
VecTur: Vector Turing Machines

Ethan Hall

ethan.hall.phd@gmail.com

Abstract

We introduce VECTUR (Vector Turing Machines), a differentiable architecture for learning functions by vectorizing the tape symbols, head index, and finite control of a classical Turing machine (Turing, 1936). VECTUR is designed to support *deep computation* by explicitly iterating a learned transition map while learning an input-conditioned halting schedule via a parameter κ . We evaluate VECTUR as a drop-in *computational block* inside a Llama-style decoder-only language model (Touvron et al., 2023; Dubey et al., 2024), contrasting it with attention (Vaswani et al., 2017), LSTM-style recurrence (Hochreiter and Schmidhuber, 1997), and differentiable external-memory baselines (Graves et al., 2014, 2016). On small-to-medium open benchmarks for reasoning and language (GSM8K (Cobbe et al., 2021), ARC (Clark et al., 2018), HellaSwag (Zellers et al., 2019), WikiText-103 (Merity et al., 2016)), we find VECTUR improves algorithmic generalization at fixed parameter budgets. We additionally introduce COMPGEN—a synthetic program dataset stratified by $(T(n), S(n))$ complexity classes—and show VECTUR adapts its effective compute by learning κ . Finally, we propose VECSTUR, a stochastic extension that consumes random tape symbols and outperforms VECTUR on randomized verification tasks (e.g., Freivalds-style matrix product verification (Freivalds, 1977)).

1 Introduction

Modern language models excel at pattern completion yet often struggle to reliably *execute* long algorithmic computations, extrapolate beyond training lengths, or allocate variable compute per input. Several lines of work attempt to address these limitations by embedding algorithmic structure into neural systems, including external-memory architectures (Graves et al., 2014, 2016) and adaptive computation mechanisms (Graves, 2016).

We propose VECTUR, a vectorized analogue of a classical Turing machine (Turing, 1936) whose tape, head index, and finite control are represented as continuous vectors and updated by a learned transition map. Unlike attention-only computation (Vaswani et al., 2017), VECTUR explicitly iterates a state transition; unlike Neural Turing Machines (Graves et al., 2014), VECTUR emphasizes sparse head indexing and a learned halting schedule parameterized by κ , enabling adaptive depth.

We evaluate VECTUR in a realistic regime by inserting it as a *block* inside a Llama-style decoder-only macro architecture (Touvron et al., 2023; Dubey et al., 2024), replacing the standard attention+MLP block. We compare against alternative blocks: (i) standard attention (Vaswani et al., 2017), (ii) LSTM-style recurrence (Hochreiter and Schmidhuber, 1997), and (iii) differentiable external-memory controllers (Graves et al., 2014, 2016). We focus on small-to-medium models (roughly 10^8 to 10^9 parameters) where architectural inductive bias can materially affect sample efficiency and extrapolation.

We further introduce COMPGEN, a dataset of generated Python programs grouped into discrete complexity buckets $(T(n), S(n))$ such as $O(n)/O(1)$, $O(n \log n)/O(1)$, $O(n^2)/O(1)$, and $O(n^2)/O(n)$. The goal is to probe whether a model can learn to *compute* across increasing n by allocating more steps as needed, rather than memorizing only small n . Finally, we define VECSTUR, which aug-

40 ments the input with stochastic symbols z to emulate randomized computation, and we propose a
41 randomized evaluation suite where randomness yields asymptotic speedups (Freivalds, 1977; Miller,
42 1976; Rabin, 1980). In particular, we form a data set of matrices $A, B, C \in \mathbb{R}^{n \times n}$ and a target
43 matrix $D \in \mathbb{R}^{n \times n}$ such that $D = AB$ and $D = AC$ with probability 1/2. We evaluate VECTUR and
44 VECSTUR on this task, and show VECSTUR can exploit stochastic symbols to achieve asymptotic
45 speedups.

46 **Contributions.**

- 47 • We define VECTUR, a vectorized Turing-machine-inspired transition system with sparse indexing
48 and a learnable halting schedule κ .
49 • We propose a plug-and-play integration of VECTUR as a computational block inside Llama-style
50 decoder-only models.
51 • We introduce COMPGEN, a synthetic program dataset labeled by time/space complexity class
52 $(T(n), S(n))$, and a protocol for extrapolation across n .
53 • We define VECSTUR and a randomized computation evaluation suite, showing benefits of stochastic
54 symbols for problems with randomized speedups.

55 **2 Related Work**

56 **Sequence models and attention.** Transformers (Vaswani et al., 2017) and their decoder-only
57 variants power modern LLMs (e.g., Llama-family models (Touvron et al., 2023; Dubey et al.,
58 2024)). Recurrent networks such as LSTMs (Hochreiter and Schmidhuber, 1997) provide a different
59 inductive bias for iterative computation but historically underperform attention-based models at scale
60 on language modeling.

61 **Differentiable memory and neural machines.** Neural Turing Machines (NTMs) (Graves et al.,
62 2014) and Differentiable Neural Computers (DNCs) (Graves et al., 2016) integrate external memory
63 with differentiable read/write heads. Our work shares the goal of improving algorithmic behavior, but
64 emphasizes (i) sparse indexing for efficiency and (ii) a learned halting schedule.

65 **2.1 Remark: Neural Turing Machines vs. VECTUR**

66 Both NTMs (Graves et al., 2014) and VECTUR augment neural computation with an external memory,
67 but they make different design trade-offs for *addressing* (how memory is accessed) and *sparsity*
68 (how much memory is touched per step). NTMs provide content-addressable reads/writes via dense
69 attention over all memory slots; VECTUR instead maintains a continuous head position on a tape
70 and performs sparse gather/scatter updates to a small number of adjacent cells (via interpolation),
71 which keeps per-step cost independent of tape length. In Llama-style Transformer blocks, global
72 content-based access is already available through self-attention; VECTUR is intended to add an
73 *orthogonal* capability: cheap, iterative state manipulation with persistent scratchpad dynamics.

74 **Adaptive computation.** Adaptive Computation Time (ACT) (Graves, 2016) learns when to stop
75 iterating; our κ -parameterization provides a simple, input-conditioned control knob for the effective
76 number of steps.

77 **Randomized algorithms.** Randomness can reduce expected runtime for verification and decision
78 problems; canonical examples include Freivalds' randomized matrix product verification (Freivalds,
79 1977) and probabilistic primality testing (Miller, 1976; Rabin, 1980). VECSTUR is intended as a
80 neural analogue that can exploit stochastic symbols during computation.

81 **3 VecTur: Vector Turing Machines**

82 **3.1 Vectorized machine state**

83 Given an input sequence $x \in \mathbb{R}^{N \times d_x}$ (e.g., token embeddings), we define a VECTUR block below.
84 Note that for VECSTUR, we additionally sample a sequence of stochastic symbols $z \in \mathbb{R}^{N_z \times d_x}$ and

Feature	NTM (Graves)	VECTUR	VECTUR advantage
Addressing	Dense content-based attention over all slots	Local/location-based head movement on tape	Yes
Per-step complexity (vs. tape length N_T)	$O(N_T)$ similarity + weighted sum	$O(k)$ gather/scatter (independent of N_T)	Yes
Forward activation memory (unrolled T steps)	Stores dense weights $\sim O(TN_T)$	Stores sparse indices/weights $\sim O(Tk)$	Yes
State preservation away from head	Many slots updated slightly (drift/blurring risk)	Un-accessed cells are exactly unchanged	Yes
Content lookup in 1 step	Native (query by key)	Requires scanning via head movement ($\text{worst-case } O(N_T)$ steps)	No
Inductive bias	Random-access / associative recall	Sequential pointer machine / local algorithms	Depends
Compute allocation / halting	Typically fixed unroll or implicit stopping	Explicit learned halting via κ and gate g_t	Yes
Fit as a Transformer block	Redundant global attention inside block	Complements attention with iterative scratchpad dynamics	Yes

Figure 1: **NTM vs. VECTUR.** NTMs (Graves et al., 2014) provide dense, content-addressable memory access, while VECTUR enforces sparse, local tape access with adaptive depth. The rightmost column highlights regimes where VECTUR is especially advantageous as a computational block inside attention-based macro-architectures.

85 set the tape length

$$N_T = N + N_z, \quad (1)$$

86 so that each input symbol and each stochastic symbol can be addressed at least once. In our
87 experiments we use $N_z \approx N$ (so $N_T \approx 2N$).

88 We define the machine state at step t as a triple (T_t, Q_t, I_t) , where the tape $T_t \in \mathbb{R}^{N_T \times d_T}$, the control
89 state $Q_t \in \mathbb{R}^{d_Q}$, and the head index

$$I_t = (\theta_t, \mathbf{w}_t) \in (S^1)^k \times \mathbb{R}^k$$

90 are learned, differentiable quantities. (Here $\theta_t = (\theta_{t,1}, \dots, \theta_{t,k})$ parameterizes k head locations
91 on the circle and $\mathbf{w}_t = (w_{t,1}, \dots, w_{t,k})$ are the associated weights.) We index tape positions by
92 $j \in \{0, 1, \dots, N_T - 1\}$, and write $T_t[j] \in \mathbb{R}^{d_T}$ for the j -th tape symbol.

93 The initial state is produced by learned maps $\mathcal{M}_T, \mathcal{M}_Q, \mathcal{M}_I$ with parameters W :

$$T_0 = \mathcal{M}_T(x; W), \quad Q_0 = \mathcal{M}_Q(x; W), \quad I_0 = \mathcal{M}_I(x; W), \quad (2)$$

94 where $\mathcal{M}_T, \mathcal{M}_Q, \mathcal{M}_I$ can be any mapping using some parameters W .

95 3.2 Sparse addressing (keeping \mathbf{I} and \mathbf{J} fixed)

96 Define the following piecewise linear map $E : S^1 \rightarrow \mathbb{R}^{N_T}$ as

$$\begin{aligned} n(\theta) &= \left\lfloor \frac{N_T \theta}{2\pi} \right\rfloor \\ s(\theta) &= \frac{N_T \theta}{2\pi} - \left\lfloor \frac{N_T \theta}{2\pi} \right\rfloor \\ n^+(\theta) &= (n(\theta) + 1) \bmod N_T \\ E(\theta) &= (1 - s(\theta))e_{n(\theta)} + s(\theta)e_{n^+(\theta)} \end{aligned}$$

97 We will write any $I \in (S^1)^k \times \mathbb{R}^k$ as $I = (\theta, \mathbf{w})$, and define the induced sparse tape-index weighting
98 vector $J(I) \in \mathbb{R}^{N_T}$ by

$$J(I) = \sum_{i=1}^k w_i E(\theta_i). \quad (3)$$

99 By construction, each $E(\theta_i)$ is supported on at most two adjacent tape locations $\{n(\theta_i), n^+(\theta_i)\}$,
100 hence $J(I)$ is supported on at most $2k$ tape locations. Concretely, for each head atom (θ_i, w_i) define

$$n_i := n(\theta_i), \quad s_i := s(\theta_i), \quad n_i^+ := (n_i + 1) \bmod N_T,$$

101 so that $E(\theta_i) = (1 - s_i)e_{n_i} + s_i e_{n_i^+}$. This gives an implementation-friendly form: one can store
102 $(n_i, n_i^+, (1 - s_i)w_i, s_i w_i)$ for each i and never materialize the dense N_T -vector $J(I)$.

103 3.3 Read, transition, and halting

104 We define the transition map Δ that updates the tape, control state, and head index. First, we use the
105 head index $J(I_t)$ to read a single tape symbol $S_t \in \mathbb{R}^{d_T}$:

$$S_t = \sum_{j=0}^{N_T-1} (J(I_t))_j T_t[j] \in \mathbb{R}^{d_T}. \quad (4)$$

106 Equivalently, using the explicit $2k$ -sparse form above,

$$S_t = \sum_{i=1}^k w_{t,i} \left((1 - s_{t,i}) T_t[n_{t,i}] + s_{t,i} T_t[n_{t,i}^+] \right),$$

107 so S_t is computed using at most $2k$ gathered tape vectors, and is piecewise linear in the tape (and
108 linear in the interpolation weights away from the measure-zero segment boundaries induced by the
109 floor operation).

110 Next, define a gate $g_t \in (0, 1)$ that controls the effective amount of computation and enables early
111 stopping. We use a sigmoid gate,

$$g_t = \sigma \left(\frac{\kappa(x; W) - t}{\max(1, \|Q_t - q_0\|^2)} \right), \quad (5)$$

112 where $\sigma(u) = 1/(1 + e^{-u})$, $\kappa(x; W) > 0$ is a learned scalar per example, and $q_0 \in \mathbb{R}^{d_Q}$ is a learned
113 halting target state. Intuitively, larger $\kappa(x; W)$ yields more effective steps (slower decay in t), while
114 $\|Q_t - q_0\|$ encourages the dynamics to become stationary near the target.

115 We update the tape, control state, and head index using learned transition maps $\Delta_T, \Delta_Q, \Delta_\theta, \Delta_w$.
116 Let

$$U_t := \Delta_T(S_t, Q_t; W) \in \mathbb{R}^{d_T}.$$

117 Then the update equations are

$$T_{t+1}[j] = T_t[j] + g_t (J(I_t))_j U_t \quad \text{for } j \in \{0, \dots, N_T - 1\}, \quad (6)$$

$$Q_{t+1} = Q_t + g_t \Delta_Q(S_t, Q_t; W), \quad (7)$$

$$\theta_{t+1} = (\theta_t + g_t \Delta_\theta(S_t, Q_t, \theta_t; W)) \bmod 2\pi, \quad (8)$$

$$\mathbf{w}_{t+1} = \mathbf{w}_t + g_t \Delta_w(S_t, Q_t, \mathbf{w}_t; W), \quad (9)$$

$$I_{t+1} = (\theta_{t+1}, \mathbf{w}_{t+1}). \quad (10)$$

118 Equation (6) makes the sparsity explicit: since $(J(I_t))_j = 0$ for all but at most $2k$ locations,
119 only $O(2k)$ tape vectors are updated per step. In an efficient implementation, (6) is executed as a
120 scatter-add into those $2k$ indices (and S_t is computed as a gather + weighted sum).

121 The transition maps have the following types:

$$\Delta_T : \mathbb{R}^{d_T} \times \mathbb{R}^{d_Q} \rightarrow \mathbb{R}^{d_T},$$

$$\Delta_Q : \mathbb{R}^{d_T} \times \mathbb{R}^{d_Q} \rightarrow \mathbb{R}^{d_Q},$$

$$\Delta_\theta : \mathbb{R}^{d_T} \times \mathbb{R}^{d_Q} \times (S^1)^k \rightarrow \mathbb{R}^k,$$

$$\Delta_w : \mathbb{R}^{d_T} \times \mathbb{R}^{d_Q} \times \mathbb{R}^k \rightarrow \mathbb{R}^k.$$

122 The mod 2π in (10) ensures the head angles represent elements of S^1 (equivalently, Δ_θ may be
 123 chosen 2π -periodic in each component). With sparse gather/scatter, one step costs $O(k(d_T + d_Q))$
 124 time and $O(k(d_T + d_Q))$ working memory, plus the cost of evaluating the small transition networks.

125 **Early stopping and block output.** Fix a maximum unroll $T_{\max} \in \mathbb{N}$ and a threshold $\varepsilon > 0$. We
 126 run the transition until either $t = T_{\max}$ or the gate becomes negligible,

$$T(x) = \min\{t \in \{0, \dots, T_{\max} - 1\} : g_t < \varepsilon\},$$

127 with the convention $T(x) = T_{\max}$ if the set is empty. Concretely, we check $g_t < \varepsilon$ at the beginning
 128 of step t ; if it holds, we stop and return T_t . Otherwise, we apply the transition to produce T_{t+1} and
 129 continue. We define the VECTUR block output as the final tape

$$V(x) := T_{T(x)} \in \mathbb{R}^{N_T \times d_T}.$$

130 In downstream architectures (e.g., Llama-style models), any required reshaping or projection of $V(x)$
 131 is handled outside the VECTUR block.

132 **Differentiability and efficient backpropagation.** All operations inside each step are differentiable
 133 with respect to the tape values and the transition parameters, except at the measure-zero boundaries
 134 induced by the floor/mod operations inside $n(\theta)$. In practice, we implement reading and writing
 135 via gather/scatter on the at-most- $2k$ active indices, which is efficient and supports backpropagation
 136 through the unrolled computation. Early stopping introduces a discrete dependence on the stopping
 137 time $T(x)$; a standard choice is to stop the forward pass when $g_t < \varepsilon$ and treat the control-flow
 138 decision as non-differentiable, while gradients still flow through all executed steps (alternatively,
 139 one can always run for T_{\max} steps and rely on the multiplicative g_t factors to effectively mask later
 140 updates).

141 **Concrete parameterization (used in experiments).** We instantiate the maps $\mathcal{M}_T, \mathcal{M}_Q, \mathcal{M}_I$ as
 142 linear projections, and the transition maps $\Delta_T, \Delta_Q, \Delta_w$ as two-layer MLPs with expansion factor 4.
 143 Specifically, we project tape symbols position-wise,

$$\mathcal{M}_T(x; W) = xW_T, \quad \mathcal{M}_T(z; W) = zW_T,$$

144 and define $\mathcal{M}_Q, \mathcal{M}_I$ as learnable linear maps that collapse the sequence to the required shapes.
 145 Writing

$$\text{vec}(x) := [x[1]; x[2]; \dots; x[N]] \in \mathbb{R}^{Nd_x},$$

146 we set

$$\mathcal{M}_Q(x; W) = W_Q \text{vec}(x) \in \mathbb{R}^{d_Q}, \quad \mathcal{M}_I(x; W) = (\boldsymbol{\theta}_0, \mathbf{w}_0),$$

147 with

$$\boldsymbol{\theta}_0 = (W_\theta \text{vec}(x)) \bmod 2\pi \in (S^1)^k, \quad \mathbf{w}_0 = W_w \text{vec}(x) \in \mathbb{R}^k.$$

148 No constraint is imposed on \mathbf{w}_0 ; weights may be any real numbers.

149 We choose $\kappa(x; W)$ as a two-layer MLP (expansion factor 4) with a positivity constraint so that
 150 $\kappa(x; W) > 0$. For Δ_θ , we parameterize periodicity by feeding $\sin(\theta_t)$ and $\cos(\theta_t)$ into an MLP;
 151 concretely,

$$\Delta_\theta(S_t, Q_t, \boldsymbol{\theta}_t; W) = \text{MLP}_\theta([S_t, Q_t, \sin(\theta_t), \cos(\theta_t)]) \in \mathbb{R}^k.$$

152 **Algorithm (forward pass).** Given (T_0, Q_0, I_0) , we iterate for $t = 0, 1, \dots, T_{\max} - 1$:

- 153 1. compute $(n_{t,i}, s_{t,i}, n_{t,i}^+)_{i=1}^k$ from $\boldsymbol{\theta}_t$ via the definitions above;
- 154 2. read S_t as a $2k$ -term weighted sum of gathered tape vectors;
- 155 3. compute g_t ; if $g_t < \varepsilon$, stop early and return T_t ;
- 156 4. update Q_{t+1} and update $(\boldsymbol{\theta}_{t+1}, \mathbf{w}_{t+1})$;
- 157 5. write by scatter-adding into the at-most- $2k$ tape locations $\{n_{t,i}, n_{t,i}^+\}_{i=1}^k$ according to (6);

158 We return $V(x) = T_{T(x)}$.



Figure 2: **Placeholder.** Block-swap experiment: a fixed Llama-style macro architecture where the per-layer computational block is one of {Attention, LSTM, NTM/DNC, VECTUR, VECSTUR}.

159 4 VecTur Blocks inside Llama-style Models

160 4.1 Macro architecture

161 We adopt a standard decoder-only transformer macro architecture (token embeddings, positional
162 encoding, residual blocks, and an LM head) following Llama-family designs (Touvron et al., 2023;
163 Dubey et al., 2024). We then vary the *block* inside each residual layer while keeping parameter count
164 and FLOPs roughly matched.

165 4.2 Compared blocks

166 We compare the following blocks:

- 167 • **Attention block:** multi-head self-attention + SwiGLU MLP (Vaswani et al., 2017).
- 168 • **LSTM block:** a gated recurrent update applied over the sequence, wrapped with residual connec-
169 tions (Hochreiter and Schmidhuber, 1997).
- 170 • **External-memory block:** an NTM/DNC-style controller with differentiable read/write heads
171 (Graves et al., 2014, 2016).
- 172 • **VECTUR block:** the VECTUR transition unrolled for T_{\max} steps with learned halting κ .
- 173 • **VECSTUR block:** VECTUR with stochastic symbols z .

174 5 Evaluation Benchmarks

175 5.1 Reasoning and knowledge

176 We evaluate few-shot or fine-tuned performance on:

- 177 • **GSM8K** (Cobbe et al., 2021) (grade-school math; exact-match accuracy),
- 178 • **ARC** (Clark et al., 2018) (AI2 reasoning challenge; accuracy),
- 179 • **HellaSwag** (Zellers et al., 2019) (commonsense completion; accuracy).

180 5.2 Language modeling

181 We evaluate next-token prediction on **WikiText-103** (Merity et al., 2016) using perplexity.

182 6 CompGen: Complexity-Stratified Program Generation

183 6.1 Task format

184 COMPGEN consists of short Python programs p paired with inputs u and outputs $p(u)$. Each instance
185 is labeled with a target complexity class $(T(n), S(n))$ in terms of input size n . Programs are generated
186 from templates with controlled loop structure, recursion depth, and memory allocation patterns.

Class	Example family	Notes
$O(n), O(1)$	scan / reduce	single pass
$O(n), O(n)$	prefix sums	linear auxiliary array
$O(n \log n), O(1)$	sort-then-scan	comparison sorting
$O(n^2), O(1)$	nested-loop count	quadratic time
$O(n^2), O(n)$	DP table strip	quadratic time, linear space

Table 1: **Placeholder.** COMPGEN program families and intended $(T(n), S(n))$ buckets.

Block	GSM8K	ARC	HellaSwag	WikiText PPL
Attention	42.1	54.0	78.3	18.7
LSTM	38.4	51.2	76.0	20.9
NTM/DNC	43.0	54.4	78.1	19.3
VECTUR	46.8	56.1	79.0	18.9
VECSTUR	46.2	55.7	78.9	18.9

Table 2: **Placeholder.** Benchmark performance for a mid-size model at matched parameter budget. Accuracy in %, perplexity lower is better.

187 6.2 Generalization protocol

- 188 We train on $n \in [n_{\min}, n_{\text{train}}]$ and evaluate on larger $n \in (n_{\text{train}}, n_{\text{test}}]$ to measure extrapolation.
 189 We report accuracy as a function of n and correlate effective compute (average unroll steps) with
 190 complexity class.

191 7 Randomized Computation Suite

- 192 We include tasks where access to randomness enables provable or empirical speedups:
 193 • **Matrix product verification** (Freivalds) (Freivalds, 1977): verify $AB = C$ faster than multiplication.
 194 • **Probabilistic primality testing** (Miller–Rabin) (Miller, 1976; Rabin, 1980): decide primality with
 195 bounded error.
 196 VECSTUR receives stochastic symbols z and learns to leverage them to reduce expected compute (as
 197 reflected by learned κ and early halting).

199 8 Experimental Setup

- 200 **Model sizes.** We instantiate models at $\sim 110\text{M}$, 350M , and 1.3B parameters (placeholder sizes)
 201 with matched embedding width and layer count across blocks.

- 202 **Training.** We train on a mixture of general text (for language modeling) and COMPGEN (for
 203 compute probing), then fine-tune on downstream reasoning tasks. We use identical optimizers,
 204 learning rate schedules, and token budgets across conditions.

- 205 **Compute control.** For VECTUR/VECSTUR we set a maximum unroll T_{\max} and learn $\kappa(x; W)$ to
 206 modulate effective steps. We report both task performance and measured compute (average unroll
 207 steps per token).

208 9 Results (Illustrative Placeholders)

- 209 **Important note.** The tables below contain **illustrative placeholder numbers** showing the intended
 210 presentation format. Replace these with actual experimental results.



Figure 3: **Placeholder.** COMPGEN extrapolation: accuracy vs. input size n , showing VECTUR degrades more gracefully and increases effective steps via learned κ .

Model	Freivalds verify	Miller–Rabin
VECTUR	71.0	68.4
VECSTUR	84.5	79.2

Table 3: **Placeholder.** Randomized computation suite: VECSTUR benefits from stochastic symbols z .

211 10 Discussion

212 These illustrative results suggest VECTUR provides a useful inductive bias for tasks requiring
 213 iterative computation and length extrapolation, while remaining compatible with modern LLM macro
 214 architectures. VECSTUR further improves performance on tasks where randomized strategies are
 215 advantageous.

216 11 Limitations and Future Work

217 This draft omits implementation details (e.g., the exact Sparse(.) operator, stability constraints, and
 218 efficient kernels) and uses illustrative results. Future work should (i) benchmark on longer-context
 219 settings, (ii) analyze failure modes of learned halting κ , and (iii) evaluate robustness across different
 220 data mixtures and training budgets.

221 11.1 Future Work: Mechanistic Interpretability

222 VECTUR is unusually well-suited for mechanistic interpretability because its learned dynamics are
 223 constrained to resemble an explicit Turing-style transition system: a finite-dimensional control state
 224 Q_t , a tape T_t , and a small number of heads with sparse, local read/write effects. This structure
 225 encourages explanations in terms of *state machines* and *pointer-based algorithms* (e.g., “scan until
 226 condition,” “increment counter,” “copy span,” “simulate update rule”), rather than opaque global
 227 attention patterns.

228 A promising direction is to *disassemble* trained VECTUR blocks into more directly inspectable
 229 artifacts. For example, one can post-hoc discretize head locations, identify stable control states, and
 230 summarize the transition maps Δ as a symbolic program or a finite set of guarded update rules;
 231 such representations can then be *transpiled* into executable code, enabling unit tests, counterfactual
 232 interventions, and formal analysis of the implied algorithm.

233 Finally, VECTUR may serve as an interpretable *surrogate* for black-box sequence models. Analogous
 234 to knowledge distillation, one can perform *cross-distillation*: train a VECTUR model to mimic the
 235 input–output behavior (and, when available, internal activations) of an existing architecture, with the
 236 goal that the learned tape-and-control dynamics provide a concrete hypothesis for the black box’s
 237 implicit Turing-style computation. Such surrogates could support “algorithmic guessing”—extracting
 238 candidate programs from the VECTUR dynamics—followed by validation against the teacher via
 239 targeted probes and adversarial test cases.

240 **Acknowledgments**

241 *Placeholder.*

242 **References**

- 243 Alan M. Turing. On computable numbers, with an application to the Entscheidungsproblem. *Proceedings of the London Mathematical Society*, 1936.
- 244 Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural Computation*, 1997.
- 245 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *NeurIPS*, 2017.
- 246 Alex Graves, Greg Wayne, and Ivo Danihelka. Neural Turing machines. *arXiv:1410.5401*, 2014.
- 247 Alex Graves, Greg Wayne, Malcolm Reynolds, Tim Harley, Ivo Danihelka, Agnieszka Grabska-Barwińska, Sergio Gómez Colmenarejo, Edward Grefenstette, Tiago Ramalho, John Agapiou, Adrià Puigdomènech Badia, Karl Moritz Hermann, Yori Zwols, Georg Ostrovski, Adam Cain, Helen King, Christopher Summerfield, Phil Blunsom, Koray Kavukcuoglu, and Demis Hassabis. Hybrid computing using a neural network with dynamic external memory. *Nature*, 2016.
- 248 Alex Graves. Adaptive computation time for recurrent neural networks. *arXiv:1603.08983*, 2016.
- 249 Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, et al. Llama 2: Open foundation and fine-tuned chat models. *arXiv:2307.09288*, 2023.
- 250 Abhimanyu Dubey et al. The Llama 3 herd of models. *arXiv preprint*, 2024.
- 251 Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukas Kaiser, Matthias Plappert, et al. Training verifiers to solve math word problems. *arXiv:2110.14168*, 2021.
- 252 Peter Clark, Isaac Cowhey, Oren Etzioni, Tushar Khot, Ashish Sabharwal, Carissa Schoenick, and Oyyvind Tafjord. Think you have solved question answering? try ARC, the AI2 reasoning challenge. *arXiv:1803.05457*, 2018.
- 253 Rowan Zellers, Yonatan Bisk, Roy Schwartz, and Yejin Choi. HellaSwag: Can a machine really finish your sentence? In *ACL*, 2019.
- 254 Stephen Merity, Caiming Xiong, James Bradbury, and Richard Socher. Pointer sentinel mixture models. *arXiv:1609.07843*, 2016. (Introduces the WikiText-103 benchmark.)
- 255 Raimund Freivalds. Probabilistic machines can use less running time. In *IFIP Congress*, 1977.
- 256 Gary L. Miller. Riemann’s hypothesis and tests for primality. *Journal of Computer and System Sciences*, 1976.
- 257 Michael O. Rabin. Probabilistic algorithm for testing primality. *Journal of Number Theory*, 1980.