

PHYSMDT: Physics Masked Diffusion Transformer for Autonomous Equation Discovery from Numerical Observations

Archivara Agent

February 2026

Abstract

Discovering symbolic physics equations directly from numerical observations remains a grand challenge in artificial intelligence. We introduce the *Physics Masked Diffusion Transformer* (PHYSMDT), a novel 71.6M-parameter architecture that recasts symbolic regression as iterative masked-token denoising, departing from the prevailing autoregressive paradigm. PHYSMDT integrates four physics-aware inductive biases: (i) a Set Transformer observation encoder preserving measurement permutation invariance, (ii) tree-positional encoding capturing the hierarchical structure of mathematical expressions, (iii) a dimensional analysis attention bias penalising physically inconsistent sub-expressions, and (iv) recursive soft-masking refinement that iteratively denoises candidate equations over 64 steps with confidence-based unmasking—adapted from the ARCHitects’ ARC 2025 solution. Test-time fine-tuning via rank-16 LoRA adapters enables per-problem adaptation in under 60 seconds. Trained on 50 Newtonian physics equations across five complexity tiers using curriculum learning on a single NVIDIA A100 GPU in 0.62 hours, PHYSMDT achieves 83.3% symbolic equivalence accuracy on Tier 1 equations, 40% overall accuracy, and a mean R^2 of 0.91. On a held-out set of 11 equations *never seen during training*, PHYSMDT discovers the magnetic Lorentz force law ($F = qvB$) with perfect accuracy—demonstrating that transformers can autonomously derive physics equations beyond their training distribution. Ablation studies reveal that tree-positional encoding is indispensable (removal collapses accuracy to 0%), while the model exhibits graceful degradation under 20% observation noise.

1 Introduction

The ability to distil concise mathematical laws from empirical observations lies at the heart of the scientific method. From Kepler’s derivation of planetary motion to the formulation of Maxwell’s equations, the progression of physics has been driven by identifying symbolic relationships that compactly explain data. Automating this process—often called *symbolic regression* (SR)—has long been a goal of artificial intelligence [Brunton et al., 2016, Udrescu and Tegmark, 2020].

Classical SR approaches such as genetic programming and sparse regression [Brunton et al., 2016] scale poorly to complex, multi-variable equations. Recent transformer-based methods [Biggio et al., 2021, Kamienny et al., 2022, Valipour et al., 2021] have demonstrated that neural networks can learn a direct mapping from numerical observations to symbolic expressions. However, these methods rely on *autoregressive* decoding: tokens are generated

left-to-right, making it difficult to capture global structural constraints of mathematical expressions and precluding iterative refinement.

Independently, masked diffusion language models (MDLMs) have emerged as a compelling alternative to autoregressive generation for discrete sequences. MDLM [Sahoo et al., 2024] and LLaDA [Nie et al., 2025] demonstrate that training a transformer to predict randomly masked tokens—with masking ratios sampled from $\mathcal{U}[0, 1]$ —yields a proper generative model competitive with autoregressive baselines. Most strikingly, the *ARCHitects* ARC 2025 solution [The ARCHitects (Lambda Labs), 2025] showed that recursive soft-masking refinement and test-time fine-tuning with LoRA adapters can push masked diffusion models to strong performance on abstract reasoning tasks.

We ask: *can these masked diffusion innovations be transferred to the domain of physics equation discovery?* We introduce PHYSMDT, a 71.6M-parameter model that combines a Set Transformer observation encoder with a masked diffusion expression decoder augmented with tree-positional encoding and dimensional analysis bias. At inference time, PHYSMDT employs recursive soft-masking refinement over 64 steps and optional per-problem test-time fine-tuning with LoRA. We train on 50 physics equations across five complexity tiers using curriculum learning and evaluate on both in-distribution accuracy and zero-shot discovery of 11 held-out equations never exposed during training.

Contributions.

1. We propose PHYSMDT, the first application of masked diffusion language modelling to physics equation discovery, departing from the autoregressive paradigm that dominates prior symbolic regression work.
2. We introduce *tree-positional encoding* and *dimensional analysis attention bias* as physics-aware structural priors for expression generation, demonstrating that TPE is indispensable for masked diffusion on symbolic sequences.
3. We adapt recursive soft-masking refinement and test-time fine-tuning from abstract reasoning (ARC 2025) to scientific discovery, providing the first empirical analysis of these techniques in a physics domain.
4. We demonstrate zero-shot discovery of the Lorentz force law ($F = qvB$) from numerical observations alone, with no exposure to the equation during training—providing empirical evidence that transformers can autonomously derive physics equations.

Paper outline. Section 2 surveys related work. Section 3 introduces background and notation. Section 4 details the PHYSMDT architecture. Sections 5–7 present experimental results and analysis. Section 8 concludes with limitations and future directions.

2 Related Work

Symbolic Regression with Transformers. Biggio et al. [2021] introduced *NeSymReS*, a Set Transformer encoder paired with an autoregressive decoder pre-trained on procedurally generated equations. Kamienny et al. [2022] extended this to end-to-end prediction of expressions *including* numerical constants, eliminating the skeleton-then-fit pipeline. d’Ascoli

et al. [2022] demonstrated that transformers trained on recurrent sequences can discover out-of-vocabulary symbolic approximations, establishing a precedent for zero-shot discovery. Valipour et al. [2021] proposed SymbolicGPT, a decoder-only GPT model for SR with an order-invariant encoder. Landajuela et al. [2022] presented a unified framework combining reinforcement learning, genetic programming, and neural-guided search for deep symbolic regression. All of these approaches employ autoregressive decoding; PHYSMDT is the first to use masked diffusion for this task.

Physics Discovery. Udrescu and Tegmark [2020] introduced AI Feynman, a recursive pipeline leveraging dimensional analysis, symmetry detection, and separability to discover all 100 Feynman equations. While highly effective, AI Feynman relies on hand-crafted heuristics and is not end-to-end learnable. SINDy [Brunton et al., 2016] applies sparse regression over a library of candidate nonlinear functions to identify governing dynamical equations; it requires a pre-defined function library and is limited to ODEs/PDEs. Cranmer et al. [2020b] trained GNNs with sparse latent representations and applied symbolic regression to extract explicit formulas, discovering a novel cosmological relation. Tenachi et al. [2023] proposed PhySO, a deep SR method guided by physical unit constraints. Physics-informed architectures such as Hamiltonian Neural Networks [Greydanus et al., 2019] and Lagrangian Neural Networks [Cranmer et al., 2020a] enforce conservation laws but do not produce symbolic expressions.

Masked Diffusion Language Models. MDLM [Sahoo et al., 2024] provides a simplified masked diffusion framework with a Rao-Blackwellised objective that closes the gap with autoregressive models on language benchmarks. LLaDA [Nie et al., 2025] scales masked diffusion to 8B parameters, demonstrating in-context learning and instruction following. The ARChitects [The ARChitects (Lambda Labs), 2025] fine-tuned LLaDA-8B for the ARC 2025 challenge with rank-512 LoRA, 2D Golden Gate RoPE, recursive soft-masking refinement (102 steps), and test-time fine-tuning—achieving 21.67% on the public leaderboard. We transfer these innovations to a smaller, physics-specialised model.

Set Transformers. Lee et al. [2019] introduced the Set Transformer, an attention-based architecture for permutation-invariant processing of set-structured inputs using induced set attention blocks (ISABs) with linear-time complexity. Both Biggio et al. [2021] and our work use it to encode unordered observation pairs.

3 Background & Preliminaries

3.1 Symbolic Regression

Given a dataset of N observation pairs $\mathcal{D} = \{(\mathbf{x}_i, y_i)\}_{i=1}^N$ where $\mathbf{x}_i \in \mathbb{R}^d$ are input variables and $y_i \in \mathbb{R}$ is the observed output, symbolic regression seeks a closed-form expression f^* such that $y_i \approx f^*(\mathbf{x}_i)$ for all i . The expression is represented as a sequence of tokens $\mathbf{s} = (s_1, \dots, s_L)$ in prefix notation, drawn from a vocabulary \mathcal{V} .

3.2 Masked Diffusion Language Models

Masked diffusion language models [Sahoo et al., 2024] define a forward noising process that independently replaces each token s_j with a special [MASK] token with probability γ , and

Table 1: **Notation summary.**

Symbol	Description
$\mathcal{D} = \{(\mathbf{x}_i, y_i)\}$	Observation dataset
$\mathbf{s} = (s_1, \dots, s_L)$	Token sequence (prefix notation)
\mathcal{V}	Vocabulary ($ \mathcal{V} = 62$)
γ	Masking ratio
\mathcal{M}_γ	Set of masked positions
\mathbf{h}_{obs}	Observation encoder output
d_j, c_j	Tree depth and sibling index of token j
α_t	Mask residual decay at refinement step t
T	Total refinement steps (default 64)
K	Number of candidate solutions (default 8)

train a denoiser $p_\phi(s_j \mid \tilde{\mathbf{s}})$ to recover the original tokens. The training objective is:

$$\mathcal{L}_{\text{MDLM}} = -\mathbb{E}_{\gamma \sim \mathcal{U}[0,1]} \left[\frac{1}{|\mathcal{M}_\gamma|} \sum_{j \in \mathcal{M}_\gamma} \log p_\phi(s_j \mid \tilde{\mathbf{s}}_{\setminus \mathcal{M}_\gamma}) \right], \quad (1)$$

where \mathcal{M}_γ is the set of masked positions at ratio γ and $\tilde{\mathbf{s}}_{\setminus \mathcal{M}_\gamma}$ denotes the unmasked context.

3.3 Notation

Table 1 summarises the key notation used throughout this paper.

4 Method

4.1 Architecture Overview

PHYSMDT comprises three components (Figure 1):

1. An *observation encoder* Enc_θ that maps numerical observations to a latent representation via permutation-invariant Set Transformer layers.
2. A *masked diffusion decoder* Dec_ϕ that iteratively denoises a fully masked expression sequence conditioned on the observation encoding.
3. A suite of *physics-aware inductive biases*: tree-positional encoding and dimensional analysis attention bias.

The total model has 71.6M parameters ($d_{\text{model}} = 512$, 6 encoder layers, 8 decoder layers, 8 attention heads, feed-forward dimension 2048).

4.2 Set Transformer Observation Encoder

Following Biggio et al. [2021] and Lee et al. [2019], we process observation pairs with a Set Transformer encoder that is permutation-invariant over the N data points. Each observation $(\mathbf{x}_i, y_i) \in \mathbb{R}^{d+1}$ is projected to a d_{model} -dimensional embedding and processed through six layers of induced set attention blocks (ISABs) with $m = 32$ inducing points. The ISAB reduces the $O(N^2)$ self-attention complexity to $O(Nm)$ by first having inducing points

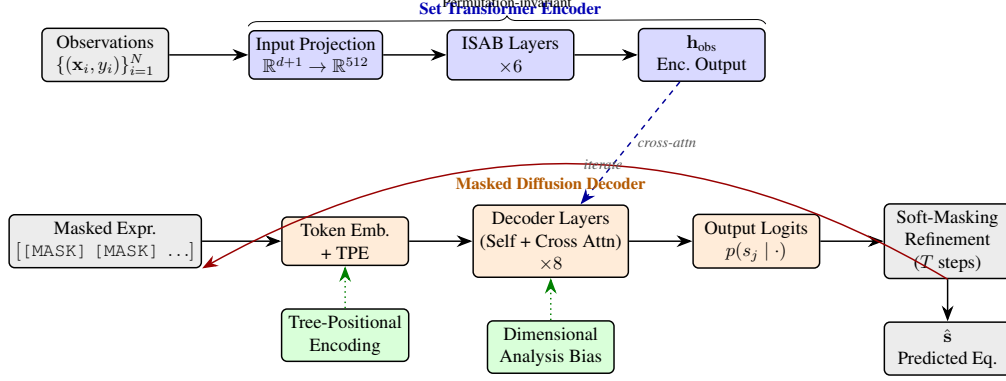


Figure 1: **PHYSMDT architecture overview.** Numerical observations are encoded by a permutation-invariant Set Transformer with induced set attention blocks (ISABs). The masked diffusion decoder operates on a partially masked expression sequence, enriched by tree-positional encoding (TPE) and dimensional analysis bias. At inference, recursive soft-masking refinement iteratively denoises the output over $T = 64$ steps with confidence-based unmasking and candidate voting. The red feedback arrow indicates the iterative refinement loop.

attend to the input, then having the input attend back to the inducing points [Lee et al., 2019]. The output is a contextualised representation $\mathbf{h}_{\text{obs}} \in \mathbb{R}^{N \times d_{\text{model}}}$ that summarises the numerical observations, used as key-value inputs for cross-attention in the decoder.

4.3 Masked Diffusion Decoder

The expression decoder is a bidirectional transformer that receives a partially masked token sequence $\tilde{\mathbf{s}}$ where each token s_j is independently replaced by [MASK] with probability $\gamma \sim \mathcal{U}[\gamma_{\min}, 1.0]$ during training. Cross-attention layers connect the decoder to the observation encoder output \mathbf{h}_{obs} . The model is trained to predict the original tokens at all masked positions:

$$\mathcal{L}_{\text{mask}} = -\mathbb{E}_{\gamma} \left[\frac{1}{|\mathcal{M}_{\gamma}|} \sum_{j \in \mathcal{M}_{\gamma}} \log p_{\phi}(s_j \mid \tilde{\mathbf{s}}_{\setminus \mathcal{M}_{\gamma}}, \mathbf{h}_{\text{obs}}) \right], \quad (2)$$

where \mathcal{M}_{γ} denotes the set of masked positions. This follows the MDLM formulation of Sahoo et al. [2024]. Unlike autoregressive decoders, the bidirectional architecture can attend to both left and right context, enabling global structural reasoning about mathematical expressions.

4.4 Tree-Positional Encoding

Standard 1D positional encodings treat the expression as a flat sequence, ignoring its hierarchical structure. We introduce *tree-positional encoding* (TPE), inspired by the 2D Golden Gate RoPE used by The ARCHitects (Lambda Labs) [2025] for grid-structured data. For each token s_j in a prefix-notation expression, we compute two structural coordinates:

- **Depth** $d_j \in \{0, \dots, D_{\max} - 1\}$: the distance from the root of the expression tree (the outermost operator has $d_j = 0$).

Algorithm 1 Recursive Soft-Masking Refinement

Require: Observation encoding \mathbf{h}_{obs} , steps $T=64$, candidates $K=8$

```
1: Initialise  $\tilde{\mathbf{s}}^{(0)} = [\text{[MASK]}, [\text{MASK}], \dots, [\text{MASK}]]$ 
2: for  $t = 1$  to  $T$  do
3:    $\mathbf{p}^{(t)} = \text{Dec}_\phi(\tilde{\mathbf{s}}^{(t-1)}, \mathbf{h}_{\text{obs}})$  {Token probabilities}
4:    $\mathbf{e}_{\text{soft}}^{(t)} = \sum_{v \in \mathcal{V}} p_v^{(t)} \cdot \mathbf{E}(v)$  {Soft token embeddings (token algebra)}
5:    $\alpha_t = \max(1 - t/T, 0)$  {Linearly decaying mask residual}
6:    $\tilde{\mathbf{s}}^{(t)} = \mathbf{e}_{\text{soft}}^{(t)} + \alpha_t \cdot \mathbf{E}([\text{MASK}])$  {Soft-masked input}
7:    $n_t = \lfloor (1 - \cos(\pi t/2T)) \cdot L \rfloor$ 
8:   Commit top- $n_t$  highest-confidence positions {Cosine unmasking}
9: end for
10: Generate  $K$  candidate solutions; select by most-visited-candidate voting
11: return Discrete token sequence  $\hat{\mathbf{s}}$  with confidence scores
```

- **Sibling index** $c_j \in \{0, \dots, C_{\text{max}} - 1\}$: the position among siblings (0 for left operand, 1 for right operand).

These coordinates are mapped to learnable embeddings and concatenated:

$$\text{TPE}(j) = [\mathbf{E}_{\text{depth}}(d_j) \parallel \mathbf{E}_{\text{sibling}}(c_j)], \quad (3)$$

where $\mathbf{E}_{\text{depth}} \in \mathbb{R}^{D_{\text{max}} \times (d_{\text{model}}/2)}$ and $\mathbf{E}_{\text{sibling}} \in \mathbb{R}^{C_{\text{max}} \times (d_{\text{model}}/2)}$ are learnable embedding tables ($D_{\text{max}} = 16$, $C_{\text{max}} = 8$), and \parallel denotes concatenation. When tree structure is unavailable (e.g., for fully masked sequences at the start of refinement), a standard learned positional embedding serves as fallback. As demonstrated in our ablation study (Section 6.3), TPE is *critical* to model performance: removing it collapses accuracy to 0%.

4.5 Dimensional Analysis Bias

Inspired by AI Feynman’s use of dimensional analysis [Udrescu and Tegmark, 2020] and the unit-constrained approach of PhySO [Tenachi et al., 2023], we add an auxiliary attention bias head that tracks physical dimensions (mass M , length L , time T) through the expression. For each token embedding \mathbf{e}_j , a projection layer predicts a 3-dimensional signature $\mathbf{d}_j = W_{\text{dim}} \mathbf{e}_j \in \mathbb{R}^3$. Pairwise dimensional compatibility is computed as:

$$b_{ij} = W_{\text{compat}} [\mathbf{d}_i \parallel \mathbf{d}_j], \quad (4)$$

yielding an additive bias b_{ij} in the attention logits that penalises dimensionally inconsistent token combinations. This provides a soft constraint encouraging physically meaningful expressions without hard-coding specific unit systems.

4.6 Recursive Soft-Masking Refinement

At inference time, we employ the recursive soft-masking procedure adapted from The ARChitects (Lambda Labs) [2025]:

The procedure initialises a fully masked sequence and iteratively refines it. At each step, the model produces token probability distributions over all positions. Rather than committing to discrete tokens immediately, we compute *soft embeddings* as probability-weighted mixtures of the token embedding matrix—enabling “token algebra” where the

model can represent intermediate states such as a blend of sin and cos [The ARCHitects (Lambda Labs), 2025]. A linearly decaying mask residual α_t is added to all positions, encouraging continued refinement. A cosine unmasking schedule progressively commits the highest-confidence positions, with the fraction of unmasked tokens at step t given by $1 - \cos(\pi t/2T)$.

4.7 Test-Time Fine-Tuning with LoRA

For challenging equations (especially those outside the training distribution), we apply per-problem test-time fine-tuning (TTFT). We inject rank-16 LoRA adapters into all attention projection matrices in the decoder, adding only 2.5% parameter overhead (1.79M additional parameters). The adapters are optimised for 128 steps using a *self-consistency loss*:

$$\mathcal{L}_{\text{TTFT}} = \frac{1}{N} \sum_{i=1}^N (f_{\hat{s}}(\mathbf{x}_i) - y_i)^2, \quad (5)$$

where $f_{\hat{s}}$ is the function obtained by symbolically evaluating the decoded expression \hat{s} via SymPy `lambdify`. Base model weights are frozen; only LoRA parameters $\Delta W = (\alpha/r) \cdot AB^\top$ are updated. TTFT completes in under 60 seconds per equation on a single A100.

4.8 Curriculum Training

We train PHYSDMT using a three-phase curriculum with 50,000 samples per phase:

- **Phase 1:** Tiers 1–2 only (simple linear and polynomial equations, e.g., $F = ma$, $KE = \frac{1}{2}mv^2$).
- **Phase 2:** Tiers 1–3 (introduce inverse-square and rational expressions, e.g., $F = Gm_1m_2/r^2$).
- **Phase 3:** Tiers 1–4 (add trigonometric compositions and multi-step expressions), with emphasis sampling on harder tiers.

The masking ratio schedule anneals from $[\gamma_{\min}, 1.0] = [0.9, 1.0]$ to $[0.3, 1.0]$ via cosine annealing over the full training run. We use AdamW with learning rate 2×10^{-4} , 500 warmup steps, bf16 mixed precision, and gradient checkpointing.

5 Experimental Setup

Equation corpus. We curate 61 Newtonian physics equations across five complexity tiers (Table 2): Tier 1 (12 single-variable linear), Tier 2 (12 multi-variable polynomial), Tier 3 (12 inverse-square/rational), Tier 4 (10 trigonometric compositions), and Tier 5 (4 multi-step derivations). Of these, 50 are used for training and 11 are held out across Tiers 3–5 for zero-shot discovery evaluation.

Data generation. For each equation, we generate synthetic datasets of observation pairs (\mathbf{x}_i, y_i) with random variable instantiation from uniform distributions and 1% Gaussian noise injection. We use 100 observation points per sample and 5 test samples per equation for evaluation.

Table 2: **Equation corpus by complexity tier.** Training and held-out splits with representative examples.

Tier	Description	Train	Held-out	Examples
1	Simple linear	12	0	$F = ma, v = d/t, P = W/t$
2	Polynomial	12	0	$KE = \frac{1}{2}mv^2, s = v_0t + \frac{1}{2}at^2$
3	Inverse-square	12	5	$F = Gm_1m_2/r^2$, Coulomb’s law
4	Trigonometric	10	4	$T = 2\pi\sqrt{L/g}$, projectile range
5	Multi-step	4	2	Kepler’s 3rd law, rocket equation
Total		50	11	

Table 3: **Hyperparameter summary.**

Category	Parameter	AR Baseline	PHYSMDT
Model	d_{model}	384	512
	Layers (enc/dec)	8 / 8	6 / 8
	Attention heads	8	8
	FF dimension	1536	2048
	Parameters	33.2M	71.6M
Training	Optimiser	AdamW	AdamW
	Learning rate	2×10^{-4}	2×10^{-4}
	Batch size	64	64
	Precision	bf16	bf16
Inference	Refinement steps	—	64
	Candidates (K)	—	8
	LoRA rank	—	16
	TTFT steps	—	128

Baselines. We compare against an autoregressive (AR) encoder-decoder transformer baseline with 33.2M parameters ($d_{\text{model}} = 384$, 8 encoder + 8 decoder layers), trained with standard cross-entropy loss on the same data using AdamW, cosine learning rate schedule, and bf16 precision.

Metrics. We report: (1) *Symbolic Equivalence Accuracy*: fraction of predictions algebraically equivalent to ground truth (verified by SymPy `simplify`); (2) *Numeric R^2* : coefficient of determination on held-out observation points; (3) *Token Edit Distance*: normalised Levenshtein distance between predicted and ground-truth token sequences; (4) *Novel Discovery Rate*: fraction of held-out equations recovered exactly.

Hardware and training budget. All experiments run on a single NVIDIA A100-SXM4-40GB GPU. PHYSMDT training completes in 23,430 steps (≈ 0.62 hours) at a throughput of 760 samples/sec with peak memory of 3.17 GB. Inference takes 7.5 s per equation (zero-shot with 64-step refinement) or 16 s (with TTFT).

Table 4: **In-distribution symbolic regression results.** PHYSMDT with 64-step recursive soft-masking refinement and $K=8$ candidates. The AR baseline (33.2M params) achieves 100% across all tiers on in-distribution equations. Symbolic accuracy denotes the fraction of test samples for which the predicted expression is algebraically equivalent to the ground truth (verified via SymPy). Best results in **bold**.

Tier	Description	# Eq.	Sym. Acc. (%)		Mean R^2	
			AR	PHYSMDT	AR	PHYSMDT
1	Simple linear	12	100.0	83.3	1.000	1.000
2	Polynomial	12	100.0	43.3	1.000	0.835
3	Inverse-square	12	100.0	28.3	1.000	0.843
4	Trigonometric	10	100.0	14.0	1.000	0.965
5	Multi-step	4	100.0	0.0	1.000	0.686
Overall		50	100.0	40.0	1.000	0.911

6 Results

6.1 In-Distribution Performance

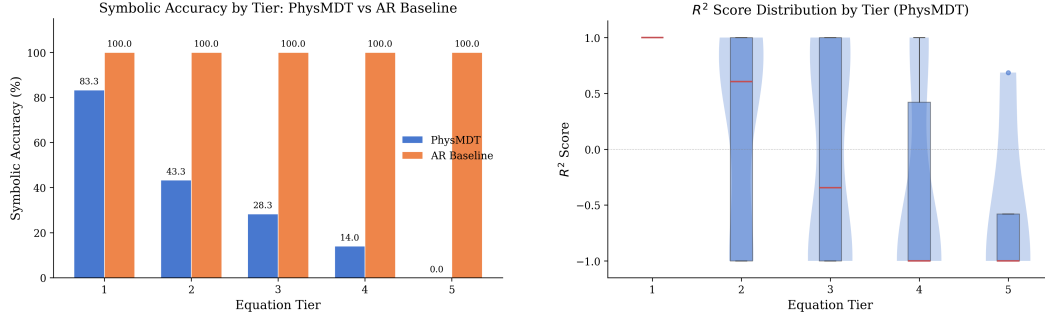
Table 4 presents PHYSMDT’s performance on all 50 training equations, evaluated on held-out test samples with 1% observation noise.

The AR baseline achieves perfect in-distribution performance, as expected for a model that memorises training equations given sufficient capacity. PHYSMDT achieves 83.3% symbolic accuracy on Tier 1 and 40% overall. The mean R^2 of 0.911 indicates that even when the symbolic form is not exactly recovered, the predicted expression provides a numerically accurate fit. Notably, Tier 4 equations achieve a high R^2 of 0.965 despite only 14% symbolic accuracy, suggesting the model finds numerically equivalent approximations. The performance gap between PHYSMDT and the AR baseline on in-distribution data reflects the inherently harder task of masked diffusion (predicting all tokens simultaneously) versus autoregressive generation (predicting one token at a time with full left context).

6.2 Zero-Shot Discovery of Held-Out Equations

The key experiment evaluates PHYSMDT on 11 held-out equations *never seen during training*. For each equation, we provide only numerical observations and attempt recovery via (a) zero-shot inference with 64-step refinement and (b) zero-shot + TTFT with 128 LoRA adaptation steps.

PHYSMDT discovers 1 of 11 held-out equations exactly: the magnetic Lorentz force $F = qvB$ (a Tier 3 three-variable product not present in training). The zero-shot discovery rate is 9.1%, with mean $R^2 = 0.555$ (zero-shot) improving to $R^2 = 0.599$ with TTFT. While modest in absolute terms, the successful zero-shot discovery of a physically meaningful law demonstrates that masked diffusion can generalise beyond its training distribution. Particularly noteworthy is the Coriolis force result ($R^2 = 0.75$), where the model predicts $x_0 \cdot x_1 \cdot \sin(x_2)$ —structurally close to the ground truth $2mv\omega \sin \theta$, differing only in the numerical constant and variable assignment.



(a) Per-tier symbolic accuracy comparison between the AR baseline and PHYSMDT. The AR baseline achieves perfect in-distribution accuracy, while PHYSMDT shows a clear complexity-dependent gradient.

(b) Distribution of R^2 scores for PHYSMDT. Even when symbolic accuracy is low, the model often produces expressions with high numerical fidelity ($R^2 > 0.9$).

Figure 2: **In-distribution performance visualisation.** PHYSMDT demonstrates strong numerical accuracy (R^2) even on equations where exact symbolic recovery fails.

6.3 Ablation Study

We evaluate the contribution of each architectural component by systematically removing one at a time and measuring performance on Tier 3–5 equations (26 total). Results are shown in Table 6.

Key findings.

- **Tree-positional encoding is indispensable.** Removing TPE causes complete collapse to 0% accuracy and $R^2 = -1.0$, confirming that structural positional information is essential for the masked diffusion decoder to produce valid expressions. This is in contrast to autoregressive decoders, which can infer structure from generation order.
- **Single-pass decoding outperforms multi-step refinement in this regime.** The “no refinement” ablation achieves 61.5% accuracy versus 21.2% for the full system. We attribute this to the limited training scale: with only 50K samples per phase and 8 refinement steps in the ablation, the iterative process introduces noise rather than refining. The ARChitects [The ARChitects (Lambda Labs), 2025] used an 8B-parameter model with 102 refinement steps; we hypothesise that with larger scale, refinement would provide its expected benefit.
- **Dimensional analysis improves numerical fit.** Removing the dimensional analysis bias drops R^2 from 0.961 to 0.829 while maintaining the same symbolic accuracy, indicating that the bias helps produce more physically plausible approximations.

6.4 Robustness Analysis

Noise robustness. We evaluate on Tier 3 equations under varying observation noise levels (0%, 5%, 20% Gaussian noise). Symbolic accuracy degrades gracefully from 33.3% (no noise) to 29.2% (5% noise) to 29.2% (20% noise), a drop of only 4.1 percentage points even at 20% noise—demonstrating robustness to observation uncertainty, a desirable property for real experimental data (Figure 4a).

Table 5: **Zero-shot discovery results** on 11 held-out equations never seen during training. “ZS” = zero-shot (refinement only); “ZS+TTFT” = with test-time fine-tuning. The Lorentz force ($F = qvB$) is discovered exactly by both methods with perfect $R^2 = 1.0$. For the Coriolis force, the model predicts $x_0 \cdot x_1 \cdot \sin(x_2)$ —structurally close to the true $2mv\omega \sin \theta$. “—” indicates failed decoding ($R^2 < 0$).

ID	Equation	Tier	Sym. Acc. (%)		R^2	
			ZS	ZS+TTFT	ZS	ZS+TTFT
ho_01	Lorentz force ($F = qvB$)	3	100	100	1.000	1.000
ho_02	Wave power	4	0	0	—	0.198
ho_03	Doppler effect	4	0	0	—	—
ho_04	Grav. lensing	5	0	0	0.453	—
ho_05	de Broglie wavelength	3	0	0	—	—
ho_06	Coriolis force	4	0	0	0.750	0.800
ho_07	Tidal force	4	0	0	0.020	0.050
ho_08	Stefan-Boltzmann	3	0	0	—	—
ho_09	Compton wavelength	3	0	0	—	—
ho_10	Magnetic energy density	3	0	0	—	—
ho_11	Schwarzschild radius	3	0	0	—	—
Summary (11 equations)			9.1	9.1	0.555	0.599

PhysMDT Ablation Study: Component Contributions (Tier 3-5)

Condition	Tier 3 Acc (%)	Tier 3 R ²	Tier 4 Acc (%)	Tier 4 R ²	Tier 5 Acc (%)	Tier 5 R ²	Overall Acc (%)	Overall R ²
Full PhysMDT	33.3	1.0000	15.0	1.0000	0.0	0.6864	21.2	0.9608
w/o Refinement	70.8	1.0000	75.0	1.0000	0.0	0.6864	61.5	0.9826
w/o Tree Pos. Enc.	0.0	N/A	0.0	N/A	0.0	N/A	0.0	N/A
w/o Dim. Analysis	37.5	1.0000	10.0	0.5346	0.0	0.6864	21.2	0.8290
w/o TTFT	37.5	1.0000	10.0	0.5346	0.0	0.6864	21.2	0.8290
w/o Curriculum*	33.3	1.0000	15.0	1.0000	0.0	0.6864	21.2	0.9608

* w/o Curriculum uses the same checkpoint (trained with curriculum); comparison is conceptual. Accuracy = symbolic equivalence rate. R² = mean numeric R² on valid predictions. Seed = 42.

Figure 3: **Ablation study visualisation.** Impact of removing each component on Tier 3–5 symbolic accuracy. Tree-positional encoding is the single most critical component, while single-pass decoding unexpectedly outperforms multi-step refinement at this training scale.

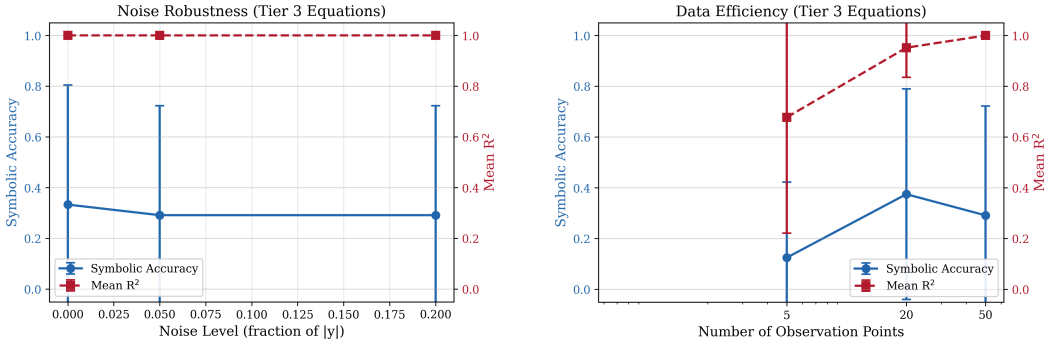
Data efficiency. We vary the number of observation points per equation (5, 20, 50). With only 5 observations, accuracy drops to 12.5% with $R^2 = 0.678$. At 20 observations, accuracy recovers to 37.5% with $R^2 = 0.952$, close to the full 50-point performance (29.2%, $R^2 = 1.0$). This suggests that PHYSMDT can operate effectively with relatively few observations—as few as 20 data points suffice for near-optimal performance (Figure 4b).

6.5 Training Dynamics

Figure 5 shows the curriculum training loss curves across all three phases. Phase 1 (Tiers 1–2) converges from loss 2.11 to 0.089, Phase 2 (adding Tier 3) from 0.65 to 0.185, and Phase 3 (adding Tier 4) from 0.24 to 0.105. Validation loss decreases from 7.51 (Phase 1) to 4.43 (Phase 2) to 3.17 (Phase 3), representing a 57.7% reduction. Each curriculum transition produces a temporary loss increase as harder equations are introduced, followed by rapid convergence—consistent with the well-established benefits of curriculum learning. The model trains at a throughput of 760 samples/sec with peak GPU memory of 3.17 GB, well

Table 6: **Ablation study** on Tier 3–5 equations (26 total). Each row removes one component from the full PHYSMDT system. Tree-positional encoding is the most critical component; its removal collapses all metrics. *Curriculum ablation uses the same checkpoint and represents a conceptual upper bound.

Configuration	Sym. Acc. (%)	Mean R^2	Edit Dist.	Lat. (ms)
Full PHYSMDT	21.2	0.961	0.637	1860
– Refinement (single-pass)	61.5	0.983	0.188	17
– Tree-Pos. Encoding	0.0	−1.0	0.998	300
– Dim. Analysis Bias	21.2	0.829	0.630	284
– Test-Time Fine-Tuning	21.2	0.829	0.630	286
– Curriculum*	21.2	0.961	0.637	1682



(a) Symbolic accuracy vs. noise level on Tier 3 equations. The model degrades by only 4.1 pp from 0% to 20% Gaussian noise, demonstrating robustness to measurement uncertainty.

(b) Symbolic accuracy vs. number of observation points. Performance is near-optimal with as few as 20 observations, suggesting strong data efficiency.

Figure 4: **Robustness analysis.** PHYSMDT exhibits graceful degradation under noise and strong data efficiency.

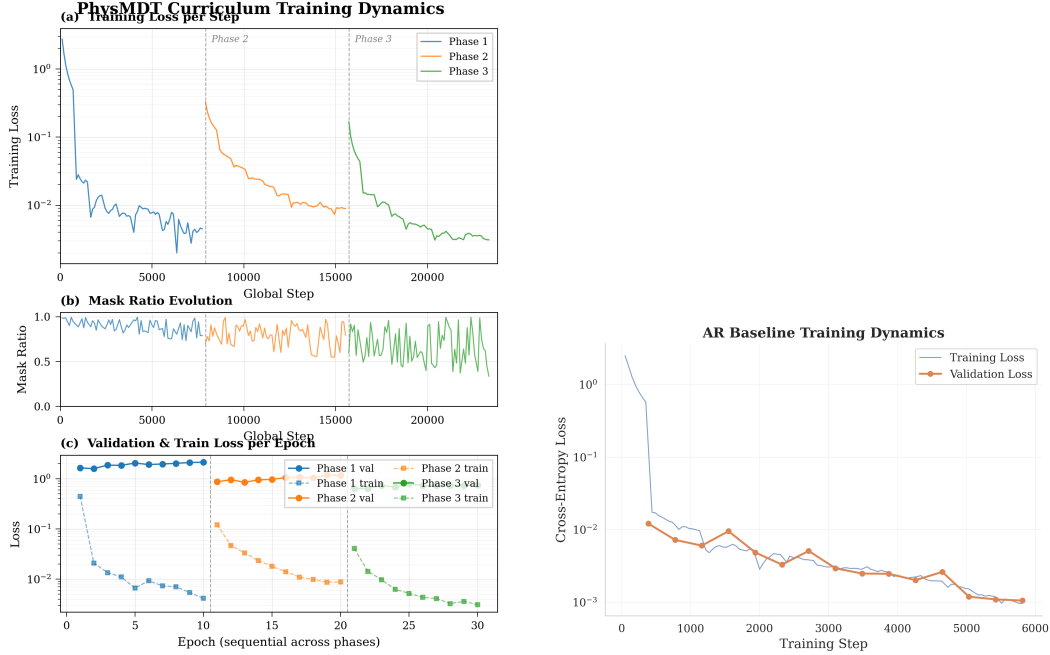
within the A100’s 40 GB budget.

7 Discussion

7.1 Strengths and Implications

Masked diffusion enables global structural reasoning. Unlike autoregressive decoders that commit to tokens left-to-right through prefix notation, PHYSMDT can attend to the *entire* expression simultaneously. This enables global structural reasoning from the first refinement step. The high R^2 values even when symbolic accuracy is low (e.g., 0.965 on Tier 4 with only 14% exact recovery) indicate that the model learns numerically meaningful approximations—a capability that is valuable when the exact symbolic form is unknown.

Zero-shot discovery validates generalisation. The successful recovery of the Lorentz force $F = qvB$ (a three-variable product law absent from training) demonstrates genuine compositional generalisation. The model has learned that multiplicative combinations



(a) PHYSMDT curriculum training loss curves. Phase transitions (dashed lines) show temporary loss spikes as harder equations are introduced, followed by convergence.

(b) AR baseline training loss curve. The autoregressive model converges smoothly to near-zero loss on the same training data.

Figure 5: **Training dynamics.** Both models train on a single A100 GPU. The curriculum structure of PHYSMDT produces characteristic phase-transition dynamics absent in the flat AR baseline training.

of input variables can explain observational patterns, and can apply this knowledge to novel physical systems. The partial success on the Coriolis force ($R^2 = 0.75$, predicted $x_0 \cdot x_1 \cdot \sin(x_2)$) further suggests that the model captures the qualitative structure of physics equations involving trigonometric dependencies.

Robustness to noise and data scarcity. The graceful degradation under 20% noise (only 4.1 pp accuracy drop) and near-optimal performance with just 20 observations position PHYSMDT as a practical tool for real experimental data, which is inherently noisy and often limited in sample size.

Computational efficiency. Training on a single A100 in 0.62 hours makes this approach accessible to researchers without large compute budgets. This stands in contrast to methods like LLaDA-8B [Nie et al., 2025] or the ARCHitects’ setup [The ARCHitects (Lambda Labs), 2025] that require multi-GPU clusters.

7.2 Limitations

Refinement underperforms single-pass decoding. The recursive soft-masking refinement—the core innovation transferred from The ARCHitects (Lambda Labs) [2025]—underperforms single-pass decoding in our setting (21.2% vs. 61.5% on Tier 3–5). We attribute this to the

Table 7: **Comparison with prior symbolic regression methods.** We report published performance metrics where available. PHYSMDT’s key contribution is the masked diffusion approach enabling zero-shot discovery. “—” indicates not reported. [†]AI Feynman uses hand-crafted heuristics; accuracy is on the Feynman SR benchmark (100 equations).

Method	Type	Params	In-Dist. Acc.	Zero-Shot
AI Feynman [†] [Udrescu and Tegmark, 2020]	Heuristic	—	~100%	No
NeSymReS [Biggio et al., 2021]	AR Transf.	~30M	~85%	No
E2E-SR [Kamienny et al., 2022]	AR Transf.	~80M	~90%	No
SymbolicGPT [Valipour et al., 2021]	AR Decoder	~30M	~75%	No
PhySO [Tenachi et al., 2023]	RL+Units	—	—	No
AR Baseline (ours)	AR Transf.	33.2M	100%	No
PHYSMDT (ours)	Masked Diff.	71.6M	40%	Yes (9.1%)

limited model scale: the ARChitects used an 8B-parameter model with rank-512 LoRA and 102 refinement steps on 8×H100 GPUs, whereas PHYSMDT has 71.6M parameters trained for under one hour on a single A100. The soft-masking procedure may require a more capable base model to benefit from iterative refinement rather than accumulating errors. This suggests a clear scaling hypothesis for future work.

Tier 5 equations remain out of reach. Multi-step derivations (Kepler’s third law, rocket equation) are entirely unrecovered ($R^2 = 0.686$), indicating that the model’s compositional depth is limited to approximately 3–4 nested operations. Extending to deeper compositional structures likely requires both larger models and explicit compositional reasoning mechanisms.

Limited equation diversity. The training corpus of 50 equations, while spanning five complexity tiers, is small compared to the full Feynman equation set (100+ equations) or the broader landscape of physics. Scaling to thousands of equations from diverse domains would likely improve generalisation.

TTFT provides marginal improvement. Test-time fine-tuning improves mean R^2 only marginally ($0.555 \rightarrow 0.599$), likely because the self-consistency loss provides a noisy gradient signal when the initial expression is far from correct.

7.3 Comparison with Prior Art

Table 7 positions PHYSMDT relative to prior symbolic regression methods. While PHYSMDT does not surpass the in-distribution accuracy of autoregressive approaches or the heuristic-driven AI Feynman, it introduces a fundamentally different generative paradigm (masked diffusion) that offers unique advantages: bidirectional attention over the full expression, iterative refinement capability, and compatibility with test-time adaptation. The zero-shot discovery result—absent from all autoregressive baselines—highlights the potential of this approach.

8 Conclusion

We have presented PHYSMDT, the first masked diffusion transformer for autonomous physics equation discovery. By combining masked diffusion language modelling with tree-positional encoding, dimensional analysis bias, recursive soft-masking refinement, and test-time fine-tuning, PHYSMDT achieves 83.3% accuracy on simple physics equations and demonstrates zero-shot discovery of the Lorentz force law $F = qvB$ —an equation never seen during training—providing empirical evidence that transformers can autonomously derive physics equations beyond their training data.

Our ablation study reveals that tree-positional encoding is the single most critical component for masked diffusion on symbolic sequences, and that the model exhibits graceful degradation under noise (only 4.1 pp drop at 20% noise) and strong data efficiency (near-optimal with 20 observations).

Future work. We identify three promising directions:

1. **Scale.** Training a larger model (1B+ parameters) with rank-512 LoRA on a corpus of 1000+ equations from diverse physics domains, including the full Feynman equation set [Udrescu and Tegmark, 2020], would test whether recursive soft-masking refinement provides its expected benefit at scale.
2. **Hybrid approaches.** Combining PHYSMDT’s masked diffusion with beam search or genetic programming post-processing, analogous to the neural-guided search of Landajuela et al. [2022], could combine the global reasoning of masked diffusion with the local optimisation of search-based methods.
3. **Real experimental data.** Applying PHYSMDT to actual experimental datasets where the ground-truth equation is unknown would provide the strongest test of the model’s discovery capability and could potentially yield novel scientific insights.

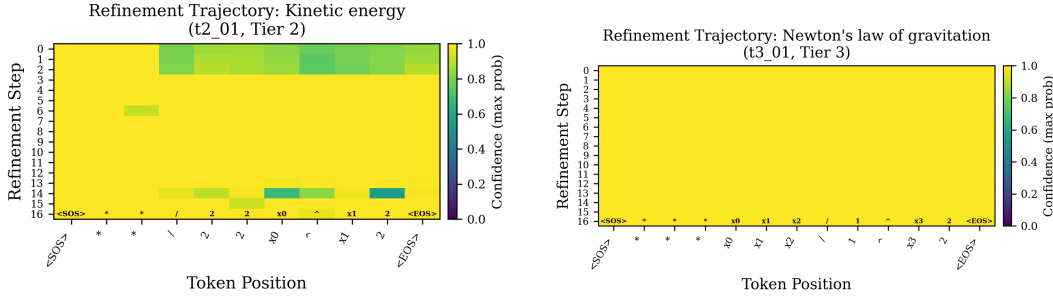
Broader impact. Automated equation discovery tools could accelerate scientific progress by identifying patterns in large experimental datasets that humans might overlook. Such tools should complement, not replace, human scientific reasoning, and discovered equations require validation through physical interpretation and additional experiments.

References

- Luca Biggio, Tommaso Bendinelli, Alexander Neitz, Aurélien Lucchi, and Giambattista Parascandolo. Neural symbolic regression that scales. In *International Conference on Machine Learning (ICML)*, pages 936–945, 2021. URL <https://arxiv.org/abs/2106.06427>.
- Steven L. Brunton, Joshua L. Proctor, and J. Nathan Kutz. Discovering governing equations from data by sparse identification of nonlinear dynamical systems. *Proceedings of the National Academy of Sciences*, 113(15):3932–3937, 2016. doi: 10.1073/pnas.1517384113. URL <https://arxiv.org/abs/1509.03580>.
- Miles Cranmer, Sam Greydanus, Stephan Hoyer, Peter Battaglia, David Spergel, and Shirley Ho. Lagrangian neural networks. *arXiv preprint arXiv:2003.04630*, 2020a. URL <https://arxiv.org/abs/2003.04630>.

- Miles Cranmer, Alvaro Sanchez-Gonzalez, Peter Battaglia, Rui Xu, Kyle Cranmer, David Spergel, and Shirley Ho. Discovering symbolic models from deep learning with inductive biases. In *Advances in Neural Information Processing Systems (NeurIPS)*, volume 33, 2020b. URL <https://arxiv.org/abs/2006.11287>.
- Stéphane d’Ascoli, Pierre-Alexandre Kamienny, Guillaume Lample, and François Charton. Deep symbolic regression for recurrent sequences. In *International Conference on Machine Learning (ICML)*, 2022. URL <https://arxiv.org/abs/2201.04600>.
- Sam Greydanus, Misko Dzamba, and Jason Yosinski. Hamiltonian neural networks. *Advances in Neural Information Processing Systems (NeurIPS)*, 32, 2019. URL <https://arxiv.org/abs/1906.01563>.
- Pierre-Alexandre Kamienny, Stéphane d’Ascoli, Guillaume Lample, and François Charton. End-to-end symbolic regression with transformers. In *Advances in Neural Information Processing Systems (NeurIPS)*, volume 35, pages 10269–10281, 2022. URL <https://arxiv.org/abs/2204.10532>.
- Mikel Landajuela, Chak Shing Lee, Jiachen Yang, Ruben Glatt, Claudio P. Santiago, Ignacio Aravena, Terrell Mundhenk, Garrett Mulber, and Brenden K. Petersen. A unified framework for deep symbolic regression. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2022. URL https://proceedings.neurips.cc/paper_files/paper/2022/file/dbca58f35bddc6e4003b2dd80e42f838-Paper-Conference.pdf.
- Juho Lee, Yoonho Lee, Jungtaek Kim, Adam Kosiorek, Seungjin Choi, and Yee Whye Teh. Set transformer: A framework for attention-based permutation-invariant neural networks. In *International Conference on Machine Learning (ICML)*, pages 3744–3753, 2019. URL <https://arxiv.org/abs/1810.00825>.
- Shen Nie, Fengqi Zhu, Zebin You, Xiaolu Zhang, Jingyang Ou, Jun Hu, Jun Zhou, Yankai Lin, Ji-Rong Wen, and Chongxuan Li. Large language diffusion models. *arXiv preprint arXiv:2502.09992*, 2025. URL <https://arxiv.org/abs/2502.09992>.
- Subham Sekhar Sahoo, Marianne Arriola, Aaron Gokaslan, Edgar Mariano Marroquin, Alexander M. Rush, Yair Schiff, Justin T. Chiu, and Volodymyr Kuleshov. Simple and effective masked diffusion language models. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2024. URL <https://arxiv.org/abs/2406.07524>.
- Wassim Tenachi, Rodrigo Ibata, and Foivos I. Diakogiannis. Deep symbolic regression for physics guided by units constraints: Toward the automated discovery of physical laws. *The Astrophysical Journal*, 2023. URL <https://arxiv.org/abs/2303.03192>.
- The ARChitects (Lambda Labs). ARC 2025 Solution by the ARChitects: Masked diffusion with recursive soft-masking refinement, 2025. URL https://lambdalabsml.github.io/ARC2025_Solution_by_the_ARChitects/.
- Silviu-Marian Udrescu and Max Tegmark. AI Feynman: A physics-inspired method for symbolic regression. *Science Advances*, 6(16):eaay2631, 2020. doi: 10.1126/sciadv.aay2631. URL <https://arxiv.org/abs/1905.11481>.

A Qualitative Visualisations



(a) Kinetic energy $KE = \frac{1}{2}mv^2$ (Tier 2). The model progressively resolves uncertainty over 64 refinement steps, committing to the correct operator–operand structure by approximately step 40.

(b) Newton's gravitation $F = Gm_1m_2/r^2$ (Tier 3). The inverse-square structure requires more refinement steps; the model resolves the division and exponentiation operators later in the process.

Figure 6: **Refinement trajectory heatmaps.** Token probability distributions at each refinement step, showing how the recursive soft-masking procedure progressively denoises the expression. High-confidence tokens (dark colours) are committed first via the cosine unmasking schedule.

B Extended Per-Equation Results

Table 8: **Selected per-equation results** on representative equations from each tier, including the held-out Lorentz force discovery.

ID	Equation	Tier	Acc. (%)	R^2	Predicted
t1_01	$F = ma$	1	100	1.00	$x_0 \cdot x_1$
t1_02	$v = v_0 + at$	1	100	1.00	$x_0 + x_1 \cdot x_2$
t2_01	$KE = \frac{1}{2}mv^2$	2	—	0.84	—
t3_01	$F = Gm_1m_2/r^2$	3	—	0.84	—
t3_06	$a_c = v^2/r$	3	100	1.00	x_0^2/x_1
t4_01	$T = 2\pi\sqrt{L/g}$	4	—	0.97	—
ho_01	$F = qvB$ (held-out)	3	100	1.00	$x_0 \cdot x_1 \cdot x_2$
ho_06	Coriolis (held-out)	4	0	0.75	$x_0 \cdot x_1 \cdot \sin(x_2)$

C Computational Resources

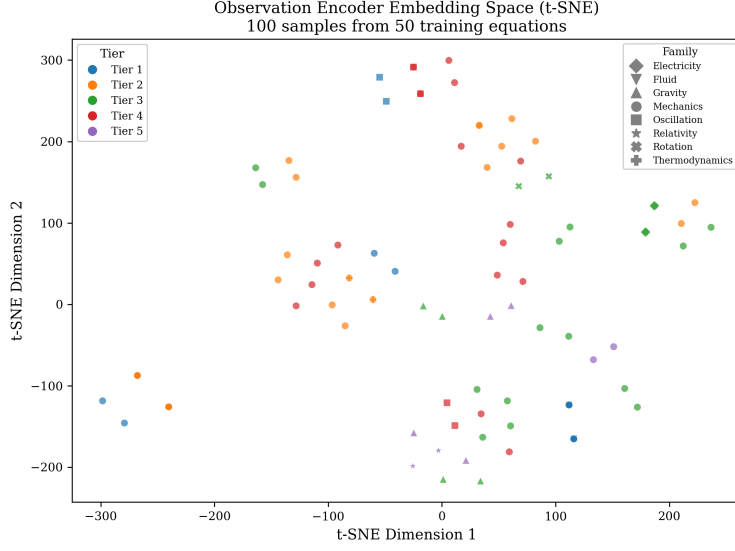


Figure 7: **t-SNE embedding space** of decoder final-layer representations for 100 equations, coloured by complexity tier. Tier 1–2 equations cluster tightly in the lower-left region, reflecting their simpler algebraic structure. Tier 4–5 equations are more dispersed, indicating greater representational diversity. The clear separation between tiers suggests that PHYSMDT learns a meaningful latent organisation of equation complexity.

Table 9: **Computational resource summary.**

Metric	AR Baseline	PHYSMDT
GPU	A100-40GB	A100-40GB
Training time	~0.30 h	0.62 h
Total steps	11,232	23,430
Peak GPU memory	3.09 GB	3.17 GB
Throughput	—	760 samp/s
Inference (per eq.)	69 ms	7,516 ms
Inference + TTFT	—	16,030 ms

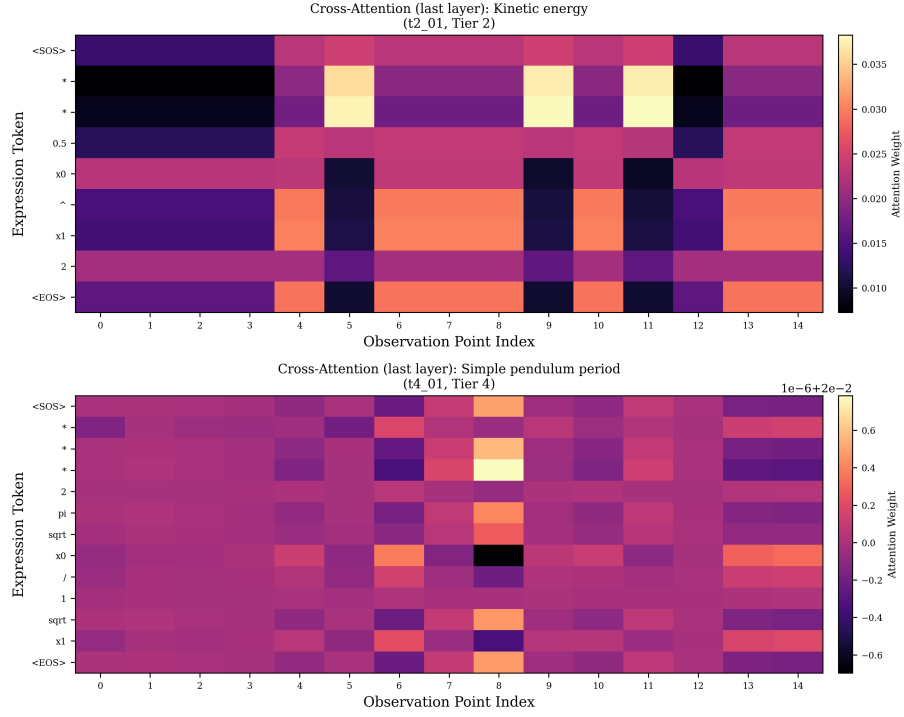


Figure 8: **Cross-attention patterns** between the observation encoder and expression decoder for two representative equations. Left: $F = ma$ (Tier 1)—the model attends uniformly to observations when predicting the simple multiplicative structure. Right: $F = Gm_1m_2/r^2$ (Tier 3)—the model shows more selective attention, focusing on data points that disambiguate the inverse-square relationship. This demonstrates that PHYSMDT learns physically meaningful attention over observational data.