

# Solve the Loop: Attractor Models for Language and Reasoning

Jacob Fein-Ashley, Paria Rashidinejad

University of Southern California

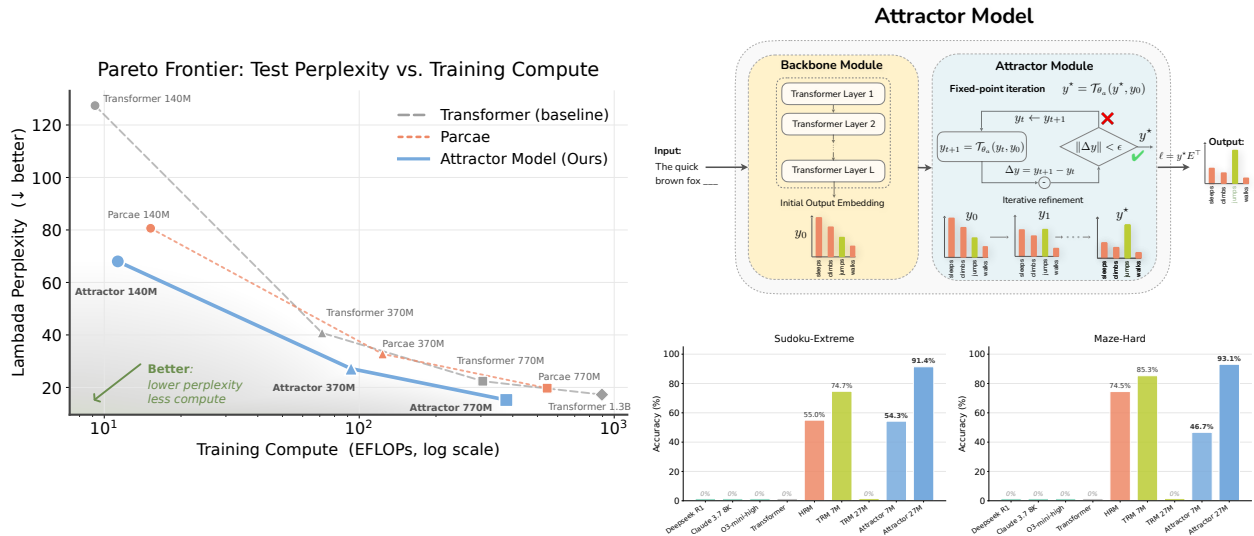
Looped Transformers offer a promising alternative to purely feed-forward computation by iteratively refining latent representations, improving language modeling and reasoning. Yet recurrent architectures remain unstable to train, costly to optimize and deploy, and constrained to small, fixed recurrence depths. We introduce *Attractor Models*, in which a backbone module first proposes output embeddings, then an attractor module refines them by solving for the fixed point, with gradients obtained through implicit differentiation. Thus, training memory remains constant in effective depth, and iterations are chosen adaptively by convergence. Empirically, Attractor Models outperform existing models across two regimes, large-scale language-model pretraining and reasoning with tiny models. In language modeling, Attractor Models deliver a *Pareto improvement* over standard Transformers and stable looped models across sizes, improving perplexity by up to 46.6% and downstream accuracy by up to 19.7% while reducing training cost. Notably, a 770M Attractor Model outperforms a 1.3B Transformer trained on twice as many tokens. On challenging reasoning tasks, we show that our model with only 27M parameters and approximately 1000 examples achieves 91.4% accuracy on Sudoku-Extreme and 93.1% on Maze-Hard, scaling favorably where frontier models like Claude and GPT o3, fail completely, and specialized recursive reasoners collapse at larger sizes. Lastly, we show that Attractor Models exhibit a novel phenomenon, which we call *equilibrium internalization*: fixed-point training places the model’s initial output embedding near equilibrium, allowing the solver to be removed at inference time with little degradation. Together, these results suggest that Attractor Models make iterative refinement scalable by turning recurrence into a computation the model can learn to internalize.

**Date:** May 12, 2026

**Website:** <https://attractor-models.github.io/>

**Code:** <https://github.com/jacobfa/Attractor>

**Correspondence:** {feinashl, paria.rashidinejad}@usc.edu



**Figure 1 Left: Lambda perplexity versus training FLOPs.** Attractor Models achieve strong language-modeling performance with less compute. **Top Right: Attractor Model architecture.** A non-recurrent backbone first maps the input to an initial output embedding proposal, which the attractor module refines by solving a fixed-point iteration before decoding the (approximate) equilibrium into the final output distribution. **Bottom Right: Performance on hard reasoning tasks.** Attractor Models outperform frontier and specialized recursive models on hard reasoning tasks.

# 1 Introduction

The modern language-modeling era has been dominated by Transformers (Vaswani et al., 2017), which produce each token through a fixed feed-forward computation. This recipe has been extraordinarily successful (OpenAI, 2023; Gemini Team, 2023; Grattafiori et al., 2024; Anthropic, 2024; Guo et al., 2025), but it leaves a basic question unresolved: should each token be the product of a single pass of computation, or should a model be able to refine its latent prediction before committing to an output? A growing body of work suggests that such refinement can be powerful. Chain-of-thought reasoning (Wei et al., 2023) can be viewed as one form of such refinement, where a model writes intermediate tokens, feeds them back into its context, and uses them to shape later predictions. Yet this routes computation through the discrete token channel and forces “thinking” to be written down, even when the effect might be to merely refine internal representations.

This limitation has inspired several lines of work on latent (or implicit) thinking and a re-emergence of architectural recurrence, which move thinking away from purely token-level generation. These include universal Transformers (Dehghani et al., 2019), looped Transformers (Giannou et al., 2023a; Yang et al., 2024a), recurrent-depth Transformers (Kohli et al., 2026), looped language models (Zhu et al., 2025c), latent reasoning methods (Geiping et al., 2025a; Saunshi et al., 2025a), and continuous chain-of-thought approaches (Hao et al., 2025; Mohtashami et al., 2023; Zhu et al., 2026). Looped architectures can, in principle, express iterative or algorithmic procedures (Yang et al., 2024b; Giannou et al., 2023a), emulate additional depth through weight sharing (Dehghani et al., 2019; Zhu et al., 2025c), reduce the context-length costs of token-level reasoning, and improve downstream generalization (Labovich, 2026b; Fan et al., 2025). Empirically, recent looped language models offer gains in language modeling and reasoning (Zhu et al., 2025c; Geiping et al., 2025a), and tiny recursive models (Wang et al., 2025; Jolicoeur-Martineau, 2025) have shown that recurrence can be useful in hard reasoning tasks in small-data regimes.

The challenge is that recurrence has proven difficult to use as a stable architectural building block. Recurrence is often accompanied by unstable training, large memory requirements that grow linearly with the number of recurrent steps, and significant, sequential compute (Geiping et al., 2025a; Zhu et al., 2025c; Prairie et al., 2026). Training recurrent networks typically requires backpropagation through time (or, depth) and carefully designed stabilization techniques; even then, latent-thinking models remain fragile and difficult to optimize (Wei et al., 2025; Özeren and Aßenmacher, 2025; Deng et al., 2026, 2025; Rizvi-Martel et al., 2026). For training, looped language models tend to require substantially more compute than comparable feed-forward models and can become memory-limited at larger recurrence depths. For instance, Geiping et al. (2025a) reports that training a recurrent model can consume raw FLOPs comparable to those of a feed-forward model ten times larger. At the opposite end of the spectrum, specialized tiny recursive reasoners exhibit a troubling “less is more” behavior and respond negatively to scaling: increasing model size can degrade or even collapse performance (Jolicoeur-Martineau, 2025).

## 1.1 Contributions

In this work, we design a general-purpose architecture for iterative refinement that is (i) stable to train, (ii) uses constant-memory in the number of refinement steps, (iii) is substantially cheaper to train than explicit unrolling, (iv) is efficient during inference, and (v) achieves a strong performance across both large-scale language modeling and hard reasoning with tiny models.

**Refine outputs by solving the loop with Attractor Models.** We introduce *Attractor Models*, a new family of architectures that treat latent refinement as a fixed-point problem in the output embedding space. The model first proposes an initial guess embedding using a non-recurrent backbone module (implemented as a Transformer in ours). A separate, typically smaller, recurrent network then refines this guess (Figure 1). Recent mechanistic analyses of looped language models demonstrate that, for the vast majority of tokens, the recurrent trajectory converges to a fixed point (Blayney et al., 2026). We build directly on this observation and instead of unrolling the loop for a predefined number of steps, we solve for the state to which the loop converges, inspired by Deep Equilibrium Models (DEQ; (Bai et al., 2019)). The name Attractor Model comes from dynamical systems, where an attractor is a set of states toward which a system evolves. In a sense, Attractor Models can be viewed as *thinking* before producing each token: the backbone proposes an initial latent prediction, the attractor module refines it to equilibrium, then decodes it into the output distribution.

**Attractor Models offer stable, constant-memory, efficient training, and adaptive refinement.** Unlike looped LMs, which finitely unroll the recurrent block, Attractor Models solve an equilibrium by treating the prediction target as a fixed-point computation. The number of refinement steps is therefore chosen adaptively according to convergence during both training and inference. We show that the memory cost during training remains constant in the number of iterations; whereas standard looped language models have a linear scaling increase with the number of loops. Our experiments demonstrate that the two-stage structure of Attractor Models, in which the backbone proposes and the attractor refines, enables stable, efficient training and strong performance.

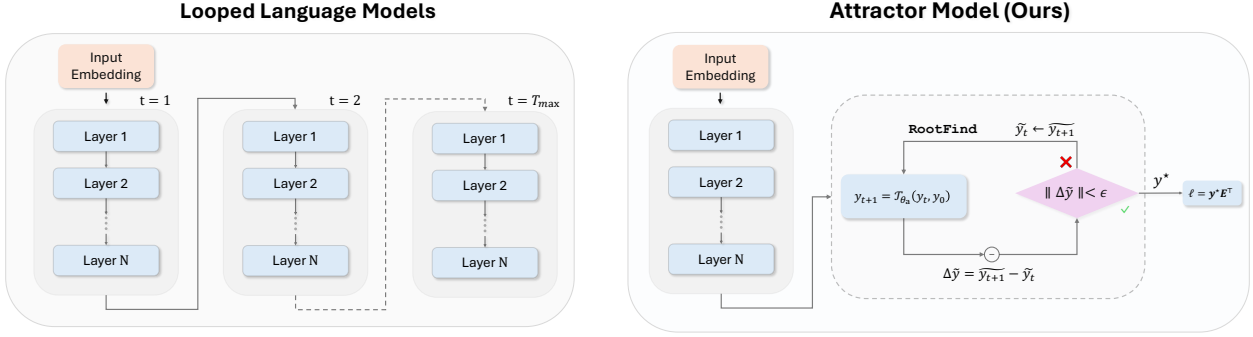
**Novel phenomenon: Equilibrium internalization.** We observe that despite the fact that Attractor models are trained only with the next-token prediction loss, they learn to make the solver unnecessary. During training, the backbone’s initial prediction moves progressively closer to the fixed point, so fewer refinement steps are needed to reach approximate equilibrium (c.f. Figures 6 and 7). We call this phenomenon *equilibrium internalization*: the model appears to self-distill the iterative refinement process into its own initial output embedding, through a form of automatic curriculum. In this sense, recurrence acts as a moving training target, teaching the backbone where its computation should converge.

**Strong performance in large-scale language modeling and hard reasoning with tiny models.** Our experiments show that Attractor Models scale across two regimes. In large-scale language modeling, Attractor Models consistently outperform standard Transformers and stable looped language models across small (140M), medium (370M), and large (770M) sizes, delivering a Pareto improvement (Figure 1). We show that our models improve validation perplexity, out-of-distribution perplexity on Lambada (Paperno et al., 2016), and downstream benchmark accuracy while using substantially less training compute than comparable looped baselines. Notably, a 770M-parameter Attractor Model outperforms a 1.3B-parameter Transformer trained on twice as many tokens. Compared to looped LM Parcae (Prairie et al., 2026), our models use up to 31% less training compute, while avoiding the memory growth associated with explicit unrolling. In hard reasoning tasks with tiny models, with only 27M parameters and approximately 1000 training examples, Attractor Models achieve 91.4% accuracy on Sudoku-Extreme and 93.1% on Maze-Hard. In this regime, standard Transformers as well as proprietary frontier models such as DeepSeek R1, Claude, and o3-mini fail completely at 0%, while specialized recursive architectures underperform our model and collapse when scaled. Attractor Models, in contrast, improve with scale.

## 2 Background: Looped Architectures

We begin with background on looped architectures. Let  $x = (x_1, \dots, x_n) \in \mathcal{V}^n$  be an input sequence over vocabulary  $\mathcal{V}$ , and let  $d$  denote the model width. Looped models can be written as a composition of three units: a prelude unit  $\tilde{x} = \mathcal{P}(x) \in \mathbb{R}^{n \times d}$ , which produces an input representation  $\tilde{x} \in \mathbb{R}^{n \times d}$ ; a weight-tied recurrent unit  $h_{t+1} = \mathcal{R}(h_t, \tilde{x})$ , which is applied repeatedly to a latent state  $h_t \in \mathbb{R}^{n \times d}$  for  $T$  steps; and a coda unit, which maps the final latent state to output probabilities  $p = \mathcal{C}(h_T) \in \Delta(\mathcal{V})^n$ . Importantly, looped architectures commonly initialize the latent state at an *uninformative value*, such as  $h_0 = 0$  or Gaussian noise  $h_0 \sim \mathcal{N}(0, \sigma^2 I)$  (Geiping et al., 2025a; Prairie et al., 2026; Bai et al., 2019). Furthermore, the recurrent step may use the input representation only at the first step (Zhu et al., 2025c). or at every recurrent step (Geiping et al., 2025a; Prairie et al., 2026). Such injection may be through addition or concatenation with the recurrent state.

Models such as Parcae (Prairie et al., 2026), Huggin (Geiping et al., 2025a), and Ouro (Zhu et al., 2025c) differ mainly in how they train, stop, or scale this looped architecture. In particular, the recurrence depth  $T$  is a central design choice in these models. It may be fixed (Jolicoeur-Martineau, 2025), sampled during training (Geiping et al., 2025a; McLeish et al.; Prairie et al., 2026), or determined by an auxiliary halting mechanism (Bae et al., 2025b; Zhu et al., 2025c). Training then minimizes an objective averaged over both the data distribution and the chosen recurrence-depth mechanism, typically by backpropagation through depth. Consequently, both training cost and gradient memory are tied to the number of recurrent steps. Furthermore, changing  $T$  at inference introduces a train–test mismatch, since the model is evaluated under a different computation graph than the one used during training, leading to degraded performance.



**Figure 2 Comparison of looped language models vs. Attractor Models.** Looped language models repeatedly apply a shared block for a finite number of steps before decoding the final state. In contrast, Attractor Models use a backbone to produce an initial output embedding, then refine it with an attractor module until the fixed-point residual is small, and decode the resulting approximate equilibrium.

### 3 Solve the Loop with Attractor Models

As discussed in the previous section, standard looped language models (Zhu et al., 2025c; Prairie et al., 2026) use weight sharing to recurrently refine a hidden state that is initialized from an uninformative value and based on input embeddings. Predictions are then read out after a finite number of loops (Prairie et al., 2026), or once an auxiliary halting head becomes confident (Graves, 2017; Zhu et al., 2025c; Prairie et al., 2026). This design carries three drawbacks: the loop count must be chosen at training time, training memory grows linearly in the number of loops, and accuracy degrades when more loops are run at inference than were seen during training (Zhu et al., 2025c). As a result, recurrence often comes with unstable training, growing memory requirements, and large sequential compute, in some cases approaching the costs of training non-recurrent models ten times larger (Geiping et al., 2025a).

Recent mechanistic analyses of looped language models (Blayney et al., 2026) reveals that for the vast majority of tokens, the recurrent trajectory eventually converges to a fixed point. This suggests that the weight-tied recurrent modules are often approximating an underlying fixed-point computation, doing so through the recursive application of the weight-tied block truncated after  $T$  steps. This observation motivates the design of Attractor Models, which we subsequently describe.

#### 3.1 Attractor Model: Backbone and Attractor Modules

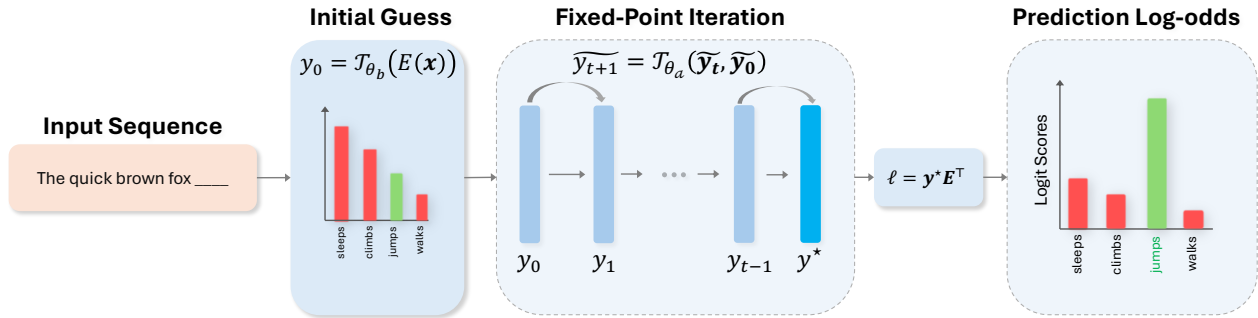
Motivated by the fixed-point behavior observed in looped models, we model recurrent refinement as an attractor. Rather than training a model to produce good predictions after a prescribed number of recurrent steps, we define the output as the equilibrium of the refinement process. Attractor Models consist of two modules: the *backbone module* (typically a larger Transformer network) first proposes a meaningful initial output embedding, and the *attractor module* (typically a smaller Transformer-based network) then refines this proposal until convergence. This makes the number of refinement steps a solver choice rather than a fixed architectural choice.

We first start by mapping the inputs  $x$  into input embeddings  $\tilde{x} = E(x) \in \mathbb{R}^{n \times d}$ , where  $E$  denotes the tied embedding/unembedding. Then, the input embedding is processed by the backbone and attractor modules as described below.

**The backbone module proposes an initial “guess ” output embedding.** The backbone module  $\mathcal{T}_{\theta_b}$  maps the input embeddings to an initial proposal:

$$\tilde{y}_0 = \mathcal{T}_{\theta_b}(\tilde{x}), \text{ where } \tilde{x} = E(x). \quad (1)$$

We use  $\tilde{y}_0$  as an initialization for the attractor module. Instead of initializing the loop from zero, noise, or an input-side representation, Attractor Models initialize the recurrent computation from a state that is already



**Figure 3 Overview of Attractor Models.** The backbone maps the input embeddings to an output-side proposal  $\tilde{y}_0 = \mathcal{T}_{\theta_b}(E(x))$ . The attractor module then refines this proposal through the proposal conditioned iteration  $\tilde{y}_{t+1} = \mathcal{T}_{\theta_a}(\tilde{y}_t, \tilde{y}_0)$  until convergence to an approximate equilibrium  $\tilde{y}^*$  or a maximum number of solver steps. The equilibrium is decoded with the tied unembedding  $\tilde{y}^* E^\top$ .

a coherent prediction embedding. In practice,  $\mathcal{T}_{\theta_b}$  is a relatively high-capacity causal Transformer, so the refinement begins near a meaningful initialization rather than 0. We find that this makes training our method stable compared to DEQ, which experiences a blow-up in the number of iterations used later in training; whereas our method stabilizes later in training (c.f. Figure 6(b)).

**The attractor module refines the output embedding.** The attractor module is a separate weight-tied refinement network  $\mathcal{T}_{\theta_a}$ . Starting from the backbone proposal  $\tilde{y}_0$ , it repeatedly refines the output embedding according to

$$\tilde{y}_{t+1} = \mathcal{T}_{\theta_a}(\tilde{y}_t, \tilde{y}_0), \text{ where } \tilde{y}_0 = \mathcal{T}_{\theta_b}(\tilde{x}). \quad (2)$$

Here, we persistently inject the initial guess  $\tilde{y}_0$  at every refinement step. This persistent injection keeps the attractor proposal-dependent and prevents it from collapsing to a proposal-independent fixed point. We ablate this conditioning mechanism in Section 4. Importantly, we *warm-start* the attractor module at an informative proposal  $\tilde{y}_0$ , in contrast to existing work that initialize the recurrent state at uninformative values such as zero or Gaussian noise (Geiping et al., 2025a; Prairie et al., 2026); see Table 6 for a comparison.

Rather than rolling out recurrent steps to reach a fixed point, we directly solve for the convergence:

$$\mathcal{A}_{\theta_a}(\tilde{y}^*, \tilde{y}_0) := \mathcal{T}_{\theta_a}(\tilde{y}^*, \tilde{y}_0) - \tilde{y}^* = 0 \Rightarrow \tilde{y}^* = \text{RootFind}(\mathcal{A}_{\theta_a}(\cdot, \tilde{y}_0); y_0). \quad (3)$$

In the forward pass, we compute this equilibrium with a root finder initialized at the backbone proposal. In our implementation, the `RootFind` algorithm uses Anderson acceleration, which combines a small window of past iterates and residuals to reach the fixed point faster than plain recursion. The solver exits when  $\|\mathcal{A}_{\theta_a}(\tilde{y}_t, \tilde{y}_0)\|_2 / \|\tilde{y}_t\|_2 < \varepsilon$  or after  $T_{\max}$  steps. Thus, the computation is controlled by the convergence of the residual rather than by a learned halting head or a preset loop count. In contrast to fixed unrolling, the number of refinement steps can therefore vary at inference time without changing the model. Finally, the equilibrium embedding is decoded with the tied unembedding.

Parameters of the Attractor Models consist of the tied embedding/unembedding matrices, the backbone module, and parameters of the attractor module:  $\theta := (\theta_a, \theta_b, E)$ . Compared to looped models, Attractor Models change both the starting point and the endpoint of recurrence: we initialize the loop from the output guess from the backbone network  $\tilde{y}_0$ , and the decoded state is the attractor  $\tilde{y}^*$  rather than a finite unroll.

### 3.2 Training and Inference of Attractor Models

We now describe the training procedure for Attractor Models. We first explain how to differentiate through the fixed-point solver using implicit differentiation, and then show how the model is optimized with the standard cross-entropy language-modeling loss applied to the output  $y^*$ .

**Backward pass and implicit differentiation.** Because Attractor Models define the output embedding as the solution to a fixed-point equation, we differentiate through the equilibrium using the implicit function theorem

(Krantz and Parks, 2002b). Let  $\mathcal{L}$  denote the training loss and let  $v = \partial\mathcal{L}/\partial\tilde{y}^*$ . Applying the implicit function theorem to  $\mathcal{A}_{\theta_a}(\tilde{y}^*, \tilde{y}_0) = 0$  gives

$$\frac{\partial\mathcal{L}}{\partial\theta} = u^\top \frac{\partial\mathcal{T}_{\theta_a}(\tilde{y}^*, \tilde{y}_0)}{\partial\theta}, \quad u = (I - J_{\tilde{y}}^\top)^{-1} v, \quad \text{where } J_{\tilde{y}} = \left. \frac{\partial\mathcal{T}_{\theta_a}(\tilde{y}, \tilde{y}_0)}{\partial\tilde{y}} \right|_{\tilde{y}=\tilde{y}^*}. \quad (4)$$

The derivative with respect to  $\theta = (\theta_a, \theta_b, E)$  includes the direct dependence on the attractor parameters  $\theta_a$ , as well as the dependence on the initialization  $\tilde{y}_0 = \mathcal{T}_{\theta_b}(E(x))$  through the backbone parameters  $\theta_b$  and the tied embedding parameters  $E$ .

Following prior work on implicit models (Geng et al., 2022; Fung et al., 2021), we use the one-step approximation  $u \approx v$ . This avoids the extra linear solve for  $u$  and reduces the backward pass to one vector–Jacobian product through  $\mathcal{A}_{\theta_a}$ . Since we do not backpropagate through every solver step, memory in the attractor block does not grow with the number of forward iterations. In Section 4, we show that the Anderson solver yields only marginal quality gains, whereas the one-step approximation enables substantially cheaper training.

*Remark.* Attractor Models are implicit equilibrium models in the spirit of DEQs (Bai et al., 2019), but the equilibrium plays a different role. Classical DEQs replace the prediction network with a hidden-state equilibrium  $z^*$  (single-layer), decoded with a separate output head. We instead keep a standard causal Transformer backbone and add an equilibrium refinement block on top of its prediction state. The fixed point lives directly in the tied embedding space, so every iterate  $\{y_0, y_1, \dots, y^*\}$  is already a representation in the output space that can be decoded. This gives three practical differences: (i) the solver is initialized from a semantically meaningful proposal  $\tilde{y}_0$  rather than from an uninformative state such as zero (as in DEQ), (ii) inference can stop according to a residual tolerance  $\varepsilon$  rather than a fixed depth or learned halting head, and (iii) DEQ shows that scaling the number of DEQ blocks can harm the performance of their method, whereas we allow for an arbitrary depth backbone transformer and show that we can use a variable number of solver blocks.

---

**Algorithm 1** Training Attractor Models.

---

**Require:** input tokens  $x$ , parameters  $\theta = (\theta_a, \theta_b, E)$ , tolerance  $\varepsilon$ , max iterations  $T_{\max}$

```

// Forward pass
1:  $\tilde{x} \leftarrow E(x)$  ▷ tied token embedding
2:  $\tilde{y}_0 \leftarrow \mathcal{T}_{\theta_b}(\tilde{x})$  ▷ backbone proposal / initialization
3: Define  $\mathcal{A}_{\theta_a}(\tilde{y}; \tilde{y}_0) \leftarrow \mathcal{T}_{\theta_a}(\tilde{y}, \tilde{y}_0) - \tilde{y}$ 
4:  $\tilde{y}^* \leftarrow \text{RootFind}(\mathcal{A}_{\theta_a}(\cdot; \tilde{y}_0); \tilde{y}_0, \varepsilon, T_{\max})$ 
5:  $p \leftarrow \text{Softmax}(\tilde{y}^* E^\top)$  ▷ tied unembedding and softmax
// Backward pass, given  $v = \partial\mathcal{L}/\partial\tilde{y}^*$ 
6: solve  $(I - J_{\tilde{y}}^\top) u = v$  via  $\text{RootFind}(u' \mapsto (I - J_{\tilde{y}}^\top) u' - v; u_0 = v)$ 
7:  $\partial\mathcal{L}/\partial\theta \leftarrow u^\top \partial\mathcal{T}_{\theta_a}(\tilde{y}^*, \tilde{y}_0)/\partial\theta$ 

```

---

**Training objective and inference.** We train Attractor Models with the standard next-token prediction cross-entropy objective applied to the fixed-point output  $y^*$ . Inference reuses the same equilibrium computation. Given an input sequence  $x$ , the backbone first produces the proposal  $\tilde{y}_0 = \mathcal{T}_{\theta_b}(E(x))$ , the attractor solver computes  $\tilde{y}^*$ , and the tied unembedding decodes  $\tilde{y}^*$  into next-token probabilities. Peak memory is bounded by a single forward through the attractor module and standard KV-caching applies in the backbone. In principle, the solver tolerance  $\varepsilon$  and maximum iteration budget  $T_{\max}$  are inference-time hyperparameters: they can be adjusted without retraining the model, turning test-time computation into a budget for approaching the learned attractor. Interestingly, however, we find that trained Attractor Models often require very little test-time refinement, as we describe in the next section.

### 3.3 Equilibrium Internalization and Stability in Attractor Models

Although Attractor Models define predictions through the equilibrium  $\tilde{y}^*$ , we observe a surprising phenomenon: after training, the backbone proposal  $\tilde{y}_0$  often lies close to the equilibrium (c.f. Figures 6 and 7). We refer to this phenomenon as *equilibrium internalization*. Intuitively, the attractor module appears to act as a moving teacher for the backbone, resulting in a form of automatic curriculum. Early in training, the proposal  $\tilde{y}_0$  may

be far from a good prediction, and the solver must perform nontrivial refinement to reach  $\tilde{y}^*$ . Since  $\tilde{y}_0$  and  $\tilde{y}^*$  live in the same tied output-embedding space, gradients through the equilibrium also train the backbone proposal to move toward the state that the solver would have found. Thus, when the backbone is sufficiently expressive, much of the prediction work can be internalized into  $\tilde{y}_0$ , leaving the attractor to perform a stable refinement.

**Stability.** Equilibrium training also biases the recurrent map toward convergent dynamics. The implicit gradient contains the inverse factor  $(I - J_{\tilde{y}}^T)^{-1}$ , which becomes ill-conditioned near non-contractive regimes. This creates a barrier against unstable fixed-point dynamics, unlike fixed-loop training, which can learn trajectories that are accurate only at a prescribed step count and fail under extra inference-time loops. We discuss this contrast in detail in Appendix B.3. We refer to Appendix B for theoretical analysis of Attractor Models and comparison with finite-loop models.

## 4 Experiments

We evaluate Attractor Models in two regimes. We first study language modeling across model sizes, comparing against parameter-matched Transformers and looped-LM baselines. We evaluate scaling behavior, downstream accuracy, and training efficiency. We also present results on hard reasoning tasks, where we test whether the same fixed-point refinement mechanism improves small models on problems that require iterative computation.

### 4.1 Attractor Models Improve Large-Scale Language Modeling

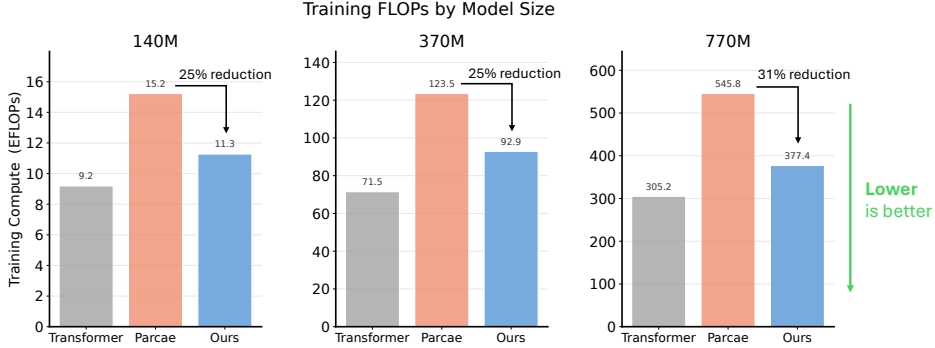
**Setup.** We follow the nanochat (Karpathy, 2025) pretraining recipe used by Parcae (Prairie et al., 2026) for its main Transformer comparison, training on FineWeb-Edu (Penedo et al., 2024). To ensure a fair comparison, all models are matched in parameter count and trained with the same data budget, optimizer, and learning-rate schedule as the Parcae baselines; the only change is the recurrent block. We compare at three scales: 140M, 370M, and 770M parameters. Parcae is a middle-looped language model with prelude, coda, and recurrent blocks. The stability in this model comes from a linear injection that bounds the spectral radius of the recurrence below one. In contrast, our architecture uses standard Transformer blocks followed by a fixed-point iteration block, and lets the solver itself control the effective depth. Detailed hyperparameter settings are in Appendix C.

**Parameter Scaling.** We demonstrate how large-scale pretraining scales with our model. We evaluate our method on 140M, 370M, and 770M parameters against a parameter-matched Transformer and Parcae, a looped-language model (Prairie et al., 2026). Our model improves monotonically with scale. Across all three sizes, our method achieves the best validation PPL, Lambada PPL, and CORE accuracy. These results show that our fixed-point refinement scales cleanly with model size, with especially large gains on Lambada, where iterative refinement substantially improves long-context prediction.

**Table 1 Parameter scaling results on large-scale language modeling.** Our Attractor Model outperforms standard Transformers and looped model PARCAE on nearly every benchmark. Our 770M model performs comparably to a standard Transformer with nearly twice as many parameters and trained on twice as many tokens. Arrows indicate the relative improvement over the parameter-matched standard Transformer baseline. Arrows indicate the relative improvement over the parameter-matched Transformer baseline.

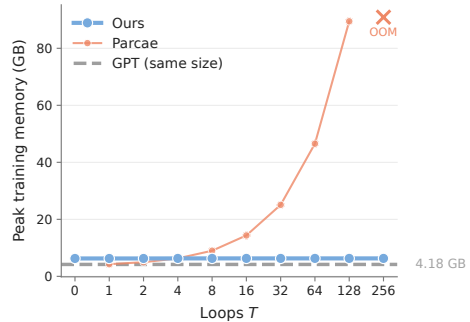
Size	Model	Val. PPL ↓	Lambada PPL ↓	Core ↑	Core-Ext. ↑
<b>140M</b>	Transformer	21.48	127.39	13.00 ± 0.15	8.80 ± 0.21
	PARCAE	19.06	80.64	14.04 ± 0.20	9.67 ± 0.28
	<b>Attractor Model (ours)</b>	<b>18.30</b> ↓14.8%	<b>68.02</b> ↓46.6%	<b>14.59 ± 0.11</b> ↑12.2%	<b>10.03 ± 0.05</b> ↑14.0%
<b>370M</b>	Transformer	15.79	40.77	17.46 ± 0.03	11.71 ± 0.22
	PARCAE	14.49	32.74	20.00 ± 0.06	<b>12.75 ± 0.31</b>
	<b>Attractor Model (ours)</b>	<b>14.03</b> ↓11.1%	<b>27.14</b> ↓33.4%	<b>20.24 ± 0.09</b> ↑15.9%	12.64 ± 0.33 ↑7.9%
<b>770M</b>	Transformer	13.08	22.37	22.42 ± 0.20	14.20 ± 0.63
	PARCAE	12.49	19.71	25.07 ± 0.33	15.19 ± 0.43
	<b>Attractor Model (ours)</b>	<b>12.09</b> ↓7.6%	<b>15.21</b> ↓32.0%	<b>26.83 ± 0.29</b> ↑19.7%	<b>15.42 ± 0.51</b> ↑8.6%
<b>1.3B</b>	Transformer	11.95	17.26	25.45 ± 0.08	15.90 ± 0.23

**Training efficiency: Lower FLOPs.** The one-step IFT backward pass keeps training memory constant in the number of solver iterations. The memory of standard looped language models like (Prairie et al., 2026) scales linearly with the number of loops. The total training FLOPs (Figure 4) follow the same trend: although every step of our recurrent block costs roughly the same as Parcae’s, our solver typically converges below  $\varepsilon$  in well under  $T_{\max}$  steps, so the realized depth during training is lower despite identical  $T_{\max}$ , yielding 25–31% lower training FLOPs across scales.



**Figure 4 Training-time efficiency.** Despite being a looped architecture, our method uses 25–31% fewer FLOPs than Parcae because the fixed-point solver often converges before  $T_{\max}$  and the backward pass uses the one-step implicit-gradient approximation.

**Training efficiency: O(1) memory.** An additional advantage of our method is that the training memory required stays constant with the number of iterations (Figure 5). This is because our implicit backward pass does not need to store the intermediate activations from every recurrent step. In contrast, standard looped language models must backpropagate through each unrolled iteration, causing memory usage to grow linearly with the number of loops.



**Figure 5 Peak training memory versus recurrent depth.** Attractor Models keep memory nearly constant, while Parcae’s memory grows with the number of loops.

**Result: Attractor Models deliver Pareto improvement on large-scale language modeling.**

Across model sizes, Attractor Models outperform parameter-matched Transformers and stable looped baselines on perplexity and downstream accuracy, while using substantially less training compute and memory than looped models. Notably, the 770M Attractor Model reaches the quality regime of a 1.3B Transformer trained on twice as many tokens.

## 4.2 Attractor Models Lead to Significant Gains in Hard Reasoning Tasks

Beyond language modeling, we train and evaluate Attractor Models on Sudoku-Extreme and Maze-Hard from (Wang et al., 2025; Jolicoeur-Martineau, 2025): two challenging reasoning benchmarks, where non-recurrent Transformers and most frontier LLMs fail.

**Setup.** We train small models with approximately 1000 training examples for each task and require predicting the full output grid in a single direct forward pass (no autoregressive decoding). We follow the TRM training protocol (Jolicoeur-Martineau, 2025), using deep-supervision steps.

**Method.** TRM carries two latents across deep-supervision steps, a current answer  $y$  and a reasoning state  $z$ , and applies a tiny two-layer network  $T(n+1)$  times per step to update them. We keep this protocol except for

the inner update. Specifically, instead of unrolling  $T(n+1)$  applications, we solve directly for the fixed point of the  $(y, z)$  update with our solver. The initialization is handled by deep supervision itself: the previous step’s  $(y, z)$  initializes the solver at the next step, with a learned embedding at step zero, so we do not use a separate backbone.

We present our results in Table 2. The fixed-depth Transformer fails on both tasks. HRM (27M) achieves 55.0% and 74.5%, respectively on Sudoku-Extreme and Maze-Hard. TRM is the strongest tiny baseline at 7M, achieving 74.7% and 85.3%, but (counter to the goal of scaling) collapses to 0% on both tasks when scaled to 27M parameters. Our model scales naturally with parameter count. We attribute the difference to the explicit fixed-point objective, which appears to provide regularization that bare iterative refinement lacks at higher capacity.

For the backward pass, we use the phantom-gradient scheme (Geng et al., 2022) rather than the one-step approximation used in the language-modeling experiments. This choice is important in the small-data reasoning regime: with only  $\sim 1,000$  training examples and much smaller networks, the solver dynamics are more sensitive, and the one-step surrogate can provide too crude a training signal. This is consistent with TRM, which reports that replacing its backward pass with a one-step approximation reduces Sudoku-Extreme accuracy from 87.4% to 56.5% (Jolicoeur-Martineau, 2025).

**Table 2 Small model and sample training results on Sudoku-Extreme and Maze-Hard.** Our Attractor Models improve when scaling the parameter count, compared to standard tiny recursive models which degrade.

Method	# Parameters	Sudoku-Extreme $\uparrow$	Maze-Hard $\uparrow$
Deepseek R1	671B	0.0%	0.0%
Claude 3.7	?	0.0%	0.0%
O3-mini-high	?	0.0%	0.0%
Transformer	27M	0.0%	0.0%
HRM	27M	55.0%	74.5%
TRM	7M	74.7%	85.3%
TRM	27M	0.0%	0.0%
<b>Attractor Model (ours)</b>	7M	54.3%	46.7%
<b>Attractor Model (ours)</b>	27M	<b>91.4%</b>	<b>93.1%</b>

This setup is still an instance of our Attractor Model framework, where an output-space representation is iteratively refined to an equilibrium, decoded, and differentiated through implicitly. The only difference is the initialization mechanism. Rather than a Transformer backbone producing  $\tilde{y}_0$  from the input, deep supervision supplies the initialization. Specifically, the previous supervision step’s  $(y, z)$  initializes the solver at the next step, with a learned embedding at step zero.

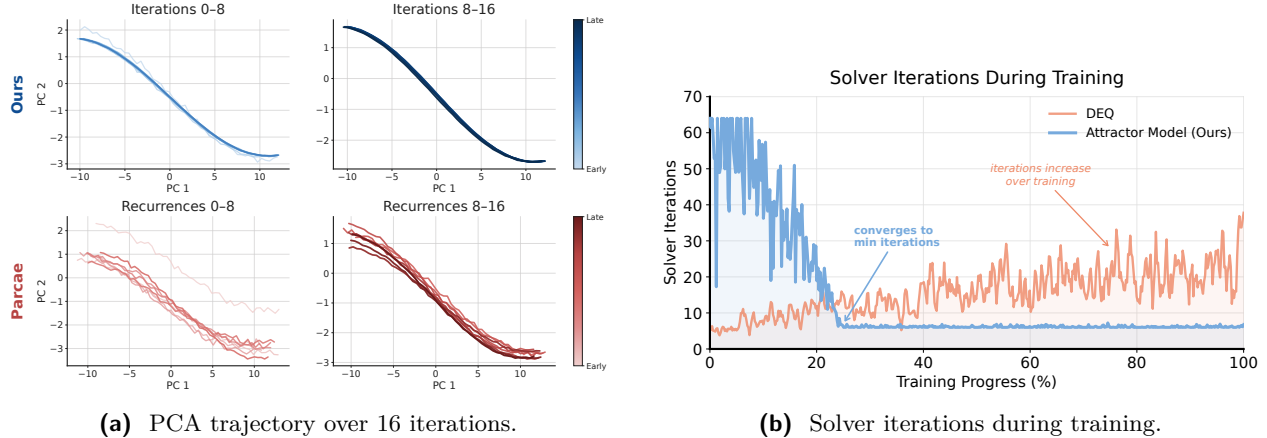
**Result: Attractor Models scale where recursive reasoners collapse.**

In hard reasoning tasks, existing recursive architectures can degrade as model capacity increases. Attractor Models avoid this failure mode. With only 27M parameters and  $\sim 1,000$  training examples, our model reaches 91.4% accuracy on Sudoku-Extreme and 93.1% on Maze-Hard, outperforming both frontier LLMs and specialized recursive baselines.

### 4.3 Equilibrium Internalization and Test-Time Behavior

**Fixed-point convergence.** Figure 6 visualizes convergence in two complementary ways. First, we project the trajectory of the final-position representation onto its first two principal components over 16 iterations. Our method rapidly contracts to a fixed point: iterations 8–16 collapse onto a single attractor. Parcae also moves toward a stable state, but does so more slowly; its trajectory remains noisier and continues to drift through later recurrences. This suggests that while finite-loop language models can exhibit fixed-point-like behavior, explicitly training the loop as a fixed-point problem produces faster and cleaner convergence. Second, we track the number of solver iterations required during training. The DEQ baseline requires increasingly many iterations as optimization progresses, consistent with the observations in the original work (Bai et al., 2019).

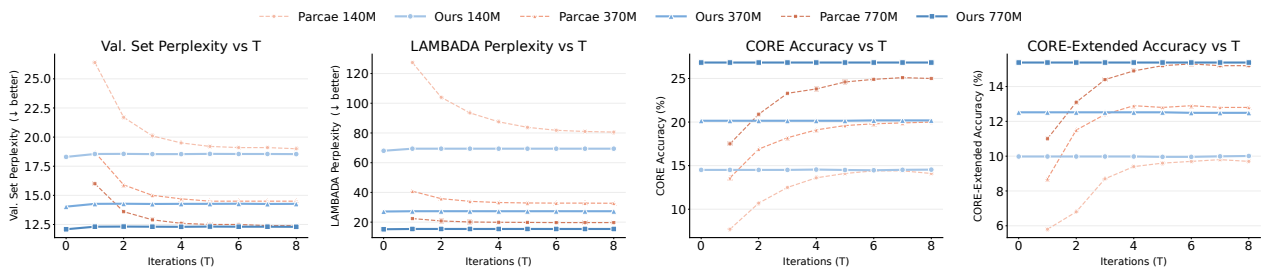
In contrast, Attractor Models rapidly converge to the minimum iteration count and remain stable. This behavior provides evidence for equilibrium internalization, where optimization shifts work from the iterative solver into the backbone proposal.



**Figure 6 Convergence behavior of Attractor Models compared to existing methods.** **Left:** PCA projection of the residual stream at the final sequence position over 16 iterations. **Right:** DEQ baseline requires increasingly many solver iterations during training, whereas Attractor Models rapidly settles to the minimum iteration count and remains stable.

**Test-time iterations vs. quality.** In Figure 7, we report validation perplexity, CORE accuracy, and CORE-Extended accuracy as we vary the number of inference iterations  $T$  while holding training fixed. For Parcae, quality improves monotonically from  $T = 1$  until it plateaus near  $T = 8$ , indicating that the model relies on repeated test-time refinement. In contrast, our method reaches peak performance at  $T = 1$  at every scale, and even  $T = 0$  (decoding the backbone proposal  $\tilde{y}_0$  directly through the tied unembedding) is already at or near the converged value.

We attribute this behavior to *equilibrium internalization*: during training, the backbone learns to produce a proposal  $\tilde{y}_0$  that already lies close to the fixed point that the solver would otherwise compute through iteration. Thus, the iterative block still shapes the learned representation during training, but at inference time the model has largely internalized the result of that refinement into its backbone initialization. Strikingly, the backbone-only readout at  $T = 0$  is already stronger than larger standard Transformers trained without attractor assistance, and matches or exceeds finite-loop baselines that require many test-time recurrences. Concretely, our 140M model at  $T = 0$  (no solver) matches or improves Parcae 140M at  $T = 8$  recurrent steps; the same pattern holds at 370M and 770M.



**Figure 7 Validation perplexity on FineWeb-Edu, Lambada, CORE accuracy, and CORE-Extended accuracy as a function of test-time iterations  $T$ .** Parcae’s quality improves monotonically up to  $T \approx \mu_{rec} = 8$ . Our method is essentially converged at  $T = 1$  (and at  $T = 0$  for the 770M model), because the solver is initialized at  $\tilde{y}_0$ , which is itself a coherent prediction.

Takeaway: Attractor Models internalize equilibrium.

Attractor Models learn to initialize a prediction near the fixed point, so inference is effectively converged before any explicit refinement step. As a result, they preserve the benefits of iterative computation during training while avoiding the need for many sequential iterations at test time, unlike standard looped models whose quality depends on repeated inference-time updates.

## 4.4 Ablations

For the ablations, we use a 60.3M-parameter version of our model trained on 1B tokens of FineWeb-Edu under the same `nanochat` configuration. All ablation runs share the same data stream, tokenizer, and optimizer; only the ablated component varies.

We isolate the contribution of (i) the location of the equilibrium and its initialization (Table 3), (ii) the additive injection of  $c$  (Table 4), (iii) and the one-step backward approximation (Table 5). All ablations train a 60.3M-parameter model on 1B tokens of FineWeb-Edu under the same `nanochat` setup, with all other hyperparameters fixed.

**Comparison to DEQ.** Table 3 compares Attractor Models against a parameter-matched DEQ (Bai et al., 2019). The DEQ baseline solves for a hidden-state equilibrium initialized from an uninformative state and decodes it with a separate output head. Tying the unembedding closes part of the gap (Val. PPL 42.18  $\rightarrow$  38.74), but the largest improvement comes from matching the equilibrium to the design of Attractor Models: the fixed point is placed directly in the tied output-embedding space and the root finder is initialized from the backbone proposal  $\tilde{y}_0 = \mathcal{T}_{\theta_b}(E(x))$  ( $\rightarrow$  34.05).

This proposal gives the solver a meaningful, input-dependent starting point rather than requiring the recurrent block to construct the representation from scratch. As a result, the Attractor Model reaches the same residual tolerance in  $1.7\times$  fewer iterations while also achieving lower perplexity. We view this as evidence of *equilibrium internalization*: the backbone learns to produce an output-side proposal  $\tilde{y}_0$  that already lies close to the eventual equilibrium  $\tilde{y}^*$ , making the subsequent attractor refinement both easier and faster.

Table 3 Comparison to DEQ at matched parameter count.

Method	Size	Equilibrium	Val. PPL $\downarrow$	Avg. Iters $\downarrow$	Core $\uparrow$
DEQ (Bai et al., 2019)	60.3M	hidden ( $z_0 = 0$ , sep. head)	42.18	14.6	$5.21 \pm 0.14$
DEQ + tied unemb.	60.3M	hidden ( $z_0 = 0$ )	38.74	13.9	$5.83 \pm 0.12$
<b>Attractor Model</b>	60.3M	output emb. ( $\tilde{y}_0 = \mathcal{T}_{\theta_b}(E(x))$ )	<b>34.05</b>	<b>8.4</b>	<b><math>6.74 \pm 0.10</math></b>

**Proposal injection.** Table 4 ablates how the backbone proposal  $\tilde{y}_0$  is provided to the attractor. We use *proposal injection* to denote how  $\tilde{y}_0$  enters the recurrent refinement map, and *persistent injection* to denote providing  $\tilde{y}_0$  at every refinement step rather than only through the initial state.

When the proposal is used only for initialization,  $\tilde{y}_{t=0} = \tilde{y}_0$  and  $\tilde{y}_{t+1} = \mathcal{T}_{\theta_a}(\tilde{y}_t)$ , the recurrent update no longer depends on the input after the first state. Consequently, the fixed point can become proposal-independent, and only 12.4% of validation tokens converge within  $T_{\max}$ . Persistent injection by concatenation,  $\tilde{y}_{t+1} = \mathcal{T}_{\theta_a}([\tilde{y}_t; \tilde{y}_0])$ , restores proposal-dependence and recovers much of the quality, but makes the refinement problem harder: convergence is slower (11.2 vs. 8.4 average iterations) and perplexity is worse (36.81 vs. 34.05). Additive proposal injection,  $\tilde{y}_{t+1} = \mathcal{T}_{\theta_a}(\tilde{y}_t, \tilde{y}_0)$ , provides the backbone proposal at every step while keeping the refinement map simple and well-conditioned, yielding the best results across all three metrics.

Table 4 Effect of input injection on convergence and quality.

Injection of $c$	Val. PPL $\downarrow$	Avg. Iters $\downarrow$	% Converged $\uparrow$
Initial only ( $\tilde{y}_0$ , no re-injection)	51.92	$T_{\max}$	12.4%
Concatenation: $\mathcal{A}_{\theta_a}([\tilde{y}_t; \tilde{y}_0])$	36.81	11.2	88.6%
<b>Additive (ours): <math>\mathcal{A}_{\theta_a}(\tilde{y}_t + \tilde{y}_0)</math></b>	<b>34.05</b>	<b>8.4</b>	<b>99.7%</b>

**Backward pass.** In table 5, we compare the one-step backward approximation against full implicit differentiation (Anderson on the linear system  $(I - J_{\tilde{y}}^T)u = v$ ) and the phantom gradient (Geng et al., 2022) unroll. Full IFT improves PPL by only 0.14 while increasing peak training memory by  $4.8\times$  and step time by  $2.7\times$ ; phantom gradient lies in between. The one-step approximation allows relatively cheap training cost while maintaining nearly all of the original performance of a full backward gradient computation.

**Table 5 Ablation for running the full IFT gradients, phantom gradients, and the one-step approximation.**

Backward Pass	Val. PPL ↓	Train Mem. ↓	Step Time ↓	Δ vs. Full
Full IFT (Anderson on $u$ )	33.91	$4.8\times$	$2.7\times$	—
Phantom gradient ( $k=3$ ) (Geng et al., 2022)	34.02	$1.8\times$	$1.4\times$	+0.11
<b>Ours: one-step (<math>u \approx v</math>) (Fung et al., 2021)</b>	<b>34.05</b>	<b><math>1.0\times</math></b>	<b><math>1.0\times</math></b>	+0.14

Takeaway: One-step gradient approximation is sufficient for Attractor Models

Using the one-step gradient approximation enables efficient training, while maintaining nearly all of the quality of using full backward gradients.

**Initialization ablation.** Table 6 isolates the effect of the solver initialization. Starting the attractor solve from an uninformative state, either zero or Gaussian noise, substantially degrades both quality and convergence. In contrast, initializing from the backbone proposal  $\tilde{y}_0 = \mathcal{T}_{\theta_b}(E(x))$  gives the solver a meaningful output-side starting point, yielding lower perplexity, fewer iterations, and higher downstream accuracy. This supports our hypothesis that the refinement should begin near a coherent prediction embedding rather than constructing one from scratch.

**Table 6 Effect of attractor solver initialization.** Initializing from the backbone proposal  $\tilde{y}_0$  clearly outperforms uninformative zero or Gaussian noise initialization.

Initialization	Val. PPL ↓	Avg. Iters ↓	% Converged ↑	Core ↑
Zero: $\tilde{y}_{\text{init}} = 0$	43.87	14.8	71.3%	$5.42 \pm 0.13$
Gaussian: $\tilde{y}_{\text{init}} \sim \mathcal{N}(0, \sigma^2 I)$	41.26	13.6	78.9%	$5.71 \pm 0.11$
<b>Backbone proposal: <math>\tilde{y}_{\text{init}} = \tilde{y}_0</math></b>	<b>34.05</b>	<b>8.4</b>	<b>99.7%</b>	<b><math>6.74 \pm 0.10</math></b>

## 5 Conclusion and Future Work

In this work, we propose Attractor Models, a new family of architectures that first produce meaningful prediction embeddings and then refine them through an attractor module by solving for a fixed point. This formulation makes recurrent refinement stable, memory-efficient, and adaptive, while avoiding the cost of explicit unrolling and achieving strong results across both large-scale language modeling and hard reasoning with tiny models. In stark contrast to prior looped language models, training Attractor Models gives rise to equilibrium internalization: the model learns to make the very refinement procedure that trained it largely unnecessary at inference time. Interesting directions for future work is to further study the equilibrium internalization phenomenon and the differences between Attractor Models and finite-loop recurrence.

**Acknowledgements.** The authors gratefully acknowledge generous support from Coefficient Giving.

## References

- Donald G. Anderson. Iterative procedures for nonlinear integral equations. *Journal of the ACM*, 12(4):547–560, 1965. doi: 10.1145/321296.321305.
- Anthropic. The claude 3 model family: Opus, Sonnet, Haiku. Model card, 2024. [https://www-cdn.anthropic.com/de8ba9b01c9ab7cbabf5c33b80b7bbc618857627/Model\\_Card\\_Claude\\_3.pdf](https://www-cdn.anthropic.com/de8ba9b01c9ab7cbabf5c33b80b7bbc618857627/Model_Card_Claude_3.pdf).

- Sangmin Bae, Adam Fisch, Hrayr Harutyunyan, Ziwei Ji, Seungyeon Kim, and Tal Schuster. Relaxed recursive transformers: Effective parameter sharing with layer-wise lora, 2025a. <https://arxiv.org/abs/2410.20672>.
- Sangmin Bae, Yujin Kim, Reza Bayat, Sungnyun Kim, Jiyouon Ha, Tal Schuster, Adam Fisch, Hrayr Harutyunyan, Ziwei Ji, Aaron Courville, and Se-Young Yun. Mixture-of-recursions: Learning dynamic recursive depths for adaptive token-level computation, 2025b. <https://arxiv.org/abs/2507.10524>.
- Shaojie Bai, J. Zico Kolter, and Vladlen Koltun. Deep equilibrium models. In *Advances in Neural Information Processing Systems 32 (NeurIPS 2019)*, 2019.
- Shaojie Bai, Vladlen Koltun, and J. Zico Kolter. Multiscale deep equilibrium models. In *Advances in Neural Information Processing Systems*, volume 33. Curran Associates, Inc., 2020. <https://proceedings.neurips.cc/paper/2020/hash/3812f9a59b634c2a9c574610eaba5bed-Abstract.html>.
- Shaojie Bai, Vladlen Koltun, and J. Zico Kolter. Stabilizing equilibrium models by jacobian regularization, 2021. <https://arxiv.org/abs/2106.14342>.
- Hugh Blayney, Álvaro Arroyo, Johan Obando-Ceron, Pablo Samuel Castro, Aaron Courville, Michael M. Bronstein, and Xiaowen Dong. A mechanistic analysis of looped reasoning language models, 2026. <https://arxiv.org/abs/2604.11791>.
- Chris Cameron, Wangzheng Wang, Nikita Ivanov, Ashmita Bhattacharyya, Didier Chételat, and Yingxue Zhang. One step forward and k steps back: Better reasoning with denoising recursion models, 2026. <https://arxiv.org/abs/2604.18839>.
- Yilong Chen, Junyuan Shang, Zhenyu Zhang, Yanxi Xie, Jiawei Sheng, Tingwen Liu, Shuohuan Wang, Yu Sun, Hua Wu, and Haifeng Wang. Inner thinking transformer: Leveraging dynamic depth scaling to foster adaptive internal thinking, 2025. <https://arxiv.org/abs/2502.13842>.
- Raj Dabre and Atsushi Fujita. Recurrent stacking of layers for compact neural machine translation models, 2018. <https://arxiv.org/abs/1807.05353>.
- Mostafa Dehghani, Stephan Gouws, Oriol Vinyals, Jakob Uszkoreit, and Łukasz Kaiser. Universal transformers, 2019. <https://arxiv.org/abs/1807.03819>.
- Jingcheng Deng, Liang Pang, Zihao Wei, Shichen Xu, Zenghao Duan, Kun Xu, Yang Song, Huawei Shen, and Xueqi Cheng. Latent reasoning in LLMs as a vocabulary-space superposition. *arXiv preprint arXiv:2510.15522*, 2025.
- Jingcheng Deng, Zihao Wei, Liang Pang, Junhong Wu, Shicheng Xu, Zenghao Duan, and Huawei Shen. Latent-GRPO: Group relative policy optimization for latent reasoning. *arXiv preprint arXiv:2604.27998*, 2026.
- Ying Fan, Yilun Du, Kannan Ramchandran, and Kangwook Lee. Looped transformers for length generalization, 2025. <https://arxiv.org/abs/2409.15647>.
- François Fleuret. The free transformer, 2025. <https://arxiv.org/abs/2510.17558>.
- Tianyu Fu, Yichen You, Zekai Chen, Guohao Dai, Huazhong Yang, and Yu Wang. Think-at-hard: Selective latent iterations to improve reasoning language models, 2026. <https://arxiv.org/abs/2511.08577>.
- Samy Wu Fung, Howard Heaton, Qiuwei Li, Daniel McKenzie, Stanley Osher, and Wotao Yin. Jfb: Jacobian-free backpropagation for implicit networks, 2021. <https://arxiv.org/abs/2103.12803>.
- Khashayar Gatmiry, Nikunj Saunshi, Sashank J Reddi, Stefanie Jegelka, and Sanjiv Kumar. Can looped transformers learn to implement multi-step gradient descent for in-context learning? In *International Conference on Machine Learning*, pages 15130–15152. PMLR, 2024a.
- Khashayar Gatmiry, Nikunj Saunshi, Sashank J. Reddi, Stefanie Jegelka, and Sanjiv Kumar. On the role of depth and looping for in-context learning with task diversity, 2024b. <https://arxiv.org/abs/2410.21698>.
- Jonas Geiping, Sean McLeish, Neel Jain, John Kirchenbauer, Siddharth Singh, Brian R. Bartoldson, Bhavya Kailkhura, Abhinav Bhattele, and Tom Goldstein. Scaling up test-time compute with latent reasoning: A recurrent depth approach, 2025a. <https://arxiv.org/abs/2502.05171>.
- Jonas Geiping, Xinyu Yang, and Guinan Su. Efficient parallel samplers for recurrent-depth models and their connection to diffusion language models, 2025b. <https://arxiv.org/abs/2510.14961>.
- Gemini Team. Gemini: A family of highly capable multimodal models. *arXiv preprint arXiv:2312.11805*, 2023.

- Zhengyang Geng, Xin-Yu Zhang, Shaojie Bai, Yisen Wang, and Zhouchen Lin. On training implicit models, 2022. <https://arxiv.org/abs/2111.05177>.
- Raj Ghugare, Michał Bortkiewicz, Alicja Ziarko, and Benjamin Eysenbach. On the role of iterative computation in reinforcement learning, 2026. <https://arxiv.org/abs/2602.05999>.
- Angeliki Giannou, Shashank Rajput, Jy-yong Sohn, Kangwook Lee, Jason D Lee, and Dimitris Papailiopoulos. Looped transformers as programmable computers. In *International Conference on Machine Learning*, pages 11398–11442. PMLR, 2023a.
- Angeliki Giannou, Shashank Rajput, Jy yong Sohn, Kangwook Lee, Jason D. Lee, and Dimitris Papailiopoulos. Looped transformers as programmable computers, 2023b. <https://arxiv.org/abs/2301.13196>.
- Sahil Goyal, Swayam Agrawal, Gautham Govind Anil, Prateek Jain, Sujoy Paul, and Aditya Kusupati. Elt: Elastic looped transformers for visual generation, 2026. <https://arxiv.org/abs/2604.09168>.
- Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, et al. The Llama 3 herd of models. In *Neural Information Processing Systems*, 2024.
- Alex Graves. Adaptive computation time for recurrent neural networks, 2017. <https://arxiv.org/abs/1603.08983>.
- Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Peiyi Wang, Qihao Zhu, Runxin Xu, Ruoyu Zhang, Shirong Ma, Xiao Bi, Xiaokang Zhang, Xingkai Yu, Yu Wu, Z. F. Wu, Zhibin Gou, Zhihong Shao, Zhuoshu Li, Ziyi Gao, Aixin Liu, Bing Xue, Bingxuan Wang, Bochao Wu, Bei Feng, Chengda Lu, Chenggang Zhao, Chengqi Deng, Chong Ruan, Damai Dai, Deli Chen, Dongjie Ji, Erhang Li, Fangyun Lin, Fucong Dai, Fuli Luo, Guangbo Hao, Guanting Chen, Guowei Li, H. Zhang, Hanwei Xu, Honghui Ding, Huazuo Gao, Hui Qu, Hui Li, Jianzhong Guo, Jiashi Li, Jingchang Chen, Jingyang Yuan, Jinhao Tu, Junjie Qiu, Junlong Li, J. L. Cai, Jiaqi Ni, Jian Liang, Jin Chen, Kai Dong, Kai Hu, Kaichao You, Kaige Gao, Kang Guan, Kexin Huang, Kuai Yu, Lean Wang, Lecong Zhang, Liang Zhao, Litong Wang, Liyue Zhang, Lei Xu, Leyi Xia, Mingchuan Zhang, Minghua Zhang, Minghui Tang, Mingxu Zhou, Meng Li, Miaojun Wang, Mingming Li, Ning Tian, Panpan Huang, Peng Zhang, Qiancheng Wang, Qinyu Chen, Qiushi Du, Ruiqi Ge, Ruisong Zhang, Ruizhe Pan, Runji Wang, R. J. Chen, R. L. Jin, Ruyi Chen, Shanghao Lu, Shangyan Zhou, Shanhuang Chen, Shengfeng Ye, Shiyu Wang, Shuiping Yu, Shunfeng Zhou, Shuting Pan, S. S. Li, Shuang Zhou, Shaoqing Wu, Tao Yun, Tian Pei, Tianyu Sun, T. Wang, Wangding Zeng, Wen Liu, Wenfeng Liang, Wenjun Gao, Wenqin Yu, Wentao Zhang, W. L. Xiao, Wei An, Xiaodong Liu, Xiaohan Wang, Xiaokang Chen, Xiaotao Nie, Xin Cheng, Xin Liu, Xin Xie, Xingchao Liu, Xinyu Yang, Xinyuan Li, Xuecheng Su, Xuheng Lin, X. Q. Li, Xiangyue Jin, Xiaojin Shen, Xiaosha Chen, Xiaowen Sun, Xiaoxiang Wang, Xinnan Song, Xinyi Zhou, Xianzu Wang, Xinxia Shan, Y. K. Li, Y. Q. Wang, Y. X. Wei, Yang Zhang, Yanhong Xu, Yao Li, Yao Zhao, Yaofeng Sun, Yaohui Wang, Yi Yu, Yichao Zhang, Yifan Shi, Yiliang Xiong, Ying He, Yishi Piao, Yisong Wang, Yixuan Tan, Yiyang Ma, Yiyuan Liu, Yongqiang Guo, Yuan Ou, Yuduan Wang, Yue Gong, Yuheng Zou, Yujia He, Yunfan Xiong, Yuxiang Luo, Yuxiang You, Yuxuan Liu, Yuyang Zhou, Y. X. Zhu, Yanping Huang, Yaohui Li, Yi Zheng, Yuchen Zhu, Yunxian Ma, Ying Tang, Yukun Zha, Yuting Yan, Z. Z. Ren, Zehui Ren, Zhangli Sha, Zhe Fu, Zhean Xu, Zhenda Xie, Zhengyan Zhang, Zhewen Hao, Zhicheng Ma, Zhigang Yan, Zhiyu Wu, Zihui Gu, Zijia Zhu, Zijun Liu, Zilin Li, Ziwei Xie, Ziyang Song, Zizheng Pan, Zhen Huang, Zhipeng Xu, Zhongyu Zhang, and Zhen Zhang. Deepseek-r1 incentivizes reasoning in llms through reinforcement learning. *Nature*, 645(8081):633–638, 2025. ISSN 1476-4687. doi: 10.1038/s41586-025-09422-z. <http://dx.doi.org/10.1038/s41586-025-09422-z>.
- Sang-Il Han. Hierarchical vs. flat iteration in shared-weight transformers, 2026. <https://arxiv.org/abs/2604.14442>.
- Shibo Hao, Sainbayar Sukhbaatar, DiJia Su, Xian Li, Zhiting Hu, Jason Weston, and Yuandong Tian. Training large language models to reason in a continuous latent space, 2025. <https://arxiv.org/abs/2412.06769>.
- Edward S. Hu, Kwangjun Ahn, Qinghua Liu, Haoran Xu, Manan Tomar, Ada Langford, Jayden Teoh, Bryon Xu, David Yan, Dinesh Jayaraman, Alex Lamb, and John Langford. The belief state transformer, 2025. <https://arxiv.org/abs/2410.23506>.
- Jianhao Huang, Zixuan Wang, and Jason D. Lee. Transformers learn to implement multi-step gradient descent with chain of thought, 2025. <https://arxiv.org/abs/2502.21212>.
- Hakan Inan, Khashayar Khosravi, and Richard Socher. Tying word vectors and word classifiers: A loss framework for language modeling, 2017. <https://arxiv.org/abs/1611.01462>.
- Ahmadreza Jeddi, Marco Ciccone, and Babak Taati. Loopformer: Elastic-depth looped transformers for latent reasoning via shortcut modulation, 2026. <https://arxiv.org/abs/2602.11451>.

Alexia Jolicoeur-Martineau. Less is more: Recursive reasoning with tiny networks, 2025. <https://arxiv.org/abs/2510.04871>.

Andrej Karpathy. nanochat: The best chatgpt that \$100 can buy, 2025. <https://github.com/karpathy/nanochat>.

Jonas Knupp, Jan Hendrik Metzen, Jeremias Bohn, Georg Groh, and Kristian Kersting. Depth-Recurrent Attention Mixtures: Giving latent reasoning the attention it deserves. *arXiv preprint arXiv:2601.21582*, 2026a.

Jonas Knupp, Jan Hendrik Metzen, Jeremias Bohn, Georg Groh, and Kristian Kersting. Depth-recurrent attention mixtures: Giving latent reasoning the attention it deserves, 2026b. <https://arxiv.org/abs/2601.21582>.

Harsh Kohli, Srinivasan Parthasarathy, Huan Sun, and Yuekun Yao. Loop, think, & generalize: Implicit reasoning in recurrent-depth transformers, 2026. <https://arxiv.org/abs/2604.07822>.

Mieszko Komisarczyk, Saurabh Mathur, Maurice Kraus, Sriraam Natarajan, and Kristian Kersting. Recursive inference machines for neural reasoning, 2026. <https://arxiv.org/abs/2603.05234>.

Steven G. Krantz and Harold R. Parks. *The Implicit Function Theorem: History, Theory, and Applications*. Birkhäuser, Boston, MA, 2002a. ISBN 978-0-8176-4285-3. doi: 10.1007/978-1-4612-0059-8.

Steven George Krantz and Harold R Parks. *The implicit function theorem: history, theory, and applications*, volume 202. Springer, 2002b.

Asher Labovich. Stability and generalization in looped transformers, 2026a. <https://arxiv.org/abs/2604.15259>.

Asher Labovich. Stability and generalization in looped transformers, 2026b. <https://arxiv.org/abs/2604.15259>.

Zhenzhong Lan, Mingda Chen, Sebastian Goodman, Kevin Gimpel, Piyush Sharma, and Radu Soricut. Albert: A lite bert for self-supervised learning of language representations, 2020. <https://arxiv.org/abs/1909.11942>.

Sean Michael McLeish, Ang Li, John Kirchenbauer, Dayal Singh Kalra, Brian R Bartoldson, Bhavya Kailkhura, Avi Schwarzschild, Jonas Geiping, Micah Goldblum, and Tom Goldstein. Teaching pretrained language models to think deeper with retrofitted recurrence. In *NeurIPS 2025 Workshop on Efficient Reasoning*.

Will Merrill and Ashish Sabharwal. A little depth goes a long way: The expressive power of log-depth transformers. *Advances in Neural Information Processing Systems*, 38:95315–95339, 2026.

William Merrill and Ashish Sabharwal. A little depth goes a long way: The expressive power of log-depth transformers, 2025. <https://arxiv.org/abs/2503.03961>.

Amirkeivan Mohtashami, Matteo Pagliardini, and Martin Jaggi. CoTFormer: More tokens with attention make up for less depth. In *Workshop on Advancing Neural Network Training: Computational Efficiency, Scalability, and Resource Optimization (WANT@ NeurIPS 2023)*, 2023.

Ibraheem Muhammad Moosa, Suhas Lohit, Ye Wang, Moitreyia Chatterjee, and Wenpeng Yin. Understanding dynamic compute allocation in recurrent transformers, 2026. <https://arxiv.org/abs/2602.08864>.

OpenAI. GPT-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023.

Enes Özeren and Matthias Aßenmacher. Reinforcement learning for latent-space thinking in LLMs. *arXiv preprint arXiv:2512.11816*, 2025.

Denis Paperno, Germán Kruszewski, Angeliki Lazaridou, Ngoc Quan Pham, Raffaella Bernardi, Sandro Pezzelle, Marco Baroni, Gemma Boleda, and Raquel Fernández. The LAMBADA dataset: Word prediction requiring a broad discourse context. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1525–1534, Berlin, Germany, August 2016. Association for Computational Linguistics. doi: 10.18653/v1/P16-1144. <https://aclanthology.org/P16-1144/>.

Francesco Pappone, Donato Crisostomi, and Emanuele Rodolà. Two-scale latent dynamics for recurrent-depth transformers, 2025. <https://arxiv.org/abs/2509.23314>.

Guilherme Penedo, Hynek Kydlíček, Loubna Ben allal, Anton Lozhkov, Margaret Mitchell, Colin Raffel, Leandro Von Werra, and Thomas Wolf. The fineweb datasets: Decanting the web for the finest text data at scale, 2024. <https://arxiv.org/abs/2406.17557>.

Hayden Prairie, Zachary Novack, Taylor Berg-Kirkpatrick, and Daniel Y. Fu. Parcae: Scaling laws for stable looped language models, 2026. <https://arxiv.org/abs/2604.12946>.

- Ofir Press and Lior Wolf. Using the output embedding to improve language models. In *Proceedings of the 15th Conference of the European Chapter of the Association for Computational Linguistics: Volume 2, Short Papers*, pages 157–163, 2017.
- Paulius Rauba, Claudio Fanconi, and Mihaela van der Schaar. Tiny autoregressive recursive models, 2026. <https://arxiv.org/abs/2603.08082>.
- Michael Rizvi-Martel, Guillaume Rabusseau, and Marius Mosbach. The illusion of superposition? A principled analysis of latent thinking in language models. *arXiv preprint arXiv:2604.06374*, 2026.
- Grigory Sapunov. Universal transformers need memory: Depth-state trade-offs in adaptive recursive reasoning, 2026a. <https://arxiv.org/abs/2604.21999>.
- Grigory Sapunov. Universal transformers need memory: Depth-state trade-offs in adaptive recursive reasoning, 2026b. <https://arxiv.org/abs/2604.21999>.
- Nikunj Saunshi, Nishanth Dikkala, Zhiyuan Li, Sanjiv Kumar, and Sashank J. Reddi. Reasoning with latent thoughts: On the power of looped transformers, 2025a. <https://arxiv.org/abs/2502.17416>.
- Nikunj Saunshi, Nishanth Dikkala, Zhiyuan Li, Sanjiv Kumar, and Sashank J. Reddi. Reasoning with latent thoughts: On the power of looped transformers, 2025b. <https://arxiv.org/abs/2502.17416>.
- Shixiang Song, He Li, Zitong Wang, Boyi Zeng, Feichen Song, Yixuan Wang, Zhiqin John Xu, Ziwei He, and Zhouhan Lin. AdaPonderLM: Gated pondering language models with token-wise adaptive depth. *arXiv preprint arXiv:2603.01914*, 2026a.
- Shixiang Song, He Li, Zitong Wang, Boyi Zeng, Feichen Song, Yixuan Wang, Zhiqin John Xu, Ziwei He, and Zhouhan Lin. Adaponderlm: Gated pondering language models with token-wise adaptive depth, 2026b. <https://arxiv.org/abs/2603.01914>.
- Sho Takase and Shun Kiyono. Lessons on parameter sharing across layers in transformers, 2023. <https://arxiv.org/abs/2104.06022>.
- Guo Tang, Shixin Jiang, Heng Chang, Nuo Chen, Yuhan Li, Huiming Fan, Jia Li, Ming Liu, and Bing Qin. Looprpt: Reinforcement pre-training for looped language models, 2026. <https://arxiv.org/abs/2603.19714>.
- Jayden Teoh, Manan Tomar, Kwangjun Ahn, Edward S. Hu, Pratyusha Sharma, Riashat Islam, Alex Lamb, and John Langford. Next-latent prediction transformers learn compact world models, 2025. <https://arxiv.org/abs/2511.05963>.
- Yalcin Tur, Jalal Naghiyev, Haoquan Fang, Wei-Chuan Tsai, Jiafei Duan, Dieter Fox, and Ranjay Krishna. Recurrent-depth vla: Implicit test-time compute scaling of vision-language-action models via latent iterative reasoning, 2026. <https://arxiv.org/abs/2602.07845>.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Advances in Neural Information Processing Systems 30*, pages 5998–6008, 2017.
- Guan Wang, Jin Li, Yuhao Sun, Xing Chen, Changling Liu, Yue Wu, Meng Lu, Sen Song, and Yasin Abbasi Yadkori. Hierarchical reasoning model, 2025. <https://arxiv.org/abs/2506.21734>.
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed Chi, Quoc Le, and Denny Zhou. Chain-of-thought prompting elicits reasoning in large language models, 2023. <https://arxiv.org/abs/2201.11903>.
- Xilin Wei, Xiaoran Liu, Yuhang Zang, Xiaoyi Dong, Yuhang Cao, Jiaqi Wang, Xipeng Qiu, and Dahua Lin. Sim-cot: Supervised implicit chain-of-thought, 2025. <https://arxiv.org/abs/2509.20317>.
- Jonathan Williams and Esin Tureci. Prioritize the process, not just the outcome: Rewarding latent thought trajectories improves reasoning in looped language models, 2026. <https://arxiv.org/abs/2602.10520>.
- Bohong Wu, Mengzhao Chen, Xiang Luo, Shen Yan, Qifan Yu, Fan Xia, Tianqi Zhang, Hongrui Zhan, Zheng Zhong, Xun Zhou, Siyuan Qiao, and Xingyan Bin. Parallel loop transformer for efficient test-time computation scaling, 2025. <https://arxiv.org/abs/2510.24824>.
- Kevin Xu and Issei Sato. On expressive power of looped transformers: Theoretical analysis and enhancement via timestep encoding, 2025a. <https://arxiv.org/abs/2410.01405>.
- Kevin Xu and Issei Sato. On expressive power of looped transformers: Theoretical analysis and enhancement via timestep encoding, 2025b. <https://arxiv.org/abs/2410.01405>.

- Kevin Xu and Issei Sato. A formal comparison between chain of thought and latent thought, 2026. <https://arxiv.org/abs/2509.25239>.
- Liu Yang, Kangwook Lee, Robert Nowak, and Dimitris Papailiopoulos. Looped transformers are better at learning learning algorithms, 2024a. <https://arxiv.org/abs/2311.12424>.
- Liu Yang, Kangwook Lee, Robert Nowak, and Dimitris Papailiopoulos. Looped transformers are better at learning learning algorithms, 2024b. <https://arxiv.org/abs/2311.12424>.
- Liu Yang, Kangwook Lee, Robert Nowak, and Dimitris Papailiopoulos. Looped transformers are better at learning learning algorithms, 2024c. <https://arxiv.org/abs/2311.12424>.
- Abbas Zeitoun, Lucas Torroba-Hennigen, and Yoon Kim. Hyperloop transformers. *arXiv preprint arXiv:2604.21254*, 2026a.
- Abbas Zeitoun, Lucas Torroba-Hennigen, and Yoon Kim. Hyperloop transformers, 2026b. <https://arxiv.org/abs/2604.21254>.
- Boyi Zeng, Shixiang Song, Siyuan Huang, Yixuan Wang, He Li, Ziwei He, Xinbing Wang, Zhiyu Li, and Zhouhan Lin. Ponderlm: Pretraining language models to ponder in continuous space, 2026. <https://arxiv.org/abs/2505.20674>.
- Xiaojing Zhang, Haifeng Wu, Gang He, Jiyang Shen, Bochen Lyu, and Zhanxing Zhu. MoDr: Mixture-of-depth-recurrent transformers for test-time reasoning. In *The Fourteenth International Conference on Learning Representations*.
- Xiaojing Zhang, Haifeng Wu, Gang He, Jiyang Shen, Bochen Lyu, and Zhanxing Zhu. Modr: Mixture-of-depth-recurrent transformers for test-time reasoning. In *The Fourteenth International Conference on Learning Representations*, 2026. <https://openreview.net/forum?id=9Pba4rcQbE>.
- Yunao Zheng, Xiaojie Wang, Lei Ren, and Chen Wei. ChainGPT: Dual-reasoning model with recurrent depth and multi-rank state updates. In *The Fourteenth International Conference on Learning Representations*.
- Yunao Zheng, Xiaojie Wang, Lei Ren, and Chen Wei. ChainGPT: Dual-reasoning model with recurrent depth and multi-rank state updates. In *The Fourteenth International Conference on Learning Representations*, 2026. <https://openreview.net/forum?id=kdZbxizwGK>.
- Cai Zhou, Chenxiao Yang, Yi Hu, Chenyu Wang, Chubin Zhang, Muhan Zhang, Lester Mackey, Tommi Jaakkola, Stephen Bates, and Dinghuai Zhang. Coevolutionary continuous discrete diffusion: Make your diffusion language model a latent reasoner. *arXiv preprint arXiv:2510.03206*, 2025.
- Hanlin Zhu, Shibo Hao, Zhiting Hu, Jiantao Jiao, Stuart Russell, and Yuandong Tian. Emergence of superposition: Unveiling the training dynamics of chain of continuous thought. *arXiv preprint arXiv:2509.23365*, 2025a.
- Hanlin Zhu, Shibo Hao, Zhiting Hu, Jiantao Jiao, Stuart Russell, and Yuandong Tian. Reasoning by superposition: A theoretical perspective on chain of continuous thought, 2025b. <https://arxiv.org/abs/2505.12514>.
- Hanlin Zhu, Shibo Hao, Zhiting Hu, Jiantao Jiao, Stuart J Russell, and Yuandong Tian. Reasoning by superposition: A theoretical perspective on chain of continuous thought. *Advances in Neural Information Processing Systems*, 38: 79931–79963, 2026.
- Rui-Jie Zhu, Zixuan Wang, Kai Hua, Tianyu Zhang, Ziniu Li, Haoran Que, Boyi Wei, Zixin Wen, Fan Yin, He Xing, Lu Li, Jiajun Shi, Kaijing Ma, Shanda Li, Taylor Kergan, Andrew Smith, Xingwei Qu, Mude Hui, Bohong Wu, Qiyang Min, Hongzhi Huang, Xun Zhou, Wei Ye, Jiaheng Liu, Jian Yang, Yunfeng Shi, Chenghua Lin, Enduo Zhao, Tianle Cai, Ge Zhang, Wenhao Huang, Yoshua Bengio, and Jason Eshraghian. Scaling latent reasoning via looped language models, 2025c. <https://arxiv.org/abs/2510.25741>.
- Łukasz Kaiser and Ilya Sutskever. Neural gpu learn algorithms, 2016. <https://arxiv.org/abs/1511.08228>.

## A Related Work

**Looped and recurrent language models.** Weight-tied recurrence has re-emerged as an alternative to deeper feed-forward stacks (Press and Wolf, 2017; Inan et al., 2017; Łukasz Kaiser and Sutskever, 2016; Lan et al., 2020; Dabre and Fujita, 2018; Takase and Kiyono, 2023; Bae et al., 2025a). Universal Transformers share a single block across depth (Dehghani et al., 2019; Sapunov, 2026a,b), while looped and recurrent language models iterate a recurrent update on a hidden state to enable latent reasoning (Zhu et al., 2025c; Prairie et al., 2026; Geiping et al., 2025a; Bae et al., 2025b; Zheng et al.; Zhang et al.; Song et al., 2026a; Knupp et al., 2026a; Zeitoun et al., 2026a; McLeish et al.; Zheng et al., 2026; Zhang et al., 2026; Song et al., 2026b; Moosa et al., 2026; Knupp et al., 2026b; Wu et al., 2025; Zeitoun et al., 2026b; Pappone et al., 2025; Rauba et al., 2026; Zeng et al., 2026; Tur et al., 2026; Komisarczyk et al., 2026; Ghugare et al., 2026; Cameron et al., 2026; Williams and Tureci, 2026; Han, 2026; Goyal et al., 2026; Tang et al., 2026). Latent (implicit) reasoning methods such as Coconut (Hao et al., 2025) and related methods (Teoh et al., 2025; Chen et al., 2025; Xu and Sato, 2026; Jeddi et al., 2026; Fu et al., 2026; Zhou et al., 2025; Geiping et al., 2025b; Hu et al., 2025; Fleuret, 2025) perform reasoning in continuous representation space rather than through chain-of-thought tokens. These approaches unroll the loop for a fixed number of steps at training time, causing training memory to grow linearly in depth, coupling inference iterations to training, and often degrading quality when iterations are extended at test time (Zhu et al., 2025c). From theoretical and mechanistic analysis perspectives, looped and continuous thinking models have shown to offer benefits over non-recurrent models (Yang et al., 2024a; Labovich, 2026b; Saunshi et al., 2025a; Xu and Sato, 2025a; Saunshi et al., 2025b; Giannou et al., 2023b; Yang et al., 2024c; Gatmiry et al., 2024b; Huang et al., 2025; Xu and Sato, 2025b; Merrill and Sabharwal, 2025; Labovich, 2026a; Zhu et al., 2025b,a; Gatmiry et al., 2024a; Merrill and Sabharwal, 2026). Mechanistic analysis shows that the recurrent updates of looped LMs in fact converge to a fixed point for the vast majority of tokens, directly motivating our formulation (Blayney et al., 2026).

**Implicit fixed-point models.** Deep Equilibrium Models (DEQs) replace finite unrolling with a fixed-point equation  $z^* = f_\theta(z^*, x)$  and train through it via implicit differentiation, decoupling effective depth from training memory (Bai et al., 2019). Solvers typically use Anderson acceleration (Anderson, 1965), with backward passes ranging from full implicit differentiation to cheap surrogates such as one-step (Fung et al., 2021) and phantom (Geng et al., 2022) gradients. Standard DEQs solve for a hidden state initialized from zero and decode through a separate output head. In Attractor Models, the equilibrium instead lives in the tied output-embedding space and is warm-started from a meaningful backbone prediction; we find this structure essential for stable training at language-modeling scale and show it gives rise to equilibrium internalization (Section 4.3).

**Tiny recursive reasoners.** On small-data algorithmic benchmarks such as Sudoku and maze solving, hierarchical and tiny recursive networks (HRM, TRM) (Wang et al., 2025; Jolicoeur-Martineau, 2025) achieve strong accuracy with few parameters, but exhibit a “less is more” behavior where performance collapses as model size grows (Jolicoeur-Martineau, 2025). Attractor Models retain the iterative-refinement benefit of these architectures while scaling cleanly with parameter count).

**Positioning.** Looped LMs couple three quantities through unrolled BPTT: inference depth, training depth, and training memory (Zhu et al., 2025c). Attractor Models decouple all three: the equilibrium is defined by a residual equation in output-embedding space, the number of solver evaluations is chosen adaptively by tolerance  $\varepsilon$ , and training memory in the recurrent block is constant in effective depth. Empirically, this combines the strengths of feed-forward and recurrent models—better perplexity than parameter-matched Transformers, 25–31% lower training FLOPs and constant memory relative to looped LMs, and clean scaling on hard reasoning where specialized recursive networks collapse.

## B Theory of Looped and Attractor Models

### B.1 Well-posedness and the implicit gradient

Fix an input  $x$  and parameters  $\theta$ . Throughout this section, we write  $F(y) \triangleq f_\theta(y, E(x))$  for brevity, so that the fixed-point equation  $y^* = f_\theta(y^*, E(x))$  becomes  $y^* = F(y^*)$ . Let  $J_F(y) \triangleq \partial F / \partial y$  denote the Jacobian of

$F$  with respect to its state input. We work in a generic norm  $\|\cdot\|$  on  $\mathbb{R}^{L \times d}$  and its induced operator norm; for concreteness, can take both to be Frobenius and spectral, respectively. Assumptions are borrowed from (Bai et al., 2020; Krantz and Parks, 2002a).

**Assumption 1** (Local contraction). *There exist a point  $\bar{y} \in \mathbb{R}^{L \times d}$ , a radius  $r > 0$ , and a constant  $L \in [0, 1)$  such that*

1.  $F$  maps the closed ball  $B_r(\bar{y}) \triangleq \{y : \|y - \bar{y}\| \leq r\}$  into itself, and
2.  $F$  is  $L$ -Lipschitz on  $B_r(\bar{y})$ :

$$\|F(y) - F(y')\| \leq L \|y - y'\| \quad \text{for all } y, y' \in B_r(\bar{y}).$$

When  $F$  is continuously differentiable, a sufficient form of (ii) is  $\sup_{y \in B_r(\bar{y})} \|J_F(y)\| \leq L$ .

**Theorem 1** (Well-posedness and implicit gradient). *Under Assumption 1, the following hold.*

1. Existence and uniqueness: *There is a unique  $y^* \in B_r(\bar{y})$  with  $y^* = F(y^*)$ .*
2. The picard iteration converges linearly: *For any  $y_0 \in B_r(\bar{y})$ , the iterates  $y_{k+1} = F(y_k)$  remain in  $B_r(\bar{y})$  and satisfy*

$$\|y_k - y^*\| \leq L^k \|y_0 - y^*\|. \quad (5)$$

3. Validity of the implicit gradient: *If  $F$  is continuously differentiable on  $B_r(\bar{y})$ , then  $I - J_F(y^*)$  is invertible,  $y^*$  depends continuously differentiably on  $\theta$  in a neighborhood of the current parameters, and*

$$\frac{\partial y^*}{\partial \theta} = (I - J_F(y^*))^{-1} \frac{\partial f_\theta}{\partial \theta} \Big|_{y^*}. \quad (6)$$

*Proof.* Parts (i) and (ii) follow from the Banach fixed-point theorem applied to the contraction  $F$  on the closed (hence complete) ball  $B_r(\bar{y}) \subset (\mathbb{R}^{L \times d}, \|\cdot\|)$ .

For (iii), Lipschitz continuity with constant  $L$  gives  $\|J_F(y^*)\| \leq L < 1$ , so the spectral radius satisfies  $\rho(J_F(y^*)) \leq \|J_F(y^*)\| < 1$ . The Neumann series  $\sum_{k \geq 0} J_F(y^*)^k$  therefore converges to  $(I - J_F(y^*))^{-1}$ , and in particular  $I - J_F(y^*)$  is invertible. Define  $G(y, \theta) \triangleq f_\theta(y, E(x)) - y$ , which is  $C^1$  in both arguments and satisfies  $G(y^*, \theta) = 0$  with

$$\frac{\partial G}{\partial y} \Big|_{y^*} = J_F(y^*) - I, \quad (7)$$

which is invertible. The implicit function theorem applied to  $G = 0$  yields a unique  $C^1$  map  $\theta \mapsto y^*(\theta)$  in a neighborhood of the current parameters, with

$$\frac{\partial y^*}{\partial \theta} = - \left( \frac{\partial G}{\partial y} \right)^{-1} \frac{\partial G}{\partial \theta} = (I - J_F(y^*))^{-1} \frac{\partial f_\theta}{\partial \theta} \Big|_{y^*}. \quad (8)$$

□

Equation (6) is exactly the gradient computed in the backward pass of Section 3. The linear solve is convergent for the same reason: since  $\|J_F^\top(y^*)\| = \|J_F(y^*)\| \leq L < 1$ , the Neumann iteration  $u \leftarrow J_F^\top(y^*) u + v$  converges geometrically with rate  $L$  to the unique solution  $u = (I - J_F^\top(y^*))^{-1} v$ .

## B.2 Looped language models are fixed-point iterators

A LoopLM (Zhu et al., 2025c) of depth  $T$  with the same weight-tied block  $f_\theta$  and warm start  $y_0 = E(x)$  produces

$$y_T^{\text{loop}} \triangleq F^{(T)}(E(x)) = \underbrace{F \circ F \circ \dots \circ F}_{T \text{ times}}(E(x)). \quad (9)$$

This is exactly  $T$  Picard iterations of our residual  $g_\theta(\cdot, x)$  warm-started from the input embedding. See the following corollary:

**Corollary 1** (LoopLM as a truncated approximation). *Suppose Assumption 1 holds and the warm start  $E(x) \in B_r(\bar{y})$ . Then for every  $T \geq 0$ ,*

$$\|y_T^{\text{loop}} - y^*\| \leq L^T \|E(x) - y^*\|. \quad (10)$$

*In particular,  $y_T^{\text{loop}} \rightarrow y^*$  as  $T \rightarrow \infty$ , and the discrepancy between a LoopLM of depth  $T$  and our model decays geometrically in  $T$ .*

*Proof.* Apply Theorem 1(ii) with  $y_0 = E(x)$ . □

Two consequences are worth highlighting:

**The limit:** Our fixed-point model is the  $T \rightarrow \infty$  limit of LoopLM with shared parameters. Training a LoopLM at any finite depth  $T$  implicitly trains an approximation to the same equilibrium  $y^*$ , with the approximation error controlled by (10).

**The iteration count at inference:** Picard iteration converges at the linear rate  $L$ . Our forward pass uses Anderson acceleration, which under standard regularity conditions converges *superlinearly* near  $y^*$  (Anderson, 1965). To reach a target residual tolerance  $\varepsilon$ , the root finder therefore typically requires fewer evaluations of  $f_\theta$  than the unroll depth  $T$  a LoopLM would need for a comparable residual. The exact crossover depends on  $L$  and on the solver hyperparameters.

**Caveat:** Assumption 1 is local and is not guaranteed to hold for an arbitrarily trained transformer block. In our experiments, we generally notice that  $\|J_F(y^*)\| \leq L < 1$ , the authors of DEQ also find that using a regularizer on the Jacobian to enforce this helps (Bai et al., 2021).

## B.3 Where looped language models fall short

### B.3.1 Attractor Models

The initialization  $\tilde{y}_0$  and the fixed point  $\tilde{y}^*$  live in the same tied output-embedding space, so either can be decoded by the unembedding  $E^\top$ . The loss depends only on  $\tilde{y}^*$ , and the implicit gradient with respect to the backbone parameters is

$$\nabla_{\theta_b} \mathcal{L} = u^\top J_{\tilde{y}_0} \nabla_{\theta_b} \tilde{y}_0, \quad u = (I - J_{\tilde{y}}^\top)^{-1} v.$$

Thus, the loss may be viewed as a function of the backbone’s output embedding: perturbations in  $\tilde{y}_0$  affect  $\mathcal{L}$  by perturbing  $\tilde{y}^*$  in the same space. The backbone is therefore optimized as a standalone next-token predictor whose embedding is decoded by  $E^\top$ .

When the backbone has strictly greater capacity than the weight-tied attractor, the joint optimum places the prediction work in the backbone. Consequently,  $\tilde{y}_0(x; \theta_b^*)$  approaches the loss-minimizing embedding. The fixed-point constraint

$$\mathcal{T}_{\theta_a^*}(\tilde{y}^*, \tilde{y}_0) = \tilde{y}^*$$

is then consistent only when  $\tilde{y}^* = \tilde{y}_0$ , yielding

$$\mathcal{T}_{\theta_a^*}(\tilde{y}_0, \tilde{y}_0) = \tilde{y}_0.$$

### B.3.2 Standard looped LMs

In standard looped language models, the latent initialization  $h_0$  is uninformative, typically zero or noise, and occupies the role of a hidden state rather than an output embedding. Decoding  $h_0 E^\top$  is therefore not meaningful.

Gradients instead flow by backpropagation through time along the unrolled trajectory  $(h_0, \dots, h_T)$ . There is no term in this gradient that drives  $h_0$  toward the loss-bearing final state  $h_T$  in embedding distance: the initialization is fixed by the architecture and is not produced by a separately trained predictor.  $\square$

### B.3.3 Implicit-Gradient Barrier

The implicit gradient

$$\nabla_{\theta} \mathcal{L}_{\infty} = v^\top (I - J_g^\top)^{-1} \partial_{\theta} g$$

contains the factor  $(I - J_g)^{-1}$ . This inverse is undefined when  $1 \in \sigma(J_g)$ , and its operator norm scales as

$$\|(I - J_g)^{-1}\| \sim \text{dist}(1, \sigma(J_g))^{-1}$$

near such a singularity.

Along any continuous gradient-descent path, the eigenvalues of  $J_g(y^*; \theta)$  vary continuously. Therefore, for the dominant eigenvalue to leave the unit disk through  $+1$ —the canonical loss-of-contraction route for residual iteration maps—it must approach 1. This forces  $\|(I - J_g)^{-1}\| \rightarrow \infty$ .

Unless  $v$  is orthogonal to the offending eigendirection, a codimension-one and hence non-generic condition, the gradient norm diverges. The descent step therefore blows up before the boundary can be crossed, confining  $\theta$  to the contractive region  $\{\rho(J_g) < 1\}$ .

For exits through complex eigenvalues, the same conclusion holds under the one-step JFB approximation, since the truncated Neumann series

$$\sum_k (J_g^\top)^k$$

diverges as  $\rho(J_g) \rightarrow 1$ .  $\square$

### B.3.4 Interpretation

Fixed-loop training has no inherent mechanism that favors contractive iterations. An optimum in which  $g_{\theta}$  perturbs the embedding for  $K$  steps and only happens to land accurately at step  $K$  is a valid fixed-loop solution. These are precisely the solutions that fail when extra loops are run at inference.

Equilibrium training cannot reach such solutions, because the implicit gradient creates a barrier of diverging gradients around non-contractive regions. As a result, the trajectory of  $\theta$  is confined to the regime in which the solver converges, the one-step gradient is a descent direction, and additional iterations remain stable.

The mechanistic observation (Blayney et al., 2026) that fixed points emerge only late in standard looped training is consistent with this picture: the loss landscape contains a basin near contractive solutions, but only equilibrium training applies pressure toward that basin from the beginning of training.

## C Hyperparameter and Experimental Settings

We use NVIDIA H200 GPUs for all of our experiments.

**Table 7** Architecture hyperparameters for each model family and scale.

Hyperparameter	Transformer			Parcae			Attractor Model		
	140M	370M	770M	140M	370M	770M	140M	370M	770M
$d_{\text{model}}$	768	1024	1280	768	1024	1280	1024	1280	1280
$d_{\text{ff}}$	3072	4096	5120	3072	4096	5120	4096	5120	5120
Attention heads	6	8	10	6	8	10	8	10	10
Total layers	6	12	18	6	12	18	8	17	37
Prelude	—	—	—	2	4	6	7	15	35
Core (recurrent)	—	—	—	2	4	6	1	2	2
Coda	—	—	—	2	4	6	0	0	0
Mean recurrence $T$	1	1	1	8	8	8	(solver-determined)		
Sequence length	2048								
Vocab size	32768								
Tie embeddings							✓		
Norm							RMSNorm ( $\epsilon = 10^{-5}$ )		
Activation							ReLU <sup>2</sup>		
QK-norm							✓		
Precision							bf16-mixed		

**Table 8** Fixed-point solver hyperparameters for Attractor Model. Parcae and Transformer do not use a solver.

Hyperparameter	140M	370M	770M
Forward solver	Anderson	Anderson	Anderson
Max iterations (fwd)	64	64	32
Min iterations (fwd)	6	8	6
Tolerance (fwd)	$3 \times 10^{-4}$	$2 \times 10^{-4}$	$1.5 \times 10^{-4}$
Anderson $m$	5	5	5
Anderson $\beta$	1.0	1.0	0.96
Backward type	one-step IFT	one-step IFT	Anderson
Max iterations (bwd)	64	64	32
Min iterations (bwd)	6	6	6
Tolerance (bwd)	$3 \times 10^{-4}$	$2 \times 10^{-4}$	$1.5 \times 10^{-4}$
Adjoint grad clip	—	—	2.0
Layer-scale init ( $\gamma_0$ )	0.75	0.75	0.73
$\gamma_{\text{max}}$	0.75	1.0	0.9
FP LR scale	0.5	0.4	0.4
FP weight decay	0.1	0.1	0.1

**Table 9** Training hyperparameters. All three model families share the same optimizer, schedule, and data configuration at each scale. Differences are noted where they occur.

Hyperparameter	140M	370M	770M
Training tokens	11.2B	29.6B	61.6B
Dataset	FineWeb-Edu 100B (shuffled)		
Tokenizer	BPE, 32768 vocab		
Optimizer	MuonAdamW		
AdamW LR	$8 \times 10^{-3}$	$8 \times 10^{-3}$	$5-6 \times 10^{-3\dagger}$
AdamW $(\beta_1, \beta_2)$	(0.8, 0.95)		
AdamW $\epsilon$	$10^{-10}$		
AdamW weight decay	0		
Muon LR	$8 \times 10^{-3}$	$8 \times 10^{-3}$	$6-8 \times 10^{-3\dagger}$
Muon momentum	0.95		
Muon weight decay	0.2 (linear decay to 0)		
Muon NS steps	5		
LR schedule	Trapezoid (0% warmup, 50% cooldown)		
Min LR	0		
Gradient clipping	1.0		
Global batch size (tokens)	524K (= $256 \times 2048$ )		
Micro batch size	32	32	8-32 <sup>‡</sup>
Gradient checkpointing	✗	✗	model-dependent <sup>‡</sup>
<code>torch.compile</code>	✓*	✓*	✓*
Strategy	DDP		
Init strategy	scaled-zero + orthogonal		
Seed	42		

<sup>†</sup>At 770M, Attractor Model uses AdamW LR  $5 \times 10^{-3}$  and Muon LR  $6 \times 10^{-3}$ ; Transformer/Parcae use  $6 \times 10^{-3}$  /  $8 \times 10^{-3}$ .

<sup>‡</sup>Gradient checkpointing enabled for 770M Attractor Model and Parcae; disabled for Transformer.

\*`torch.compile` enabled for Transformer and Parcae; disabled for Attractor Model (implicit-gradient hooks are not compile-compatible).

**Table 10** Parcae recurrence hyperparameters.

Hyperparameter	140M	370M	770M
Injection type	diagonal		
State init	like-init		
Mean recurrence	8	8	8
Mean backprop depth	4	4	4
Sampling scheme	poisson-truncated-full		
Iteration method	per-sequence		
Prelude norm	✓		