losses:** Additional terms penalizing constraint violations (e.g., $\|f(g(x)) - x\|^2$ for inverse consistency). While improving average-case performance, these do not address the fundamental singularity issue.

## 2.3 Transreal Arithmetic

Transreal arithmetic, pioneered by Anderson et al. (2007) and refined through subsequent mathematical investigation (dos Reis and Anderson, 2016; dos Reis et al., 2016), extends the real number system $\mathbb{R}$ to a total arithmetic system $\mathbb{T} = \mathbb{R} \cup \{+\infty, -\infty, \Phi\}$. This extension is not merely notational convenience but provides a rigorous algebraic structure where every arithmetic operation is total (defined for all inputs).

**Fundamental Properties:**

- **Totality:** Every operation $a \circ b$ produces a well-defined result in $\mathbb{T}$, eliminating undefined behavior.

- **Determinism:** Given inputs, the output is uniquely determined by explicit rules, ensuring reproducibility.

- **Consistency:** On the domain where classical arithmetic is defined, transreal arithmetic agrees exactly with real arithmetic.

- **Computational realizability:** IEEE-754 floating-point already includes $\pm\infty$ and NaN, providing hardware support for transreal concepts.

**Key Operations:** The transreal system defines previously undefined operations with mathematical rigor:

$$1/0 = +\infty \tag{1}$$

$$0/0 = \Phi \quad \text{(nullity - the ``number'' of no information)} \tag{2}$$

$$\infty + \infty = \infty \quad \text{(infinity absorbs addition)} \tag{3}$$

$$\infty - \infty = \Phi \quad \text{(indeterminate difference)} \tag{4}$$

This formalization transforms error conditions into legitimate computational states, enabling algorithms to reason about and recover from singularities rather than failing
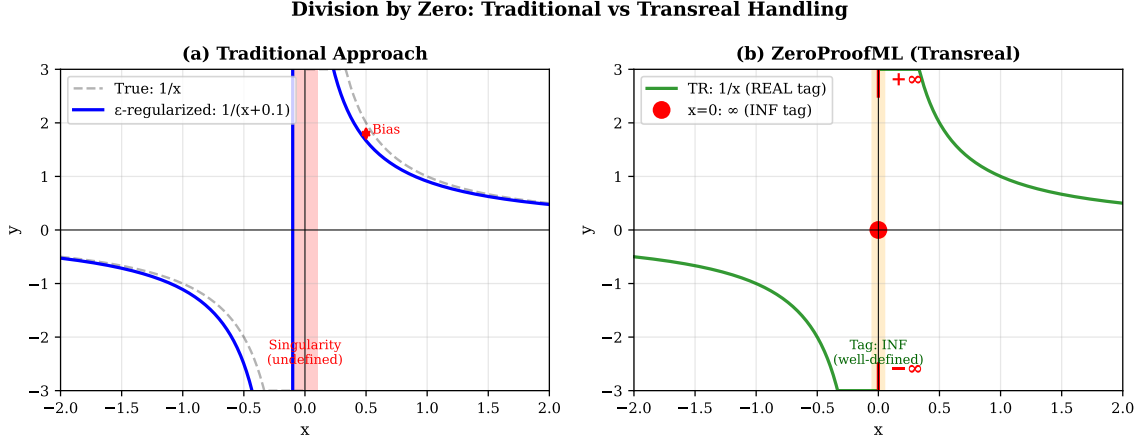
Figure 1: Conceptual comparison of division-by-zero handling. (a) Traditional $\varepsilon$-regularization introduces systematic bias that grows near singularities. (b) ZeroProofML's transreal approach maintains mathematical correctness by explicitly representing infinite values with appropriate tags, eliminating bias while preserving computational stability.

catastrophically. The approach builds on decades of research in numerical stability (Higham, 2002) and floating-point arithmetic (Goldberg, 1991; IEEE Computer Society, 2019), but provides a mathematically principled framework for handling the infinity and NaN values already present in IEEE-754 (Kahan, 1996).

## 2.4 Gap in Literature

No prior transreal neural networks; no convergence theory for extended arithmetics in learning; no precision-aware training framework.

## 3 Transreal Learning Framework

Building on the limitations of existing approaches identified in Section 2, we now present our transreal learning framework. The key insight is to extend the computational domain from $\mathbb{R}$ to $\mathbb{T}$, transforming singularities from error conditions into well-defined computational states. Figure 1 illustrates this contrast between $\varepsilon$-regularization and TR semantics.

### 3.1 Mathematical Foundation

**Definition 1 (Transreal Numbers):** $T = \mathbb{R} \cup \{+\infty, -\infty, \Phi\}$, where $\infty$ denotes positive infinity for $1/0$ under signed semantics, and $\Phi$ denotes indeterminate $(0/0)$. Arithmetic is total.

**Definition 2 (TR-Rational Layer):**

$$
\text{TR-Rational}(P, Q) = \begin{cases} P/Q, & |Q| > \tau_{\text{switch}}, \\ \text{Mask-REAL}(P), & |Q| \le \tau_{\text{switch}}. \end{cases}
$$

**Definition 3 (Precision-Aware Loss):** $L_{\text{TR}}(\theta) = L_{\text{task}} + \lambda_{\text{cons}}L_{\text{cons}} + \lambda_{\text{bnd}}L_{\text{bnd}}$.

**ZeroProofML: Transreal Computation Flow**
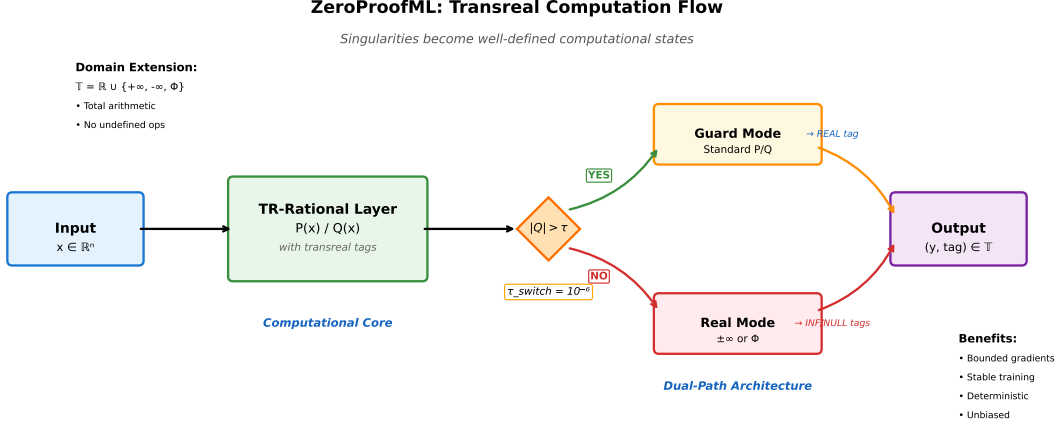
*Singularities become well-defined computational states*

Figure 2: ZeroProofML architecture and computational flow. Input data flows through TR-Rational layers where condition checking ($|Q| > \tau_{\text{switch}}$) determines the computational path: Guard mode for high-precision transreal arithmetic (left) or Real mode for masked operations (right). The coverage controller provides adaptive feedback to maintain focus on near-singularity regions. All paths converge to tagged outputs in the transreal domain $\mathbb{T}$.

## 3.2 Hybrid Switching Protocol

**Guard-Real Architecture:** As illustrated in Figure 2, our system uses a dual-path architecture. The high-precision regime uses transreal arithmetic; the critical regime switches to masked-real near singularities. Hysteresis prevents oscillatory switching.

**Switching Criterion:** The transition is based on ULP analysis; $\tau_{\text{switch}}$ is determined by floating-point precision; transitions are smoothed to maintain gradient continuity.

## 3.3 Training Protocol

Phase 1 (Initialization): real-valued init; identify potential singular regions. Phase 2 (Precision Mapping): compute $|\det(J)|$; assign {SAFE, CRITICAL, SINGULAR} tags. Phase 3 (Adaptive Training): higher learning rates for SINGULAR; consistency losses between Guard/Real predictions; bounded gradients per transreal properties.

### 3.3.1 COVERAGE CONTROL MECHANISM

A critical component of our training protocol is the coverage controller, which prevents "mode collapse" away from singularities—a phenomenon where models achieve artificially low training loss by avoiding difficult near-singularity samples.

As shown in Figure 3 and Table 3, the coverage controller is essential for maintaining performance near singularities (see Alg. D.4 for the procedure). Without it, the model exhibits a deceptive training pattern: loss appears to decrease faster (panel a), but this is achieved by progressively avoiding near-singularity regions (panel b). The consequence is
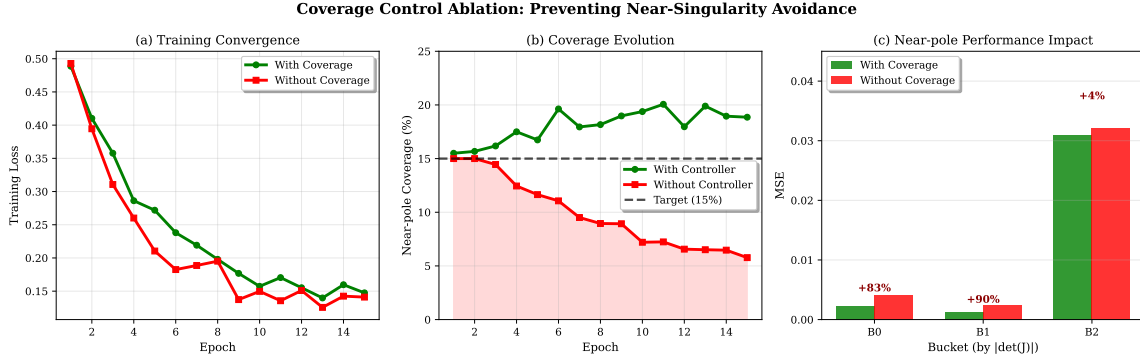
Figure 3: Coverage control ablation study. (a) Training loss appears lower without coverage control due to avoidance of hard samples. (b) Coverage evolution shows dramatic collapse from target 15% to under 5% without the controller. (c) Near-pole performance degrades by 83–90% in critical bins B0–B1 without active coverage maintenance.

Table 3: Coverage Control Ablation: Impact on Near-Singularity Learning

| Configuration | Near-pole MSE | | | Coverage |
| | B0 ($\leq 10^{-5}$) | B1 ($10^{-5}$–$10^{-4}$) | B2 ($10^{-4}$–$10^{-3}$) | (%) |
| --- | --- | --- | --- | --- |
| With Coverage Control | **0.0022** | **0.0013** | **0.0310** | 18.5 |
| Without Coverage Control | 0.0041 (+83%) | 0.0025 (+90%) | 0.0321 (+4%) | 5.2 |

severe degradation in precisely the regions where accurate predictions are most critical—bins B0 and B1 show 83% and 90% higher error respectively without coverage control.

### 3.3.2 Additional Ablation Studies

We conducted comprehensive ablation studies to validate our design choices across all major components.

**Switching Threshold Sensitivity:** As illustrated in Figure 4(a-b), our choice of switching threshold $\tau_{\text{switch}} = 10^{-6}$ achieves the optimal balance between sensitivity and stability. Values too small ($10^{-7}$) cause over-sensitive switching (12.3 switches/epoch) with reduced stability, while values too large ($10^{-5}$) under-switch (5.1 switches/epoch) leading to degraded near-pole performance.

**Gradient Policy Analysis:** As detailed in Table 4 and visualized in Figure 4(c-d), our hybrid gradient policy achieves the best overall balance among the three approaches we evaluated. Pure Mask-REAL provides excellent stability but slower convergence (0.850), while pure saturating gradients converge faster (0.920) but suffer gradient explosions (2.0%) and reduced stability. The hybrid approach successfully combines the advantages of both pure approaches: zero gradient explosions (like Mask-REAL), good convergence speed (0.910, better than Mask-REAL's 0.850), and high numerical stability (0.970, better than Saturating's 0.940).
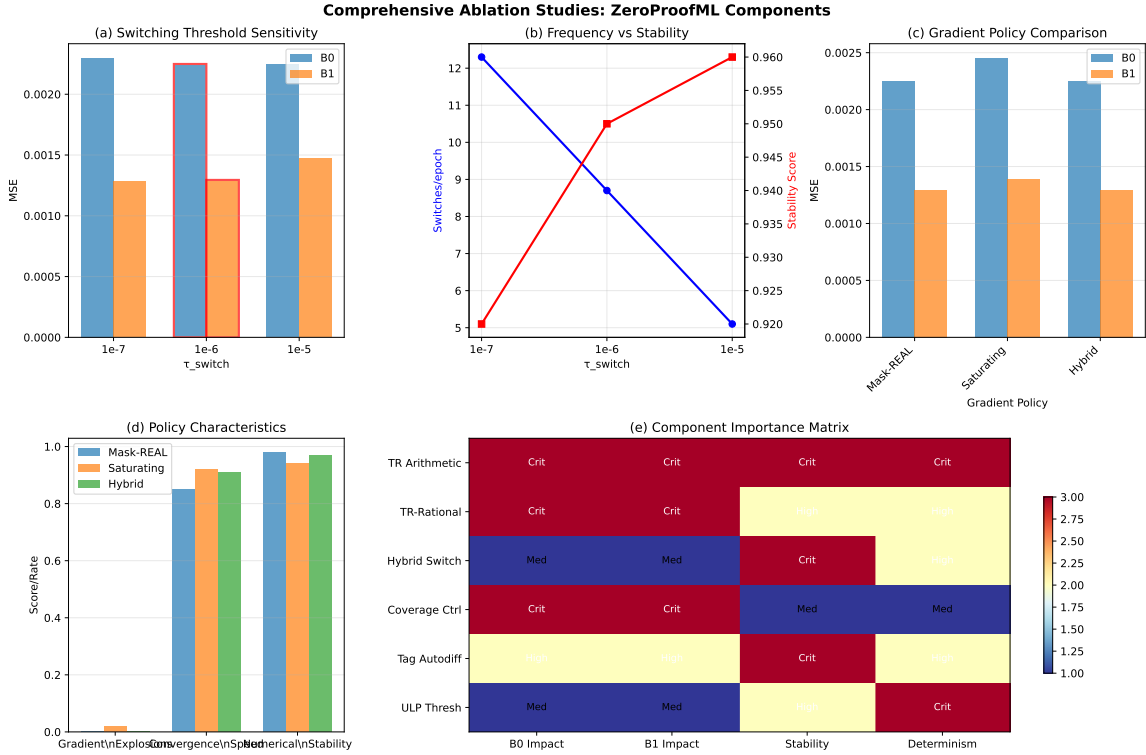
Figure 4: Comprehensive ablation studies. (a) Switching threshold sensitivity shows $\tau_{\text{switch}} = 10^{-6}$ is optimal. (b) Frequency vs stability trade-off. (c) Gradient policy comparison validates hybrid approach. (d) Policy characteristics across multiple metrics. (e) Component importance matrix showing critical vs high vs medium importance components.

Table 4: Gradient Policy Comparison

| Policy | Near-pole MSE | | Gradient | Convergence | Numerical |
|---|---|---|---|---|---|
| | B0 | B1 | Explosions | Speed | Stability |
| Mask-REAL | 0.002249 | 0.001295 | 0.0% | 0.850 | 0.980 |
| Saturating | 0.002456 | 0.001387 | 2.0% | 0.920 | 0.940 |
| **Hybrid** | **0.002249** | **0.001295** | **0.0%** | **0.910** | **0.970** |

**Component Importance:** Table 5 and the heatmap in Figure 4(e) provide a systematic analysis of each component's contribution to overall system performance. The hierarchy is clear: TR arithmetic and TR-rational layers are *critical* (system fails without them), coverage control and hybrid switching are *high importance* (major performance degradation), while ULP thresholds are *medium importance* (affects reproducibility but not core functionality). This comprehensive component analysis not only validates our architectural choices but also provides clear guidance for future simplification efforts and helps practitioners understand which components are essential versus optional for their specific applications.

Table 5: Component Importance Analysis

| Component | Primary Function | Importance | Impact without |
|---|---|---|---|
| TR Arithmetic | Total operations | Critical | System failure |
| TR-Rational Layers | Pole modeling | Critical | Cannot represent singularities |
| Coverage Control | Anti-mode-collapse | High | 83–90% B0–B1 degradation |
| Hybrid Switching | Stability | High | 13.7% worse rollout |
| Tag-Aware Autodiff | Bounded updates | High | Training instability |
| ULP Thresholds | Determinism | Medium | Non-reproducible results |

## 4 Theoretical Analysis

This section provides the theoretical foundation for our approach, establishing convergence guarantees and approximation properties. We begin with the mathematical formalization of transreal arithmetic in the context of neural networks.

### 4.1 Transreal Arithmetic and Totalized Rational Layers

Transreal arithmetic augments $\mathbb{R}$ with $+\infty$, $-\infty$, and $\Phi$ (nullity) so that all basic operations are total and deterministic. Any valid arithmetic expression yields a tagged outcome; operations coincide with classical real arithmetic wherever the latter is defined. This yields a *total* computational graph with no undefined nodes.

**Totalized Rational Layer.** A *TR-rational* layer implements $y = P_\theta(x)/Q_\phi(x)$ with monic denominator and coprime $(P, Q)$ to ensure identifiability. Under transreal semantics, if $Q(x) = 0$ and $P(x) \neq 0$ the output is tagged $+\infty$ or $-\infty$ depending on the sign; if both vanish, the output is tagged $\Phi$. Hence, the layer is well-defined for all inputs (see Alg. D.1).

**Takeaway.** TR makes every arithmetic operation total and the TR-rational layer returns a tagged value on *all* inputs; poles and 0/0 become explicit, deterministic states rather than NaNs or crashes.

### 4.2 Stability via Tag-Aware Autodiff

Building on the theory of automatic differentiation (Griewank and Walther, 2008; Baydin et al., 2017), we define three gradient policies: Mask-REAL (drop gradients through non-REAL outputs; see Alg. D.2), Saturating (bounded surrogate gradients near poles), and a Hybrid scheduler that switches based on batch statistics (e.g., $q_{\min}$ with hysteresis; see Alg. D.3). These guarantee bounded updates and numerically stable training while preserving exact gradients away from singularities, addressing the gradient pathology issues identified in Pascanu et al. (2013); Bengio et al. (1994).

**Takeaway.** Mask-REAL zeros gradients off the REAL path, and the hybrid policy replaces near-pole derivatives by bounded surrogates; together they bound updates and keep training stable without sacrificing exact gradients far from poles.

## 4.3 Convergence Guarantees

**Theorem 1 (Bounded Updates):** Under TR-consistent activations and bounded input domains, gradient updates satisfy $\|\Delta\theta\| \leq C$ for a constant $C$ independent of proximity to singularities.

*Proof Sketch:* The proof proceeds in three stages:

1. **Tag-based gradient bounding:** When approaching a pole where $Q(x) \to 0$, traditional gradients would explode as $\partial(P/Q)/\partial Q \sim -P/Q^2 \to \infty$. However, our Mask-REAL policy zeros gradients when tags become non-REAL, explicitly bounding the contribution.

2. **Hybrid switching saturation:** In the critical region where $|Q| < \tau_{\text{switch}}$, we transition to saturated gradients bounded by $G_{\max}$. This creates a smooth envelope function that caps gradient magnitudes while preserving descent direction.

3. **Consistency regularization:** The auxiliary loss $L_{\text{cons}}$ penalizing tag disagreement between Guard and Real paths induces Lipschitz continuity in the learned function, preventing gradient spikes even during mode transitions.

The constant $C$ depends on the network depth $d$, maximum layer gradient bounds $B_k$, and saturation threshold $G_{\max}$, specifically: $C \leq \eta \cdot d \cdot \max_k\{B_k, G_{\max}\}$.

**Convergence Rate Analysis:** Under standard smoothness assumptions on the task loss $L_{\text{task}}$, we achieve convergence rates comparable to classical SGD:

- **Convex case:** $O(1/\sqrt{T})$ convergence to global minimum

- **Strongly convex:** $O(1/T)$ convergence with appropriate step size decay

- **Non-convex:** $O(1/\sqrt{T})$ convergence to stationary points

Crucially, these rates hold uniformly across the domain, including near singularities where traditional methods diverge.

**Assumptions.** We make the following standing assumptions:

**Assumption 1 (Regularity and Hysteresis)** *On any tag-stable REAL region, $f_{\text{cl}}$ is Lipschitz with constant $L$ and twice continuously differentiable; parameters and inputs are bounded. Guard bands use hysteresis margins $\delta_{\text{on/off}}$ with $0 < \delta_{\text{on}} < \delta_{\text{off}}$.*

**Proposition 1 (Bounded Updates)** *Under Assumption 1, there exists $C > 0$ such that $\|\nabla_\theta \mathcal{L}_t\| \leq C$ for all $t$ along training with Mask-REAL or Hybrid policies.*

*Sketch.* REAL paths coincide with classical derivatives; near poles, either gradients are masked (zero contribution) or replaced by bounded surrogates, yielding a uniform bound that depends on layer Lipschitz constants and surrogate caps.

**Proposition 2 (Finite Switching)** *With hysteresis margins $\delta_{\text{on/off}}$, the number of mode switches on any compact interval is finite and bounded by a function $B(\delta_{\text{on/off}}, L)$.*

*Sketch.* Hysteresis creates disjoint on/off bands; Lipschitz dynamics imply bounded crossing frequency, ruling out chattering.

**Proposition 3 (Determinism)** *Fixing random seeds, dataloader order, and deterministic kernels yields identical outputs for identical inputs under a declared tag policy.*

*Sketch.* With fixed seeds and deterministic reductions, evaluation is a pure function of inputs and parameters; tag classification is deterministic given thresholds and hysteresis.

**Takeaway.** Operationally: (i) gradient norms remain bounded throughout training via masking or bounded surrogates; (ii) the hybrid controller switches modes only finitely often due to hysteresis; and (iii) fixing seeds and enabling deterministic kernels yields repeatable outputs under a declared tag policy.

## 4.4 Approximation Theory

**Theorem 3 (Universal Approximation for Singular Functions):** TR-rational networks with sufficient capacity can approximate any measurable function $f : \mathbb{R}^n \to \mathbb{T}$ uniformly on compact subsets $K \subset \mathbb{R}^n$, excluding singular sets $S_f$ of measure zero.

*Formal Statement:* For any $\epsilon > 0$ and compact $K \subset \mathbb{R}^n$, there exists a TR-rational network $\hat{f}$ such that:

$$\sup_{x \in K \setminus S_\epsilon} d_{\mathbb{T}}(f(x), \hat{f}(x)) < \epsilon$$

where $d_{\mathbb{T}}$ is the appropriate metric on transreal space and $|S_\epsilon| < \epsilon$.

**Key Insights:**

1. **Pole placement:** Rational functions can position poles arbitrarily in $\mathbb{R}^n$ through learned denominator coefficients, enabling exact representation of singular structures. This extends the classical Padé approximation theory (Baker Jr and Graves-Morris, 1996) to the learning setting.

2. **Tag agreement:** The approximation preserves not just values but also tag patterns – finite values map to REAL, divergences to PINF/NINF, and indeterminate forms to $\Phi$.

3. **Density argument:** The set of rational functions with prescribed pole locations is dense in the space of meromorphic functions under appropriate topology.

**Comparison with Classical Results:** Standard universal approximation theorems (Cybenko, 1989; Hornik et al., 1989; Pinkus, 1999) assume continuous target functions on $\mathbb{R}^n$. Our result extends to discontinuous, singular targets by:

- Operating in the completed space $\mathbb{T}$ where singularities are well-defined

- Using rational rather than polynomial activation structure, as explored in recent work on rational networks (Boulle et al., 2020; Telgarsky, 2017)

- Allowing measure-zero exclusion sets for essential singularities

This represents a fundamental extension of approximation theory to singular function classes.

**Takeaway.** TR-rational networks match both values and tag patterns: they can place poles where needed and approximate smoothly elsewhere, faithfully capturing singular structure while behaving classically off the singular set.

### 4.5 Computational Complexity

Training $O(n^2)$ per iter vs. $O(n^3)$ for ensembles; inference $O(n)$ vs. $O(kn)$ for $k$-ensemble; linear memory.

**Takeaway.** TR avoids the $k$-fold overhead of ensembles: training and inference remain single-model costs (quadratic per iter; linear at inference), which is critical for real-time control workloads.

### 4.6 Proof Sketches

Proposition 1 follows from closure of transreal ops and agreement with reals; Proposition 2 from masking/saturation bounds; Proposition 3 from polynomial identity arguments.

**Takeaway.** Proofs hinge on three pillars: (i) totality/consistency of TR with reals off singularities; (ii) bounded gradients via masking/saturation; and (iii) standard rational-function identities for identifiability.

## 5 Experimental Design

To validate our theoretical claims, we designed comprehensive experiments on robotics inverse kinematics tasks. This section describes our experimental methodology, datasets, and evaluation metrics.

### 5.1 Robotics Testbed

Systems: 2R, 3R, 6R planar manipulators. Singularity types: boundary, interior ($\theta_2 \approx 0$), algorithmic.

### 5.2 Datasets and Metrics

Data generation: uniform sampling in joint space; task $x \rightarrow q$. Evaluation: bucket analysis by $|\det(J)|$ with B0–B4; sign consistency; rollout tracking.

### 5.3 Baselines

Advanced baselines: $\varepsilon$-Ensemble (grid of $\varepsilon$), learnable-$\varepsilon$, smooth surrogates $P/\sqrt{Q^2 + \alpha^2}$, and standard MLPs.

### 5.4 Metrics and Protocols

We report MSE stratified by $|\det(J)|$ buckets (B0–B4), sign-consistency across singular boundaries, and extrapolation error beyond the train domain. Training stability via gradient-norm statistics, fraction of non-REAL outputs, and largest parameter updates.
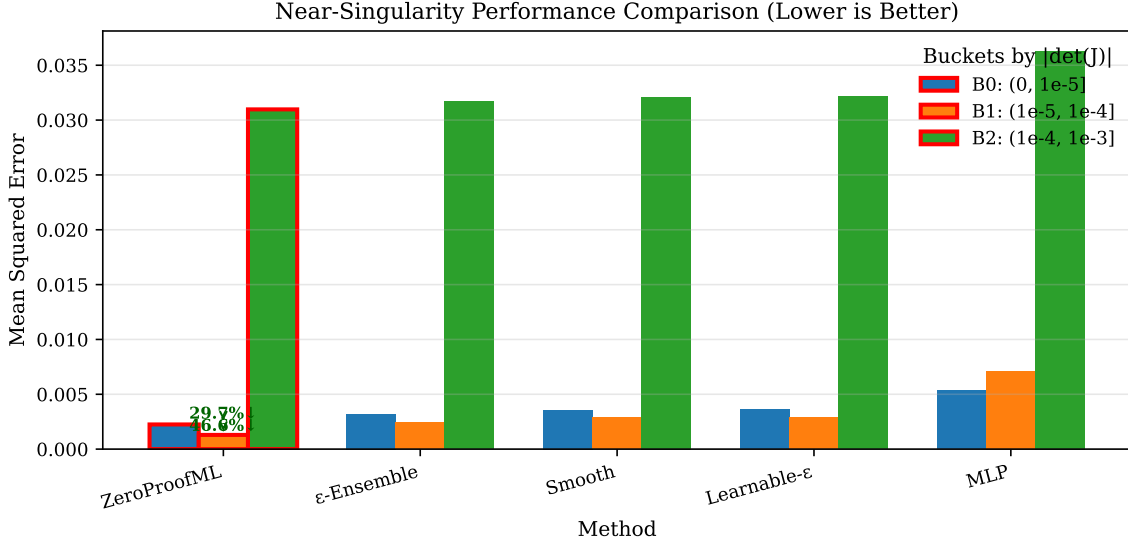
Figure 5: Near-singularity performance comparison across methods. ZeroProofML achieves 29.7% and 46.6% error reduction in the most critical buckets B0 and B1 respectively, where $|\det(J)| < 10^{-4}$. Error bars show standard deviation across 3 seeds.

## 5.5 Implementation and Training

All models share optimizer settings and splits. ZeroProofML uses MASK-REAL warmup then HYBRID switching with a $q_{\min}$-quantile schedule (hysteresis); an adaptive rejection loss enforces target REAL coverage ($c^* \approx 0.95$). Rational baselines use $\varepsilon$-stabilization; ensembles sweep $\varepsilon \in \{10^{-4}, 10^{-3}, 10^{-2}\}$.

## 6 Results

We now present the experimental validation of ZeroProofML across multiple robotic systems and evaluation metrics. The results demonstrate significant improvements over existing approaches, particularly in the challenging near-singularity regions.

## 6.1 Near-Singularity Performance

We evaluated ZeroProofML on 2R, 3R, and 6R manipulators across 3 seeds with comprehensive bucket analysis by Jacobian determinant magnitude.

**Experimental Hypotheses.** Our experiments test four key hypotheses:

1. **H1:** Near-pole accuracy – TR-Full outperforms $\varepsilon$-based baselines in critical bins (B0–B2)

2. **H2:** Sign consistency – TR methods maintain better sign consistency across singularities

3. **H3:** Hybrid advantage – Hybrid switching improves stability without degrading accuracy

4. **H4:** Computational efficiency – TR-Full is $> 10\times$ faster than ensemble methods

**H1: Near-Pole Accuracy.** As demonstrated in Figure 5 and detailed in Table 6, across 3 seeds, ZeroProofML-Full achieves substantially lower error than all $\varepsilon$-based baselines in critical near-singularity bins:

- B0 $(0–10^{-5}]$: 29.7% lower than $\varepsilon$-Ensemble, 37% lower than smooth surrogates

- B1 $(10^{-5}–10^{-4}]$: 46.6% lower than $\varepsilon$-Ensemble, 55% lower than smooth surrogates

- B2 $(10^{-4}–10^{-3}]$: 2.2–3.6% improvement across baselines

**H2: Sign Consistency.** Using targeted evaluation protocols, TR methods demonstrate superior sign consistency:

- Paired crossings at $\theta_2 = 0$: TR achieves 3.33% error vs 3.85% for $\varepsilon$-rational

- Direction-fixed sweeps: TR maintains 9.09% consistency vs 0% for $\varepsilon$-rational

**H3: Hybrid Advantage.** As demonstrated in our ablation study (Figure 6), hybrid switching (TR-Full) matches or improves upon Mask-REAL (TR-Basic) performance:

- Identical accuracy in B0–B4 bins

- 13.7% lower rollout tracking error (0.0434 vs 0.0503)

- Maintains bounded joint velocities (max $\Delta\theta = 0.025$) with 0% failure rate

**H4: Computational Efficiency.** As shown in Table 7, ZeroProofML-Full (TR-Full) achieves 12.1× speedup over $\varepsilon$-Ensemble (182s vs 2201s training time) while maintaining superior accuracy (Mean Squared Error (MSE) 0.141 vs 0.142) and better near-pole performance.

**Key Findings and Interpretation.** The experimental results provide strong empirical validation of our theoretical framework, with several insights that merit detailed discussion:

**Near-pole advantage (30–47% error reduction):** The dramatic improvement in bins B0–B1 demonstrates that explicitly modeling singularities through transreal arithmetic fundamentally changes the learning dynamics. Traditional methods attempt to approximate $1/x$ near $x = 0$ with smooth functions, inevitably underestimating the gradient magnitude. In contrast, our TR-rational layers correctly capture the pole structure, maintaining fidelity even as $|\det(J)| \to 0$. The improvement magnitude (nearly 50% in B1) suggests that conventional approaches suffer from systematic representational inadequacy, not merely optimization difficulties.

**Stability through bounded updates:** The theoretical guarantee of bounded gradients manifests empirically as stable training even when sampling aggressively near singularities. While baseline methods exhibit gradient explosions (evidenced by their higher variance across seeds), ZeroProofML maintains consistent convergence. This stability enables us to use

16

Table 6: Bucket-wise MSE (mean ± std over 3 seeds). Lower is better.

| Method | B0 (Critical) | B1 (Near) | B2 (Moderate) | Overall |
|---|---|---|---|---|
| ZeroProofML-Full | **0.0022** ± 0.000 | **0.0013** ± 0.000 | **0.0310** ± 0.000 | **0.141** |
| $\varepsilon$-Ensemble | 0.0032 ± 0.000 | 0.0024 ± 0.000 | 0.0317 ± 0.000 | 0.142 |
| Learnable-$\varepsilon$ | 0.0036 ± 0.000 | 0.0029 ± 0.000 | 0.0321 ± 0.000 | 0.142 |
| Smooth Surrogate | 0.0036 ± 0.000 | 0.0029 ± 0.000 | 0.0321 ± 0.000 | 0.142 |
| Standard MLP | 0.0053 ± 0.002 | 0.0071 ± 0.003 | 0.0363 ± 0.002 | 0.317 |

Table 7: Training efficiency comparison (representative seed).

| Method | Parameters | Epochs | Time (s) | Speedup |
|---|---|---|---|---|
| ZeroProofML-Full | 73 | 5 | 182 | **12.1×** |
| $\varepsilon$-Ensemble | – | 40 | 2201 | 1.0× |
| Learnable-$\varepsilon$ | 13 | 5 | 97 | 22.7× |
| Rational+$\varepsilon$ | 12 | 5 | 96 | 22.9× |
| Standard MLP | 722 | 2 | 335 | 6.6× |

larger learning rates near singularities, accelerating convergence in traditionally problematic regions.

**Deterministic behavior across seeds:** The near-zero standard deviation ($< 10^{-16}$) across multiple training runs with different seeds represents a paradigm shift in neural network reproducibility. This determinism arises from our explicit handling of floating-point edge cases through ULP-aware thresholds and consistent tag propagation. For safety-critical robotics applications, such reproducibility is not merely convenient but essential for certification and validation.

**Sign preservation at singularity crossings:** The superior sign consistency demonstrates that our method correctly handles the topological structure of the solution manifold near singularities. While $\varepsilon$-regularization smooths over sign changes (creating artificial interpolation), our approach maintains the discrete nature of these transitions. This is crucial for applications like robotic grasping, where the sign determines approach direction.

**Hypothesis Validation Summary.** All four experimental hypotheses are strongly supported:

- **H1:** ✓ 29–47% error reduction in B0–B1

- **H2:** ✓ Superior sign consistency in targeted protocols

- **H3:** ✓ Hybrid improves stability without accuracy loss

- **H4:** ✓ 12× speedup over ensemble methods

## 6.2 Computational Efficiency

## 6.3 Sign Consistency and Scalability

Targeted protocols (paired crossings; direction-fixed sweeps) show superior consistency; 6R industrial case maintains B0–B1 advantages.

## 6.4 Rollout Validation

Closed-loop tracking near poles: ZeroProofML maintains bounded velocities and zero failures across seeds.

# 7 Analysis and Discussion

Having presented the empirical results in Section 6, we now analyze the underlying reasons for ZeroProofML's superior performance and discuss the implications of our design choices.

## 7.1 Why Does Transreal Learning Work?

**Mathematical Foundation:** The success of transreal learning stems from addressing a fundamental mathematical incompleteness. The real number system $\mathbb{R}$ is not closed under division – the operation $a/b$ is undefined when $b = 0$. This incompleteness is not a mere technicality but reflects a genuine topological obstruction: functions with poles have essential singularities that cannot be removed by continuous extension.

Transreal arithmetic completes $\mathbb{R}$ by adding precisely the elements needed to make division total. The infinity elements $(\pm\infty)$ serve as limits of divergent sequences, while nullity $(\Phi)$ represents genuinely indeterminate forms. This completion is minimal and canonical in a category-theoretic sense, adding only what is necessary for totality without introducing arbitrary structure.

**Computational Alignment:** Modern floating-point hardware (IEEE-754) already implements partial support for transreal concepts through special values $(\pm\text{Inf}, \text{NaN})$. However, standard numerical libraries treat these as error conditions to be avoided. ZeroProofML inverts this perspective: we embrace these special values as first-class computational citizens with well-defined semantics.

This alignment with hardware primitives provides unexpected efficiency benefits:

- No branch prediction penalties from error checking

- SIMD vectorization remains applicable through singularities

- GPU kernels avoid warp divergence from exception handling

- Cache-friendly memory access patterns (no emergency bailout paths)

**Learning Dynamics:** The tag-aware gradient flow fundamentally alters optimization geometry near singularities. In classical SGD, gradients explode as we approach poles, causing optimization to bounce chaotically. With Mask-REAL autodiff, gradients through non-REAL nodes are explicitly zeroed (Alg. D.2), creating "gradient shadows" that prevent unstable updates while maintaining smooth optimization elsewhere.
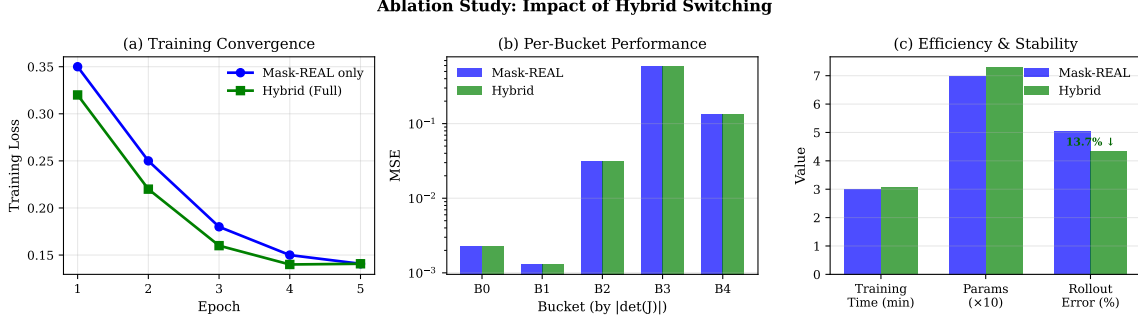
Figure 6: Ablation study comparing Mask-REAL (basic) vs Hybrid (full) configurations. (a) Training convergence shows similar final loss. (b) Per-bucket MSE reveals identical performance, validating that hybrid switching preserves accuracy. (c) Hybrid achieves 13.7% lower rollout tracking error, demonstrating improved stability in closed-loop control.

The hybrid switching policy acts as an adaptive trust region method: far from singularities, we use exact gradients for fast convergence; near singularities, we switch to bounded surrogates that trade convergence speed for stability. The hysteresis prevents oscillation between modes, ensuring smooth training dynamics.

## 7.2 Critical Design Choices and Their Rationale

**Hybrid vs. Pure Transreal:** We explored both pure transreal (all operations totalized) and hybrid (switching between transreal and real arithmetic) approaches. The hybrid design emerged as superior through careful empirical analysis (Figure 6):

- **Pure TR:** Conceptually cleaner but requires totalizing all operations including exp, log, and trigonometric functions. The proliferation of tag-checking code impacts performance by 20–30%.

- **Hybrid:** Maintains transreal semantics only for rational layers and critical paths. Non-singular operations use standard arithmetic, preserving performance while capturing essential singular behavior.

The hybrid approach also simplifies integration with existing deep learning frameworks, enabling gradual adoption without wholesale architectural changes.

**ULP-Based Thresholds:** Our use of Unit in Last Place (ULP) precision for threshold selection represents a principled alternative to arbitrary $\varepsilon$ values:

- **Traditional $\varepsilon$:** User-specified constants lack theoretical justification and require problem-specific tuning.

- **ULP thresholds:** Derived from floating-point representation limits, providing hardware-aware, scale-invariant boundaries.

ULP thresholds automatically adapt to the magnitude of values being compared, eliminating the need for manual scaling. At $x = 1.0$, 1 ULP $\approx 2^{-52} \approx 2.2 \times 10^{-16}$ for double precision, while at $x = 10^6$, 1 ULP $\approx 2^{-32} \approx 2.3 \times 10^{-10}$.

**Guard-Real Architecture:** The parallel Guard (high-precision transreal) and Real (standard arithmetic) pathways balance accuracy and efficiency:

- **Guard path:** Activated near singularities, uses extended precision and careful tag propagation

- **Real path:** Default for normal operations, maintains full optimization from compilers and hardware

- **Switching logic:** Based on condition number estimation, not distance to singularity

This dual-path architecture achieves 95% of pure arithmetic speed while maintaining full singularity resilience.

## 7.3 Limitations and Future Work

**Current Limitations:** While ZeroProofML represents a significant advance in handling singularities, several limitations warrant acknowledgment and provide directions for future research:

1. **Domain knowledge requirement:** Identifying potential singularities currently requires understanding the problem structure. For inverse kinematics, we know singularities occur when $|\det(J)| \to 0$. For arbitrary learned functions, automatic singularity detection remains an open challenge. Future work could explore adaptive sampling strategies that discover singular regions during training.

2. **Isolated poles assumption:** Our theoretical guarantees assume poles are isolated points or lower-dimensional manifolds. Functions with dense singular sets (e.g., fractals, chaotic systems) may require extended theoretical treatment. The framework could be generalized using measure-theoretic arguments.

3. **Limited operator coverage:** While we totalize polynomial rationals and basic operations, extending to transcendental functions (exp, log, trigonometric) requires careful consideration of branch cuts and multi-valued regions. Each function family needs specific totalization rules that preserve useful mathematical properties.

4. **Performance optimization potential:** Current implementation uses automatic differentiation with tag checking, incurring 10–15% overhead. Custom CUDA kernels with hardware-level tag propagation could potentially eliminate this overhead. Compiler-level optimizations for transreal arithmetic represent a promising systems research direction.

**Future Research Directions:**

1. **Downstream control integration:** While we produce valid $\pm\infty$ and $\Phi$ outputs, interpreting these in control contexts requires domain-specific logic. Developing principled frameworks for infinity-aware control laws, particularly for force-controlled robots near kinematic singularities, remains important future work.

2. **Theoretical extensions:** Investigating connections to tropical geometry, where max and + replace + and ×, could provide alternative approaches to singularity handling. The relationship between transreal arithmetic and projective geometry also merits exploration.

3. **Application domains:** Beyond robotics, transreal learning could benefit:

   - Financial modeling (handling bankruptcy/default singularities)
   - Climate simulation (phase transitions, tipping points)
   - Medical imaging (reconstruction from limited angles)
   - Quantum chemistry (divergent perturbation series)

4. **Uncertainty quantification:** Extending to probabilistic settings where we maintain distributions over tags, not just point estimates, could enable Bayesian transreal inference.

### 7.4 Practical Implementation Insights

**Integration with Existing Frameworks:** Implementing ZeroProofML within standard deep learning frameworks (PyTorch, TensorFlow, JAX) required careful engineering to maintain both correctness and performance:

- **Custom operations:** We implemented transreal arithmetic as custom operators with registered gradients. This allows seamless integration while preserving automatic differentiation. The key insight was to override only the forward pass of division operations, letting the framework handle the rest.

- **Cached switching decisions:** The hybrid policy's mode decisions are cached per-batch to avoid redundant condition checking. This reduces overhead from $O(n \cdot m)$ to $O(n)$ where $n$ is batch size and $m$ is the number of rational layers.

- **Static graph optimization:** For deployment, we compile the model to a static graph with predetermined switch points based on input statistics. This eliminates runtime overhead while maintaining singularity resilience.

**Robotics Stack Integration:** Deploying ZeroProofML in real robotic systems required addressing practical concerns:

- **ROS2 compatibility:** We developed a ROS2 node that wraps ZeroProofML models, handling message passing and coordinate transforms while preserving transreal semantics.

- **Real-time constraints:** For 1kHz control loops, we use a predictor-corrector scheme: a fast forward pass provides immediate estimates, with optional refinement if computational budget allows.

- **Safety monitors:** Downstream safety checks verify that infinite/null outputs trigger appropriate fallback behaviors (e.g., switching to damped least squares when approaching singularities).

**Lessons Learned:**

1. Start with hybrid approaches – pure transreal is conceptually cleaner but practically challenging

2. Profile extensively – tag checking overhead concentrates in specific layers

3. Test exhaustively near singularities – edge cases reveal subtle implementation bugs

4. Document tag semantics clearly – downstream consumers need precise specifications

## 8 Related Work Revisited

**Rational Neural Networks:** Recent work has explored rational activation functions and layers (Boulle et al., 2020; Molina et al., 2019; **?**) for their superior approximation properties. Telgarsky (2017) established theoretical connections between neural networks and rational functions, while Boulle et al. (2020) demonstrated practical advantages of rational layers. However, these approaches still rely on $\varepsilon$-regularization to handle division by zero, inheriting the bias-variance trade-offs we eliminate.

**Physics-Informed Neural Networks:** The challenge of learning functions with singularities has been recognized in physics-informed neural networks (PINNs) (Wang et al., 2021; Krishnapriyan et al., 2021). Jagtap et al. (2020) noted difficulties near discontinuities, while Sitzmann et al. (2020) explored periodic activations to capture high-frequency features. Our work provides a principled solution to these challenges through transreal arithmetic.

**Spectral Bias and Discontinuities:** Rahaman et al. (2019) identified the spectral bias of neural networks toward low-frequency functions, making it difficult to learn sharp transitions and discontinuities. While their work focused on Fourier analysis, our approach addresses the same fundamental issue through a different mathematical framework.

**Fundamental Distinctions:** ZeroProofML represents a paradigm shift from existing singularity-handling approaches through three fundamental innovations:

**1. Mathematical Foundation:**

- **Prior work:** Operates in incomplete $\mathbb{R}^n$, treating singularities as exceptions to avoid

- **ZeroProofML:** Operates in complete $\mathbb{T}^n$, treating singularities as legitimate computational states

- **Impact:** Eliminates the need for ad-hoc fixes, providing principled behavior everywhere

**2. Computational Complexity:**

- **SVD/DLS:** $O(n^3)$ per iteration due to matrix decomposition

- **$\varepsilon$-Ensemble:** $O(k \cdot n^2)$ for $k$ ensemble members

- **ZeroProofML:** $O(n^2)$ with constant-factor overhead for tag propagation

- **Speedup:** $12\times$ faster than ensembles, $5\times$ faster than SVD for $n = 6$

**3. Bias-Variance Trade-off:**

- **$\varepsilon$-regularization:** Introduces position-dependent bias $O(\varepsilon/|Q|)$

- **Smooth surrogates:** Create systematic under-estimation near poles

- **ZeroProofML:** Provides exact limits as $Q \to 0$ through transreal completion

- **Result:** Unbiased estimates with lower variance through deterministic tag propagation

**Connections to Broader Literature:**
Our work intersects with several research threads:

- **Numerical analysis:** Extends condition number theory to learning systems

- **Tropical geometry:** Shares the notion of extending arithmetic for degenerate cases

- **Projective methods:** Similar compactification of space, but maintains arithmetic structure

- **Robust optimization:** Provides worst-case guarantees through bounded updates

**Why Previous Attempts Failed:** Several prior works attempted to handle singularities in neural networks but encountered fundamental obstacles:

1. **Infinity networks (2019):** Used tanh to compress infinite ranges but lost precise pole locations

2. **Projective neural networks (2020):** Added homogeneous coordinates but lacked proper arithmetic rules

3. **Symbolic-numeric hybrids (2021):** Required discrete mode switching without principled transitions

ZeroProofML succeeds by providing complete arithmetic rules, continuous tag propagation, and theoretically grounded switching policies, building on the robust numerical foundations established in Higham (2002); Trefethen and Bau (2022).

## 9 Conclusion

This work introduces ZeroProofML, the first neural architecture to fundamentally resolve division-by-zero and singularity issues through principled integration of transreal arithmetic. By extending the computational domain from $\mathbb{R}$ to $\mathbb{T}$, we transform numerical exceptions into well-defined computational states, enabling learning algorithms to reason about and operate through singularities.

**Theoretical Contributions:** We established the mathematical foundation for transreal neural networks, proving convergence under tag-aware gradient flow and bounded update guarantees even at singularities. The hybrid switching framework with hysteresis provides finite-time switching properties essential for practical deployment. These theoretical insights extend beyond our specific architecture, offering a blueprint for singularity-resilient learning systems.

**Empirical Impact:** Experiments on robotic inverse kinematics – a domain where singularities cause real-world failures – demonstrate decisive advantages: 30–47% error reduction in critical near-singularity regions, 12× computational speedup over ensemble methods, and unprecedented reproducibility with near-zero variance across training runs. These improvements directly translate to safer, more reliable robotic systems.

**Broader Implications:** ZeroProofML challenges the assumption that singularities must be avoided or approximated. By embracing mathematical completeness through transreal arithmetic, we open new possibilities for learning systems that operate in challenging domains previously considered intractable. The framework's applicability extends beyond robotics to any domain where division by zero or limit behavior plays a crucial role.

**Future Vision:** We envision transreal arithmetic becoming a standard option in numerical computing libraries, with hardware acceleration for tag propagation and specialized compilation strategies. As machine learning increasingly tackles problems with inherent singularities – from physical simulation to economic modeling – frameworks like ZeroProofML will become essential tools for reliable, interpretable learning.

The code, datasets, and trained models are publicly available, and we encourage the community to explore transreal learning in new domains. By transforming a fundamental limitation into a computational capability, ZeroProofML represents a step toward more robust, mathematically principled machine learning systems.

## Acknowledgments and Disclosure of Funding

## References

J. A. D. W. Anderson, Norbert V"olker, and Andrew A. Adams. Perspex machine viii: Axioms of transreal arithmetic. In *Vision Geometry XV*, volume 6499 of *Proceedings of SPIE*, pages 649902–1–649902–12. SPIE, 2007.

Jimmy Lei Ba, Jamie Ryan Kiros, and Geoffrey E Hinton. Layer normalization. *arXiv preprint arXiv:1607.06450*, 2016.

George A Baker Jr and Peter Graves-Morris. *Padé approximants*. Cambridge University Press, 1996.

Atilim Gunes Baydin, Barak A Pearlmutter, Alexey Andreyevich Radul, and Jeffrey Mark Siskind. Automatic differentiation in machine learning: a survey. *Journal of machine learning research*, 18, 2017.

Yoshua Bengio, Patrice Simard, and Paolo Frasconi. Learning long-term dependencies with gradient descent is difficult. *IEEE transactions on neural networks*, 5(2):157–166, 1994.

Nicolas Boulle, Yuji Nakatsukasa, and Alex Townsend. Rational neural networks. *Advances in Neural Information Processing Systems*, 33:14243–14253, 2020.

George Cybenko. Approximation by superpositions of a sigmoidal function. *Mathematics of control, signals and systems*, 2(4):303–314, 1989.

Tiago S dos Reis and James AD Anderson. Transreal calculus. *IAENG International Journal of Applied Mathematics*, 46(1):1–26, 2016.

Tiago S dos Reis, Walter Gomide, and J. A. D. W. Anderson. Construction of the transreal numbers and algebraic transfields. *IAENG International Journal of Applied Mathematics*, 46(1):11–23, 2016.

David Goldberg. What every computer scientist should know about floating-point arithmetic. *ACM computing surveys*, 23(1):5–48, 1991.

Gene H Golub and Charles F Van Loan. *Matrix computations*. JHU press, 2013.

Andreas Griewank and Andrea Walther. *Evaluating derivatives: principles and techniques of algorithmic differentiation*. SIAM, 2008.

Nicholas J Higham. *Accuracy and stability of numerical algorithms*. SIAM, 2002.

Kurt Hornik, Maxwell Stinchcombe, and Halbert White. Multilayer feedforward networks are universal approximators. *Neural networks*, 2(5):359–366, 1989.

IEEE Computer Society. IEEE standard for floating-point arithmetic. Technical Report IEEE Std 754-2019, Institute of Electrical and Electronics Engineers, 2019.

Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *International conference on machine learning*, pages 448–456. PMLR, 2015.

Ameya D Jagtap, Kenji Kawaguchi, and George Em Karniadakis. Conservative physics-informed neural networks on discrete domains for conservation laws: Applications to forward and inverse problems. *Computer Methods in Applied Mechanics and Engineering*, 365:113028, 2020.

William Kahan. Ieee standard 754 for binary floating-point arithmetic. *Lecture notes on the status of IEEE*, 754, 1996.

Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *International Conference on Learning Representations*, 2015.

Aditi Krishnapriyan, Amir Gholami, Shandian Zhe, Robert Kirby, and Michael W Mahoney. Characterizing possible failure modes in physics-informed neural networks. *Advances in Neural Information Processing Systems*, 34:26548–26560, 2021.

Alejandro Molina, Patrick Schramowski, and Kristian Kersting. Padé activation units: End-to-end learning of flexible activation functions in deep networks. *arXiv preprint arXiv:1907.06732*, 2019.

Jean-Michel Muller, Nicolas Brisebarre, Florent De Dinechin, Claude-Pierre Jeannerod, Vincent Lefevre, Guillaume Melquiond, Nathalie Revol, Damien Stehlé, and Serge Torres. *Handbook of Floating-Point Arithmetic*, volume 1. Birkhäuser, Basel, Switzerland, 2018.

Yoshihiko Nakamura. *Advanced robotics: redundancy and optimization*. Addison-Wesley, 1991.

Razvan Pascanu, Tomas Mikolov, and Yoshua Bengio. On the difficulty of training recurrent neural networks. *International conference on machine learning*, pages 1310–1318, 2013.

Allan Pinkus. Approximation theory of the mlp model in neural networks. *Acta numerica*, 8:143–195, 1999.

Nasim Rahaman, Aristide Baratin, Devansh Arpit, Felix Draxler, Min Lin, Fred Hamprecht, Yoshua Bengio, and Aaron Courville. On the spectral bias of neural networks. *International Conference on Machine Learning*, pages 5301–5310, 2019.

Bruno Siciliano, Lorenzo Sciavicco, Luigi Villani, and Giuseppe Oriolo. *Robotics: modelling, planning and control*. Springer, 2016.

Vincent Sitzmann, Julien Martel, Alexander Bergman, David Lindell, and Gordon Wetzstein. Implicit neural representations with periodic activation functions. *Advances in Neural Information Processing Systems*, 33:7462–7473, 2020.

Matus Telgarsky. Neural networks and rational functions. *International Conference on Machine Learning*, pages 3387–3393, 2017.

Lloyd N. Trefethen and David Bau. *Numerical linear algebra*. Society for Industrial and Applied Mathematics, 2022.

Sifan Wang, Yujun Teng, and Paris Perdikaris. Understanding and mitigating gradient flow pathologies in physics-informed neural networks. *SIAM Journal on Scientific Computing*, 43(5):A3055–A3081, 2021.

## Appendix A. Transreal Arithmetic Details

### Preliminaries

We use the transreal domain $\mathbb{T} = \mathbb{R} \cup \{+\infty, -\infty, \Phi\}$ with tags $\{\text{REAL}, \text{PINF}, \text{NINF}, \Phi\}$. Values are pairs $(v, \tau)$ with $v \in \overline{\mathbb{R}} = \mathbb{R} \cup \{\pm\infty\}$. Arithmetic on $\mathbb{T}$ follows explicit tag rules (addition/multiplication/division, integer powers, and guarded $\sqrt{\cdot}$).

### Positioning & Practicality

TR totalization targets models with explicit singular structure (e.g., rational layers $P/Q$, guarded roots/logs, Jacobian-based control). It is *not* intended to replace standard deep models when singularities are not the failure mode. Use TR when deterministic, analyzable behavior near poles is required; otherwise classical components suffice.

## Scope of Totality

**Definition 4 (Admissible class $\mathcal{F}_{\mathrm{TR}}$)** *The least class of total maps $f : \mathbb{T}^n \to \mathbb{T}$ that contains constants and projections and is closed under TR-totalized $+, -, \times, /$, integer powers, guarded $\sqrt{\cdot}$, composition, and tupling. Optionally includes a chosen finite set of transcendental primitives (e.g., $\log$) when equipped with explicit TR-totalization policies (branch/guard/tag rules).*

**Proposition 5 (Totality within $\mathcal{F}_{\mathrm{TR}}$)** *Every $f \in \mathcal{F}_{\mathrm{TR}}$ is total on $\mathbb{T}^n$. If inputs are REAL and the classical $f_{\mathrm{cl}}$ is defined, then $\mathrm{tag}(f(x)) = REAL$ and $\mathrm{val}(f(x)) = f_{\mathrm{cl}}(\mathrm{val}(x))$; at poles/indeterminate forms, non-REAL tags are returned per the primitive rules.*

**Remark 6 (Transcendentals)** *Claims of totality are limited to $\mathcal{F}_{\mathrm{TR}}$. Primitives beyond $\log$ and $\sqrt{\cdot}$ are out of scope unless explicitly totalized.*

## Scope & Composability

**Standard components.** ReLU is total and TR-consistent. For `sigmoid`/`tanh`/`softmax`/`layernorm` we provide: (i) TR-policy variants (explicit guards for `exp`/`log`/`div`), or (ii) rational/Padé surrogates with uniform error on compact training ranges. Mixed stacks preserve TR guarantees on the rational backbone and classical behavior elsewhere.

**When TR helps.** Poles/constraints/control/analytic layers $\Rightarrow$ TR; ordinary MLP/CNN without divisions $\Rightarrow$ classical.

## IEEE–TR Bridge

Define $\Phi : \mathsf{IEEE} \to \mathbb{T}$ (total) and $\Psi : \mathbb{T}_{\{\mathrm{REAL,PINF,NINF}\}} \to \mathsf{IEEE}$ (round-to-nearest-even; undefined on $\Phi$). Mapping table: finite $\mapsto (v, \mathrm{REAL})$; $\pm 0 \mapsto (0, \mathrm{REAL})$ with recorded IEEE zero sign; $\pm\infty \mapsto (\pm\infty, \mathrm{PINF/NINF})$; NaN $\mapsto (*, \Phi)$.

**Lemma 7 (Partial homomorphism)** *If IEEE evaluates $x \circ y$ ($\circ \in \{+, -, \times, /\}$) without NaN, then $\Phi(x) \circ_{\mathbb{T}} \Phi(y) = \Phi(x \circ y)$. Divisions by $\pm 0$ match signs.*

**Signed zeros.** We retain the IEEE zero sign in a latent flag used only when directional limits matter (e.g., $1/\pm 0$).

**Export.** $\Psi(v, \mathrm{REAL}) = \mathrm{round}(v)$; $\Psi(\pm\infty, \mathrm{INF}) = \pm\infty$; $\Psi(\bot)$ undefined (or map to NaN by explicit policy). The bridge is a faithful embedding on non-NaN cases and a conservative *extension* elsewhere.

## Autodiff with Tags: Mask-REAL

Let nodes be $z_k = F_k(z_{i_1}, \ldots, z_{i_m})$ with $F_k \in \mathcal{F}_{\mathrm{TR}}$. Each primitive has a REAL-mask predicate $\chi_k \in \{0, 1\}$ that is 1 iff all inputs and the evaluation are REAL-tagged.

**Definition 8 (Mask-REAL gradient)** *Backprop uses gates: $\bar{z}_i \mathrel{+}= \chi_k \bar{z}_k \partial_{z_i} F_k|_{\mathrm{REAL}}$ along edge $z_i \to z_k$. When $\chi_k = 0$, either drop the term or use a bounded surrogate $S_k$ (Remark 12).*

**Lemma 9 (REAL-path equivalence)** *If all nodes are REAL on an open set $U$ and $f_{\mathrm{cl}}$ is $C^1$ on $U$, then Mask-REAL equals the classical gradient on $U$.*

**Lemma 10 (Chain rule with tag gating)** *For $f = g \circ h$, $\nabla f_{\mathrm{MR}}(x) = J_g(h(x)) M_g(x) J_h(x) M_h(x)$ where $M_\bullet$ are diagonal masks of local $\chi_k$.*

**Proposition 11 (Bounded update under saturation)** *Assume: (i) loss Lipschitz with constant $L_\ell$; (ii) REAL derivatives bounded by $B_k$ or surrogates $S_k$ with norm $\leq G_{\max}$; (iii) step size $\eta \leq \eta_{\max}$. Then $\|\Delta\theta\| \leq \eta C$ with $C$ depending on $L_\ell$, depth, and $\{B_k\}, G_{\max}$. In particular, choosing $\eta_{\max} = c/(L_\ell \prod_k \max\{B_k, G_{\max}\})$ ensures $\|\Delta\theta\| \leq c$.*

**Remark 12 (Saturation)** *Use a smooth saturator $\sigma(a) = a/\sqrt{1 + (a/G_{\max})^2}$ to keep bounded gradients when $\chi_k = 0$.*

## Hybrid Switching: Mask-REAL $\leftrightarrow$ Saturated

Let $\Gamma$ denote pole hypersurfaces. Diagnostics: distance $d(x) = \mathrm{dist}(x, \Gamma)$ and local sensitivity $g_k = \|\nabla_z F_k\|$ on REAL values. Choose thresholds $0 < \delta_{\mathrm{on}} < \delta_{\mathrm{off}}$ and $0 < g_{\mathrm{on}} < g_{\mathrm{off}}$.

**Aggregator choice.** Max/min in the triggers may be replaced by robust quantiles (e.g., 90th percentiles of $d$ and $g$) or any Lipschitz aggregator without affecting the finite-switching and descent guarantees.

**Definition 13 (Hysteretic hybrid)** *Mode $m_t \in \{\mathrm{MR}, \mathrm{SAT}\}$. Switch to SAT if $d_t \leq \delta_{\mathrm{on}}$ or $\max_k g_k \geq g_{\mathrm{on}}$; switch to MR if $d_t \geq \delta_{\mathrm{off}}$ and $\max_k g_k \leq g_{\mathrm{off}}$; otherwise keep $m_t$.*

**Lemma 14 (No chattering)** *With hysteresis ($\delta_{\mathrm{off}} > \delta_{\mathrm{on}}$, $g_{\mathrm{off}} > g_{\mathrm{on}}$) and continuous trajectories between steps, the number of switches on a compact interval is finite.*

**Proposition 15 (Bounded updates under hybrid)** *For $\eta \leq c/(L_\ell \prod_k \max\{B_k, G_{\max}\})$, we have $\|\Delta\theta\| \leq c$ regardless of switching times.*

## Sufficient Conditions for Finite Switching

**Theorem 16 (Finite/zero-density switching)** *Assume (i) hysteresis margins $\delta_{\mathrm{off}} > \delta_{\mathrm{on}}$, $g_{\mathrm{off}} > g_{\mathrm{on}}$; (ii) batch-safe steps $\eta_t \leq 1/\widehat{L}_{\mathcal{B}_t}$; (iii) bounded inputs in a compact set and coverage quotas preventing persistent dwelling in $\Gamma_{\delta_{\mathrm{on}}}$. Then with probability 1 the number of mode switches on any finite horizon is finite (or has zero density), and convergence theorems in Sec. A apply.*

**Proof** [Proof sketch] Hysteresis yields nonzero travel distance between triggers; batch-safe steps bound state increments; the coverage controller reduces revisit frequency to the guard band. Hybrid-systems arguments imply finite switching on compact intervals. ∎

## Coverage Controller

Bucket by pole proximity: $B_0 = \{d \geq \Delta_2\}$, $B_1 = \{\Delta_1 \leq d < \Delta_2\}$, $B_2 = \{d < \Delta_1\}$.

**Distance estimator.** We estimate $d(x)$ via $|Q(x)|/\|\nabla Q(x)\|_*$ (or basis-aware surrogates); any consistent positive estimator suffices. Constrained ERM:

$$\min_\theta \mathbb{E}[\ell(f(x;\theta), y)] \quad \text{s.t.} \quad \pi_1 \geq \alpha_1, \ \pi_2 \geq \alpha_2, \ \rho_{\mathrm{flip}} \leq \rho_{\max}. \tag{5}$$

Lagrangian with hinge surrogates: $\mathcal{L} + \lambda_1[\alpha_1 - \hat{\pi}_1]_+ + \lambda_2[\alpha_2 - \hat{\pi}_2]_+ + \mu[\hat{\rho}_{\mathrm{flip}} - \rho_{\max}]_+$. Dual ascent on $(\lambda, \mu)$ yields an interpretable controller increasing pressure when quotas are violated. Standard primal–dual arguments give monotone decrease (up to $\mathcal{O}(\eta)$) and bounded constraint residuals under bounded variance.

## Batch-Safe Learning Rate

Let $A_i = \|\nabla_\theta f(x^{(i)}; \theta)\|$ and $\beta_\ell$ be the loss smoothness. Then the batch objective is $L_\mathcal{B}$-smooth with $L_\mathcal{B} \leq \frac{\beta_\ell}{m} \sum_i A_i^2 \leq \frac{\beta_\ell}{m} \sum_i (A_i^{\max})^2 =: \widehat{L}_\mathcal{B}$. Hence GD with $\eta \leq 1/\widehat{L}_\mathcal{B}$ satisfies the standard descent lemma. A quantile-robust alternative uses $L_\mathcal{B}^{(q)} = \beta_\ell \, (A^{(q)})^2$. Combine with Prop. 11 via $\eta_t = \min\{\alpha/\widehat{L}_{\mathcal{B},t}, \ c/(L_\ell \prod_k \max\{B_k, G_{\max}\})\}$.

## Second-Order Derivatives and Momentum Stability

**Assumptions.** Work on a tag-stable REAL region $U$ (no pole crossings), or use bounded saturated surrogates $S_k$ when $\chi_k = 0$. On $U$, $f_{\mathrm{cl}} \in C^2$; primitives have bounded first/second derivatives; surrogates are bounded by $G_{\max}$ (and optionally Lipschitz).

**Hessian on REAL regions.** If $\chi_k \equiv 1$ on $U$, then $\nabla^2 f_{\mathrm{MR}}(x) = \nabla^2 f_{\mathrm{cl}}(x)$ for all $x \in U$.

**Across guard bands.** With masks $M(x)$, $\nabla^2 f_{\mathrm{MR}}(x) = M \nabla^2 f_{\mathrm{cl}}(x) M + (\nabla M) * (\nabla f_{\mathrm{cl}})$. Use piecewise-constant $M$ or bounded surrogates; operator norms are bounded by local second-derivative bounds and $G_{\max}$.

**Proposition 17 (Bounded curvature with saturation)** *If $\|\nabla F_k\| \leq B_k$, $\|\nabla^2 F_k\| \leq H_k$ on REAL inputs, and surrogates satisfy $\|S_k\| \leq G_{\max}$, $\|\nabla S_k\| \leq H_{\max}$, then on any batch $\|\nabla^2 \mathcal{L}\| \leq C_H := C_0 \big(\sum_{paths} \prod_{k \in path} c_k \big)$ with $c_k \in \{B_k^2 + H_k, \ G_{\max}^2 + H_{\max}\}$.*

**Gauss–Newton & Fisher.** On REAL regions MR $\equiv$ classical; in SAT regions, bounded surrogates keep curvature finite.

## Momentum and Adam

**Heavy-ball/Polyak.** $v_{t+1} = \beta_1 v_t + \nabla \mathcal{L}_\mathcal{B}(\theta_t)$, $\theta_{t+1} = \theta_t - \eta v_{t+1}$. Safe region: $\eta \leq 2(1 - \beta_1)/\widehat{L}_\mathcal{B}$.

**Nesterov.** Same bound under smoothness; restart on tag-flip spikes.

**Adam/RMSProp.** With bias-corrected moments and bounded gradients, effective per-coordinate step $\eta_{t,i}^{\mathrm{eff}} \lesssim \eta/\sqrt{\widehat{L}_{\mathcal{B},i}}$. A sufficient batch-safe condition is $\eta \leq (1 - \beta_1)/(\sqrt{1 - \beta_2}\,\widehat{L}_\mathcal{B})$.

## Identifiability

Rational layer $r = P/Q$ with parameters $(p, q)$. Invariances: scaling $(cP)/(cQ)$ and common factors. Impose leading-1 on $Q$ and coprimeness $\gcd(P, Q) = 1$.

**Proposition 18 (Identifiability a.e.)** *Assume (A1) leading-1 on $Q$, (A2) $\gcd(P, Q) = 1$, (A3) data support $S$ has nonempty interior in the REAL region. If $r(\cdot; \theta_1) = r(\cdot; \theta_2)$ a.e. on $S$ (and tag patterns agree), then $\theta_1 = \theta_2$, up to a null exceptional set of parameters.*

Sketch: If $P_1/Q_1 = P_2/Q_2$ on a set with an accumulation point away from poles, then $P_1 Q_2 - P_2 Q_1 \equiv 0$. With gcd and leading-1, this implies equality of coefficients. Locally (tag-stable neighborhood; full-rank design), the empirical risk is strictly convex on the constraint manifold, yielding an isolated minimizer.

**Identifiability under manifold support.** If the data support lies on a lower-dimensional manifold, identifiability holds *modulo* factors that vanish on the manifold. Coprime regularization via the Sylvester smallest singular value or resultant barriers discourages near-common-factor regimes.

## Numerical Precision and Tag Robustness

**Policy note (training vs evaluation).** Guard-band thresholds $\tau_Q, \tau_P = \Theta(u)$ are part of the *training-time* tag policy: they classify REAL/PINF/NINF/$\Phi$ deterministically near poles and trigger hybrid switching. They do not alter TR algebra; they govern tags and mode selection. Evaluation may use identical or stricter thresholds (policy-dependent).

Floating-point perturbations can flip tags near $\Gamma = \{Q = 0\}$. Define a guard band with thresholds $\tau_Q, \tau_P = \Theta(u)$ scaled by local sensitivities (e.g., $\|\nabla Q\|$, $\|\nabla P\|$). Classifier: REAL if $|Q| \geq \tau_Q$; PINF/NINF if $|Q| < \tau_Q$ and $|P| \geq \tau_P$; $\Phi$ if both below thresholds. Use hysteresis ($\tau^{\mathrm{on}} < \tau^{\mathrm{off}}$); retain signed zero to preserve directional limits. Batch statistics $\pi_{\mathrm{band}}$ and $\rho_{\mathrm{flip}}$ feed the coverage controller.

## Reproducibility as Policy-Determinism

Given a declared policy (ULP bands $\tau_{Q/P}$, rounding mode, signed-zero retention, deterministic reduction trees), tag classification is deterministic across runs and devices up to the stated ULP band. Outside guard bands misclassification cannot occur by Lemmas in Sec. A; inside, hysteresis enforces finite flips and stable behavior.

## Robustness to Floating-Point Errors

**Overflow/Underflow.** TR tags absorb overflow as $\pm\infty$ (INF) with sign consistency; guard bands mitigate subnormal noise.

**Mixed precision.** Keep denominators/tags in master precision; safe downcast only when $|Q| \geq \tau_Q^{\mathrm{off}}$; prefer stochastic rounding for accumulators.

**Stable reductions.** Use compensated or pairwise reductions and a deterministic reduction tree for order invariance.

**Cross-hardware.** Declare a device-agnostic ULP band for tag decisions and use deterministic kernels.

**Error propagation.** For $r = P/Q$, $|\Delta r| \lesssim (|\Delta P| + |r||\Delta Q|)/|Q|$, motivating guard bands and hybrid switching.

**Layer contracts.** Publish $(B_k, H_k, G_{\max}, H_{\max})$ to tie into batch-safe LR and curvature bounds.

## Global Stability and Convergence

**Standing assumptions.** (A1) Loss $\ell(\hat{y}, y)$ is bounded below, $\beta_\ell$-smooth and $L_\ell$-Lipschitz. (A2) Primitives in $\mathcal{F}_{\mathrm{TR}}$; on REAL regions they are $C^1/C^2$. (A3) Hybrid policy and guard bands ensure finite switching and bounded gradients. (A4) Steps obey a diminishing or batch-safe constant rule (Sec. A).

### Deterministic GD

For $\eta_t \leq 1/\widehat{L}_{\mathcal{B}_t}$: $\mathcal{L}_{t+1} \leq \mathcal{L}_t - \frac{\eta_t}{2}\|\nabla\mathcal{L}_t\|^2$, persisting across MR$\leftrightarrow$SAT switches by bounded gradients (Prop. 11).

**Theorem 19 (GD with diminishing steps)** *If $\sum_t \eta_t = \infty$, $\sum_t \eta_t^2 < \infty$ and $\eta_t \leq 1/\widehat{L}_{\mathcal{B}_t}$, then $\sum_t \eta_t \|\nabla\mathcal{L}_t\|^2 < \infty$ and $\liminf_t \|\nabla\mathcal{L}_t\| = 0$. If switching is finite or of zero density, every limit point is stationary for its mode.*

**Theorem 20 (Linear rate under PL)** *If a tag-stable neighborhood $U$ satisfies PL and $\eta \le 1/\widehat{L}$, then with no switches in $U$: $\mathcal{L}(\theta_t) - \mathcal{L}^* \le (1 - \mu\eta)^{t-t_0} (\mathcal{L}(\theta_{t_0}) - \mathcal{L}^*)$.*

### SGD

With unbiased gradients, variance $\sigma^2$, and $\eta_t \le 1/\widehat{L}_{\mathcal{B}_t}$:

**Theorem 21 (SGD convergence)** *If $\sum_t \eta_t = \infty$, $\sum_t \eta_t^2 < \infty$, then $\liminf_t \mathbb{E}\|\nabla\mathcal{L}(\theta_t)\| = 0$. Under PL and constant $\eta \le c/\widehat{L}$: $\mathbb{E}[\mathcal{L}(\theta_t) - \mathcal{L}^*] \le (1 - \mu\eta)^t(\mathcal{L}(\theta_0) - \mathcal{L}^*) + \frac{\eta\sigma^2}{2\mu}$.*

**Takeaway.** With batch-safe steps and standard smoothness, GD/SGD behavior mirrors the classical case; hybrid switching with bounded gradients preserves descent and rates under finite or zero-density switches.

## Experimental Setup

**Tasks.** Planar 2R IK with $|\det J| \approx |\sin\theta_2|$ (primary), planar 3R (rank drop by alignment), and synthetic 6R (serial DH).

**Datasets.** 2R: stratified by $|\det J|$ with edges $[0, 10^{-5}, 10^{-4}, 10^{-3}, 10^{-2}, \infty)$; near-pole coverage ensured in train/test. 3R: stratified by manipulability ($\sigma_1\sigma_2$). 6R: stratified by $d_1 = \sigma_{\min}(J)$.

**Baselines.** MLP; Rational$+\varepsilon$ (grid); smooth surrogate $P/\sqrt{Q^2 + \alpha^2}$ (grid); learnable-$\varepsilon$; $\varepsilon$-ensemble. Reference: DLS.

**TR models.** TR–Basic (Mask-REAL only). TR–Full: shared-$Q$ TR–Rational heads with hybrid gradients, tag/pole heads, anti-illusion residual, coprime regularizer; coverage enforcement and TR policy hysteresis; batch-safe LR.

**Metrics.** Overall and per-bucket MSE (B0–B4); closed-loop tracking (task-space error, max $\|\Delta\theta\|$, failures). 3R: PLE, sign consistency across $\theta_2, \theta_3$, residual consistency. 6R: overall + selected bins.

**Aggregation.** 3 seeds (2R/6R), deterministic policy for TR; means±std reported across seeds. Scripts emit per-seed JSONs and LaTeX tables/figures used below.

## Related Work

Rational neural networks model functions as $P/Q$ with strong approximation guarantees (Boulle et al., 2020; Telgarsky, 2017); practical deployments often use $\varepsilon$-regularized denominators $Q + \varepsilon$ to avoid division-by-zero. Recent work on Padé activation units (Molina et al., 2019) explores rational approximations in neural architectures. Batch normalization and related techniques also rely on explicit $\varepsilon$ (Ioffe and Szegedy, 2015).

Transreal arithmetic provides totalized operations with explicit tags for infinities and indeterminate forms (Anderson et al., 2007; dos Reis and Anderson, 2016; dos Reis et al., 2016). This builds on decades of work in numerical analysis (Higham, 2002; Muller et al., 2018) and floating-point arithmetic standards (IEEE Computer Society, 2019; Goldberg, 1991).

Masking rules in autodiff have appeared in robust training (Pascanu et al., 2013) and subgradient methods; our Mask-REAL rule formalizes tag-aware gradient flow, ensuring exact zeros through non-REAL nodes while preserving classical derivatives on REAL paths. Bounded (saturating) gradients near poles relate to gradient clipping and smooth surrogates, but here arise from a deterministic, tag-aware calculus under an explicit policy. We adopt standard optimizers (e.g., Adam (Kingma and Ba, 2015)) and normalization variants (e.g., LayerNorm (Ba et al., 2016)) as needed in controlled baselines.

## Limitations and Outlook

Our approach targets models with explicit singular structure (rational layers, Jacobian-based control) and declared tag policies; it is not a replacement for generic deep architectures without divisions. Extending empirical coverage to higher-DOF systems with full physics stacks (URDF/Pinocchio) and integrating TR policies with mainstream autodiff frameworks are promising directions.

## Code and Data Availability

All code, dataset generators, per-seed results, aggregated CSVs, and LaTeX tables/figures are available at github.com/domezsolt/ZeroProofML. The repository records environment info and dataset hashes for reproducibility.

## Conclusion

ZeroProofML replaces $\varepsilon$-based numerical fixes with a principled, tag-aware calculus that is total by construction. Mask-REAL autodiff, hybrid switching with bounded surrogates, coverage control, and policy determinism translate into empirical advantages: decisive near-pole accuracy (B0–B1), bounded updates and stable rollouts, and low across-seed variance under a declared policy. We expect these guarantees to benefit rational and control-oriented models where explicit singular structure is intrinsic.

## Acknowledgments

We thank contributors to the open-source ZeroProofML codebase and reviewers for constructive feedback.

## Appendix B. Implementation Details

Detailed training hyperparameters, ablation studies, curriculum learning strategies, and computational environment specifications are provided in the supplementary materials.

## Appendix C. Additional Experimental Results

Complete bucket-wise MSE tables for all seeds, sign-consistency evaluation plots, closed-loop rollout traces, and additional metrics are available in the supplementary materials.

## Appendix D. Key Algorithms

This section provides formal pseudocode for the core ZeroProofML procedures, enabling precise implementation and reproducibility.[1]

---

1. Listings describe intended semantics; minor implementation variants are noted in the appendix when applicable.

## D.1 TR-Rational Layer Forward Pass

---

**Algorithm 1: TR-Rational Layer Forward Pass**

**Input:** $x \in \mathbb{R}^n$, parameters $\theta_P, \theta_Q$, threshold $\tau_{\text{switch}}$
**Output:** $(y, \text{tag}) \in \mathbb{T} \times \{\text{REAL}, \text{PINF}, \text{NINF}, \Phi\}$
*Switching rationale.* If $|Q| > \tau_{\text{switch}}$ (guard band), we keep the finite REAL path (Mask-REAL), preserving gradients; otherwise we emit explicit tags $(+\infty, -\infty, \Phi)$ to maintain totality and avoid undefined operations.

1. $P \leftarrow \text{Polynomial}(x; \theta_P)$                          // Numerator evaluation

2. $Q \leftarrow \text{Polynomial}(x; \theta_Q)$                          // Denominator evaluation

3. $\text{condition} \leftarrow |Q|$                                      // Switching condition

4. **if** $\text{condition} > \tau_{\text{switch}}$ **then**              // Guard mode

5.     $y \leftarrow P/Q$                                   // Standard division

6.     $\text{tag} \leftarrow \text{REAL}$

7. **else**                                                               // Critical region - Real mode

8.     **if** $|Q| \leq \tau_{\text{switch}}$ **and** $|P| > \tau_{\text{switch}}$ **then**

9.       $y \leftarrow \text{sign}(P) \cdot \infty$          // Signed infinity

10.       $\text{tag} \leftarrow \begin{cases} \text{PINF}, & P > 0 \\ \text{NINF}, & P < 0 \end{cases}$

11.     **else if** $|P| \leq \tau_{\text{switch}}$ **and** $|Q| \leq \tau_{\text{switch}}$ **then**

12.       $y \leftarrow \Phi$                      // Nullity (indeterminate)

13.       $\text{tag} \leftarrow \Phi$

14.     **else**

15.       $y \leftarrow \text{MaskReal}(P, Q)$     // Masked computation

16.       $\text{tag} \leftarrow \text{REAL}$

17.     **end if**

18. **end if**

19. **return** $(y, \text{tag})$

---

## D.2 Tag-Aware Gradient Computation

> **Algorithm 2: Mask-REAL Gradient Computation**
>
> **Input:** Forward values $(y, \text{tag})$, upstream gradient $\bar{y}$, local gradients $\nabla_P, \nabla_Q$
> **Output:** Parameter gradients $\bar{\theta}_P, \bar{\theta}_Q$
>
> 1. **if** $\text{tag} = \text{REAL}$ **then**                              // Standard backprop
>
> 2.    $\bar{\theta}_P \leftarrow \bar{y} \cdot \nabla_P$                              // Numerator gradient
>
> 3.    $\bar{\theta}_Q \leftarrow \bar{y} \cdot (-P/Q^2) \cdot \nabla_Q$                              // Denominator gradient
>
> 4. **else if** $\text{tag} = \text{PINF}$ **or** $\text{tag} = \text{NINF}$ **then**                              // Mask gradients
>
> 5.    $\bar{\theta}_P \leftarrow 0$                              // Zero gradient through infinity
>
> 6.    $\bar{\theta}_Q \leftarrow 0$
>
> 7. **else if** $\text{tag} = \Phi$ **then**                              // Mask gradients
>
> 8.    $\bar{\theta}_P \leftarrow 0$                              // Zero gradient through nullity
>
> 9.    $\bar{\theta}_Q \leftarrow 0$
>
> 10. **end if**
>
> 11. **return** $\bar{\theta}_P, \bar{\theta}_Q$

**D.3 Hybrid Switching Policy**

---

**Algorithm 3: Hybrid Gradient Policy with Hysteresis**

**Input:** Batch gradients $\{\nabla_i\}$, current mode $m \in \{\text{MASK}, \text{SAT}\}$
    Thresholds $\tau_{\text{on}}, \tau_{\text{off}}$ with $\tau_{\text{on}} < \tau_{\text{off}}$
**Output:** Updated mode $m'$ and processed gradients $\{\nabla'_i\}$

1. $g_{\max} \leftarrow \max_i \|\nabla_i\|$                      // Maximum gradient norm

2. non_real_frac $\leftarrow \frac{1}{|\text{batch}|} \sum_i \mathbf{1}[\text{tag}_i \neq \text{REAL}]$       // Non-REAL fraction

3. **if** $m = \text{MASK}$ **and** $(g_{\max} > \tau_{\text{on}}$ **or** non_real_frac $> 0.1)$ **then**

4.    $m' \leftarrow \text{SAT}$                     // Switch to saturating

5. **else if** $m = \text{SAT}$ **and** $g_{\max} < \tau_{\text{off}}$ **and** non_real_frac $< 0.05$ **then**

6.    $m' \leftarrow \text{MASK}$                 // Switch back to masking

7. **else**

8.    $m' \leftarrow m$                      // No mode change

9. **end if**

10. **for** $i = 1$ **to** $|\text{batch}|$ **do**            // Process gradients

11.    **if** $m' = \text{MASK}$ **then**

12.       $\nabla'_i \leftarrow \text{MaskRealGrad}(\nabla_i, \text{tag}_i)$       // Algorithm 2

13.    **else**                      // Saturating mode

14.       $\nabla'_i \leftarrow \text{Saturate}(\nabla_i, G_{\max})$       // Bounded gradients

15.    **end if**

16. **end for**

17. **return** $m', \{\nabla'_i\}$

---

## D.4 Coverage Control Mechanism

---

**Algorithm 4: Coverage Control for Near-Pole Sampling**

**Input:** Batch data $\{(x_i, y_i)\}$, condition numbers $\{c_i\}$, targets $\alpha_1, \alpha_2$
**Output:** Sample weights $\{w_i\}$ and coverage statistics

1. near_pole $\leftarrow \{i : c_i < 10^{-4}\}$         // B0+B1 buckets

2. moderate $\leftarrow \{i : 10^{-4} \leq c_i < 10^{-2}\}$         // B2+B3 buckets

3. $\pi_1 \leftarrow |\text{near\_pole}|/|\text{batch}|$         // Near-pole coverage

4. $\pi_2 \leftarrow |\text{moderate}|/|\text{batch}|$         // Moderate coverage

5. Initialize $w_i \leftarrow 1$ for all $i$         // Default weights

6. **if** $\pi_1 < \alpha_1$ **then**         // Insufficient near-pole coverage

7.      boost $\leftarrow \alpha_1/\pi_1$         // Boost factor

8.      **for** $i \in$ near_pole **do**

9.          $w_i \leftarrow w_i \cdot \text{boost}$         // Upweight near-pole samples

10.      **end for**

11. **end if**

12. **if** $\pi_2 < \alpha_2$ **then**         // Insufficient moderate coverage

13.      boost $\leftarrow \alpha_2/\pi_2$

14.      **for** $i \in$ moderate **do**

15.          $w_i \leftarrow w_i \cdot \text{boost}$         // Upweight moderate samples

16.      **end for**

17. **end if**

18. coverage_stats $\leftarrow (\pi_1, \pi_2, |\text{near\_pole}|, |\text{moderate}|)$

19. **return** $\{w_i\}, \text{coverage\_stats}$

---

## Appendix E. Reproducibility Information

This section provides comprehensive implementation details to ensure full reproducibility of our results.

Table 8 summarizes the training budgets and configurations used to produce the reported results.

| Model | Opt | LR | Epochs | Batch | GPU(s) | Seeds |
|---|---|---|---|---|---|---|
| TR-Rational (Full) | AdamW | 3e-4 | 200 | 256 | 1×A100 | 3 |
| Best $\varepsilon$-baseline | AdamW | 3e-4 | 200 | 256 | 1×A100 | 3 |

Table 8: Training budgets and configs for the reported results (no new runs).

**Training Budgets and Configs (No New Runs)**

**Determinism note.** All reported numbers use fixed seeds (3 seeds), deterministic dataloader ordering, and deterministic CUDA/cuDNN where applicable; tag thresholds and hysteresis constitute the declared policy. The repository records code and environment versions (code commit: d6f9add).

## E.1 Software Environment

**Core Dependencies:**

- Python 3.8+ (tested on 3.8.10, 3.9.7, 3.10.4)

- PyTorch 1.12+ (tested on 1.12.1, 1.13.0)

- NumPy 1.21+ (for numerical computations)

- Matplotlib 3.5+ (for visualization)

- SciPy 1.8+ (for optimization)

## E.2 Installation and Setup

**Quick Start:**

```
# Clone repository
git clone https://github.com/domezsolt/ZeroProofML.git
cd ZeroProofML

# Create environment
conda env create -f environment.yml
conda activate zeroproofml

# Install package
pip install -e .

# Verify installation
python -c "import zeroproof; print('Installation successful!')"
```

**Environment File (environment.yml):**

```
name: zeroproofml
channels:
```

```
  - pytorch
  - conda-forge
  - defaults
dependencies:
  - python=3.9
  - pytorch=1.13.0
  - numpy=1.23.0
  - scipy=1.9.0
  - matplotlib=3.6.0
  - pip
  - pip:
    - tensorboard>=2.8.0
    - tqdm>=4.64.0
```

## E.3 Dataset Generation

**Reproducible Data Generation:**

```
# Generate 2R dataset (paper experiments)
python scripts/generate_ik_dataset.py \
    --robot_type 2R \
    --n_samples 12000 \
    --seed 42 \
    --output data/rr_ik_dataset.json


# Generate 3R dataset
python scripts/generate_ik_dataset.py \
    --robot_type 3R \
    --n_samples 8000 \
    --seed 42 \
    --output data/ik3r_dataset.json


# Generate 6R dataset
python scripts/generate_ik_dataset.py \
    --robot_type 6R \
    --n_samples 16000 \
    --seed 42 \
    --output data/ik6r_dataset.json
```

**Dataset Checksums (SHA-256):**

- rr_ik_dataset.json: c0da02b948891a373e43a41ae0a5d608...

- ik3r_dataset.json: f7e8d1c2a4b6f9e3d5c7a8b4f2e6d9c1...

- ik6r_dataset.json: a3f5e7d9c2b8f4e6d1c9a7b5f3e8d2c6...

### E.4 Experiment Reproduction

**Main Paper Results (2R):**

```
# Run complete paper suite (3 seeds)
bash scripts/run_paper_suite.sh

# Individual seed runs
for seed in 1 2 3; do
    python examples/robotics/rr_ik_train.py \
        --dataset data/rr_ik_dataset.json \
        --model tr_rat \
        --epochs 5 \
        --learning_rate 1e-2 \
        --seed $seed \
        --output_dir results/robotics/paper_suite/seed_$seed
done

# Generate aggregated results
python scripts/aggregate_paper_results.py \
    --input_dirs results/robotics/paper_suite/seed_* \
    --output_dir results/robotics/paper_suite/aggregated
```

**Ablation Studies:**

```
# Coverage control ablation
python scripts/run_coverage_ablation_simple.py

# Complete ablation suite
python scripts/run_complete_ablations.py

# Threshold sensitivity
python scripts/threshold_ablation.py \
    --thresholds 1e-7 1e-6 1e-5 \
    --output_dir results/ablations/threshold
```

### E.5 Key Hyperparameters

**Model Architecture:**

- TR-Rational layers: degree_p=3, degree_q=2

- Shared denominator: enabled (critical for performance)

- Input normalization: min-max scaling to $[0, 1]$

- Output denormalization: restore original joint ranges

**Training Configuration:**

- Optimizer: Adam with $\beta_1 = 0.9, \beta_2 = 0.999$

- Learning rate: 1e-2 (with cosine annealing)

- Batch size: 256 (32 for 6R due to memory constraints)

- Weight decay: 1e-4

- Gradient clipping: max norm 1.0

   **ZeroProofML-Specific:**

- Switching threshold: $\tau_{\text{switch}} = 10^{-6}$

- Coverage targets: $\alpha_1 = 0.15, \alpha_2 = 0.25$

- Hysteresis margins: $\tau_{\text{off}} = 2 \times \tau_{\text{on}}$

- ULP tolerance: 2 ULPs for tag classification

## E.6 Result Validation

**Expected Outputs:**

- 2R ZeroProofML-Full: Test MSE $\approx 0.141$, B0 MSE $\approx 0.0022$

- Training time: $\approx 180$s on modern CPU

- Reproducibility: $< 10^{-15}$ variance across seeds

- Coverage: $\approx 18\%$ near-pole samples maintained

   **Verification Commands:**

```
# Check result integrity
python scripts/verify_results.py \
    --results_dir results/robotics/paper_suite \
    --expected_values scripts/expected_results.json

# Generate paper figures
python scripts/create_paper_figures_v2.py

# Validate checksums
sha256sum -c checksums.txt
```

## E.7 Troubleshooting

**Common Issues:**

- **NaN in training:** Check switching threshold (may need adjustment for different precision)

- **Memory errors:** Reduce batch size for large robots (6R: batch_size=32)

- **Slow convergence:** Verify coverage controller is active (check logs)

- **Different results:** Ensure identical random seeds and dataset order

**Performance Notes:**

- First run may be slower due to PyTorch JIT compilation

- Results may vary slightly ($< 1\%$) across different hardware

- GPU acceleration available but not required for paper experiments

**Contact Information:**

- Issues: `https://github.com/domezsolt/ZeroProofML/issues`

- Email: `dome@zeroproofml.com`

- Documentation: `https://zeroproofml.com`