
VecTur: Vector Turing Machines

Ethan Hall

ethan.hall.phd@gmail.com

Abstract

We introduce VECTUR (Vector Turing Machines), a differentiable, Turing-machine-inspired transition system that represents tape symbols, head position, and finite control as continuous vectors (Turing, 1936). Conceptually, VECTUR continues a well-trodden line of differentiable memory machines (e.g., Neural Turing Machines and Differentiable Neural Computers) (Graves et al., 2014, 2016); our focus is a modern, sparse implementation that avoids dense content-based access at every step. VECTUR maintains a continuous “head on a circle” (S^1) and performs *strictly local*, $2k$ -sparse gather/scatter updates via interpolation, encouraging pointer-machine-style computation and mitigating the “memory blurring” failure mode of early dense-access NTMs. VECTUR also supports *deep computation* (Dehghani et al., 2019) by explicitly iterating a learned transition map and using an ACT-style learned halting gate (Graves, 2016) (parameterized by κ); we treat this halting parameterization as a practical heuristic rather than a core novelty claim. 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 (Kaiser and Sutskever, 2016; Press et al., 2022) at fixed parameter budgets. We additionally introduce COMPGEN—a synthetic program dataset stratified by $(T(n), S(n))$ complexity classes—as a utility benchmark for probing compute allocation. Finally, we propose VECSTUR, a stochastic extension that consumes random tape symbols and targets *randomized algorithms* as a computational resource: VECSTUR outperforms VECTUR on randomized verification tasks (e.g., Freivalds-style matrix product verification (Freivalds, 1977)).

28

1 Introduction

29 Modern language models excel at pattern completion yet often struggle to reliably *execute* long
30 algorithmic computations, extrapolate beyond training lengths (Press et al., 2022), or allocate variable
31 compute per input (Graves, 2016; Dehghani et al., 2019). Several lines of work attempt to address
32 these limitations by embedding algorithmic structure into neural systems, including external-memory
33 architectures (Graves et al., 2014, 2016) and adaptive computation mechanisms (Graves, 2016).

34 We propose VECTUR, a vectorized analogue of a classical Turing machine (Turing, 1936) whose
35 tape, head index, and finite control are represented as continuous vectors and updated by a learned
36 transition map. We do *not* claim to have invented differentiable Turing machines; rather, we revisit
37 this classic idea in a form that better matches modern LLM systems constraints. Classic NTMs
38 (Graves et al., 2014) relied on dense, content-based attention over the entire memory at each step,
39 which is $O(N)$ in memory size and can introduce diffuse “blurring” updates. In contrast, VECTUR
40 maintains a continuous head position on a circular tape (S^1) and enforces *strictly local* access: each

41 step reads and writes via interpolation over only $2k$ tape indices (sparse gather/scatter), yielding
42 per-step cost independent of tape length and an inductive bias closer to pointer machines (Vinyals
43 et al., 2015). For adaptive depth, VECTUR includes an ACT-style halting mechanism (Graves, 2016);
44 our particular κ parameterization is presented as a practical, input-conditioned control knob rather
45 than a conceptual departure from ACT.

46 **Why this helps for long context.** Self-attention provides direct interaction between all token pairs
47 but incurs $O(N^2)$ compute in sequence length N . Linear-time alternatives avoid this quadratic cost,
48 but typically do so by enforcing a *linear progression* of computation through the sequence: in an
49 LSTM, information from the past must be compressed into a fixed-size hidden state, and each new
50 token updates that state once in time order. More recent linear-time memory perspectives (e.g., the
51 MIRAS viewpoint connecting retention, memorization, and online optimization (Behrouz et al.,
52 2025)) likewise maintain and update memory in lockstep with the token stream. In contrast, VECTUR
53 decouples interaction from all-pairs attention by introducing an explicit *Turing index* (the head index
54 I_t) whose motion over the tape is learned. Across transition steps, head motion allows the model
55 to create relationships between *arbitrary* tokens by moving to (or copying from) the corresponding
56 tape locations and combining them through the finite-control state, without computing a dense
57 $N \times N$ similarity matrix. Crucially, this enables *non-linear* token processing: the controller can
58 revisit selected tokens or intermediate scratchpad cells multiple times, in an order determined by the
59 evolving machine state, while leaving the rest of the tape unchanged by construction. This selective
60 persistence—deciding what to write to the tape, what to overwrite, and what to ignore—matches the
61 operational structure of Turing-style computation (state + tape) more directly than recurrence with a
62 single compressed memory vector.

63 We evaluate VECTUR in a realistic regime by inserting it as a *block* inside a Llama-style decoder-only
64 macro architecture (Touvron et al., 2023; Dubey et al., 2024), replacing the standard attention+MLP
65 block. We compare against alternative blocks: (i) standard attention (Vaswani et al., 2017), (ii)
66 LSTM-style recurrence (Hochreiter and Schmidhuber, 1997), and (iii) differentiable external-memory
67 controllers (Graves et al., 2014, 2016). We focus on small-to-medium models (roughly 10^8 to
68 10^9 parameters) where architectural inductive bias can materially affect sample efficiency and
69 extrapolation (Kaiser and Sutskever, 2016; Press et al., 2022).

70 We further introduce COMPGEN, a dataset of generated Python programs grouped into discrete com-
71 plexity buckets $(T(n), S(n))$ such as $O(n)/O(1)$, $O(n \log n)/O(1)$, $O(n^2)/O(1)$, and $O(n^2)/O(n)$.
72 The goal is to probe whether a model can learn to *compute* across increasing n by allocating more
73 steps as needed, rather than memorizing only small n . Finally, we define VECSTUR, which aug-
74 ments the input with stochastic symbols z to emulate randomized computation, and we propose a
75 randomized evaluation suite where randomness yields asymptotic speedups (Freivalds, 1977; Miller,
76 1976; Rabin, 1980). In particular, we form a data set of matrices $A, B, C \in \mathbb{R}^{n \times n}$ and a target
77 matrix $D \in \mathbb{R}^{n \times n}$ such that $D = AB$ and $D = AC$ with probability $1/2$. We evaluate VECTUR and
78 VECSTUR on this task, and show VECSTUR can exploit stochastic symbols to achieve asymptotic
79 speedups.

80 **Contributions.**

- 81 • We define VECTUR, a Turing-style transition system with *strictly local*, $2k$ -sparse gather/scatter
82 tape access (continuous head on S^1), addressing efficiency and “memory blurring” issues associated
83 with dense-access NTMs (Graves et al., 2014).
- 84 • We present a plug-and-play integration of VECTUR as a *computational sub-layer* inside Llama-
85 style decoder-only models, treating iterative algorithmic computation as a composable block rather
86 than a separate retrieval module.
- 87 • We define VECSTUR and a randomized computation evaluation suite, highlighting randomness
88 as a computational resource (e.g., Freivalds-style verification (Freivalds, 1977)) rather than mere
89 noise.
- 90 • We introduce COMPGEN, a synthetic program dataset labeled by time/space complexity class
91 $(T(n), S(n))$, as a utility benchmark for extrapolation and compute-allocation probing.

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 (Vinyals et al., 2015)	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.

92 2 Related Work

93 **Sequence models and attention.** Transformers (Vaswani et al., 2017) and their decoder-only
 94 variants power modern LLMs (e.g., GPT-3 (Brown et al., 2020) and Llama-family models (Touvron
 95 et al., 2023; Dubey et al., 2024)). Recurrent networks such as LSTMs (Hochreiter and Schmidhuber,
 96 1997) provide a different inductive bias for iterative computation but historically underperform
 97 attention-based models at scale on language modeling.

98 **Differentiable memory and neural machines.** Neural Turing Machines (NTMs) (Graves et al.,
 99 2014) and Differentiable Neural Computers (DNCs) (Graves et al., 2016) integrate external memory
 100 with differentiable read/write heads. Our work shares the goal of improving algorithmic behavior,
 101 but emphasizes sparse, strictly local access (pointer-machine-style) and composable integration as a
 102 modern Transformer block; we include learned halting primarily as an ACT-style compute control
 103 mechanism.

104 **2.1 Remark: Neural Turing Machines vs. VECTUR**

105 Both NTMs (Graves et al., 2014) and VECTUR augment neural computation with an external memory,
 106 but they make different design trade-offs for *addressing* (how memory is accessed) and *sparsity*
 107 (how much memory is touched per step). NTMs provide content-addressable reads/writes via dense
 108 attention over all memory slots; VECTUR instead maintains a continuous head position on a tape and
 109 performs sparse gather/scatter updates to a small number of adjacent cells (via interpolation), which
 110 keeps per-step cost independent of tape length. Dense access is expressive but $O(N)$ in memory size
 111 and can induce diffuse “memory blurring” updates when many slots receive small writes; VECTUR
 112 preserves untouched cells exactly by construction. In Llama-style Transformer blocks, global content-
 113 based access is already available through self-attention; VECTUR is intended to add an *orthogonal*
 114 capability: cheap, iterative state manipulation with persistent scratchpad dynamics.

115 **Adaptive computation.** Adaptive Computation Time (ACT) (Graves, 2016) learns when to stop
 116 iterating, with later refinements such as PonderNet (Banino et al., 2021). Our learned gate g_t and κ -
 117 parameterization should be viewed as an ACT-style variant that provides a simple, input-conditioned
 118 control knob for effective depth; we do not position halting as the primary conceptual novelty.

119 **Test-time training and online optimization.** Recent work reframes sequence modeling as a form
 120 of *online learning* or nested optimization carried out during inference, including end-to-end test-
 121 time training for long-context language modeling (Tandon et al., 2025) and the MIRAS framework

122 connecting attention, retention, and online optimization (Behrouz et al., 2025). This line of work
 123 motivates the viewpoint that “System 2” computation can be injected *inside* a model by adding
 124 inner-loop dynamics as a composable block within the forward pass, rather than only via external
 125 deliberation or separate modules.

126 **Randomized algorithms.** Randomness can reduce expected runtime for verification and decision
 127 problems; canonical examples include Freivalds’ randomized matrix product verification (Freivalds,
 128 1977) and probabilistic primality testing (Miller, 1976; Rabin, 1980). VECSTUR is intended as a
 129 neural analogue that can exploit stochastic symbols during computation.

130 3 VecTur: Vector Turing Machines

131 3.1 Vectorized machine state

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

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

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

137 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
 138 state $Q_t \in \mathbb{R}^{d_Q}$, and the head index

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

139 are learned, differentiable quantities. (Here $\theta_t = (\theta_{t,1}, \dots, \theta_{t,k})$ parameterizes k head locations
 140 on the circle and $\mathbf{w}_t = (w_{t,1}, \dots, w_{t,k})$ are the associated weights.) We index tape positions by
 141 $j \in \{0, 1, \dots, N_T - 1\}$, and write $T_t[j] \in \mathbb{R}^{d_T}$ for the j -th tape symbol.

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

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

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

144 3.2 Sparse addressing (keeping I and J fixed)

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

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

146 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
 147 vector $J(I) \in \mathbb{R}^{N_T}$ by

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

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

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

150 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
 151 $(n_i, n_i^+, (1 - s_i)w_i, s_i w_i)$ for each i and never materialize the dense N_T -vector $J(I)$.

152 **3.3 Read, transition, and halting**

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

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

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

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

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

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

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

161 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
 162 halting target state. Intuitively, $\kappa(x; W)$ acts as a *decay-rate multiplier*: smaller $\kappa(x; W)$ yields a
 163 slower decay in t (more effective steps), while larger $\kappa(x; W)$ yields a faster decay (fewer effective
 164 steps). The factor $\|Q_t - q_0\|$ encourages the dynamics to become stationary near the target.

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

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

167 Then the update equations are

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

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

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

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

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

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

171 The transition maps have the following types:

$$\begin{aligned} \Delta_T : \mathbb{R}^{d_T} \times \mathbb{R}^{d_Q} &\rightarrow \mathbb{R}^{d_T}, \\ \Delta_Q : \mathbb{R}^{d_T} \times \mathbb{R}^{d_Q} &\rightarrow \mathbb{R}^{d_Q}, \\ \Delta_\theta : \mathbb{R}^{d_T} \times \mathbb{R}^{d_Q} \times (S^1)^k &\rightarrow \mathbb{R}^k, \\ \Delta_w : \mathbb{R}^{d_T} \times \mathbb{R}^{d_Q} \times \mathbb{R}^k &\rightarrow \mathbb{R}^k. \end{aligned}$$

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

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

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

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

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

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

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

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

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

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

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

196 we set

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

197 with

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

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

199 We choose $\kappa(x; W)$ as a two-layer MLP (expansion factor 4) with a positivity constraint so that
 200 $\kappa(x; W) > 0$. For Δ_θ , we parameterize periodicity by feeding $\sin(\theta_t)$ and $\cos(\theta_t)$ into an MLP;
 201 concretely,

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

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

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

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

209 4 VecTur Blocks inside Llama-style Models

210 4.1 Macro architecture

211 We adopt a standard decoder-only transformer macro architecture (token embeddings, positional
 212 encoding (Su et al., 2021), residual blocks, and an LM head) following Llama-family designs

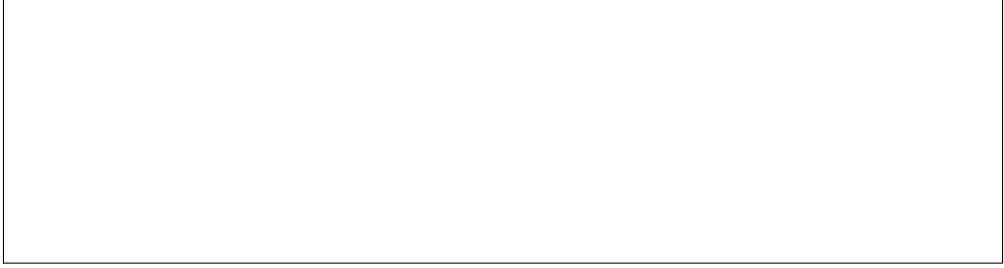


Figure 2: **Placeholder.** Block-swap experiment: a fixed Llama-style macro architecture where the per-layer computational block is one of {Attention, LSTM, NTM/DNC, VECTUR, VECSTUR}.

(Touvron et al., 2023; Dubey et al., 2024). We then vary the *block* inside each residual layer while keeping parameter count and FLOPs roughly matched. This “block as inner loop” framing is inspired by recent work that integrates deliberate, multi-step computation into the forward pass via online learning or test-time adaptation, notably TTT-style test-time training (Tandon et al., 2025) and MIRAS-style online optimization views of sequence models (Behrouz et al., 2025). In that spirit, we view VECTUR as an explicit, constrained “System 2” transition system embedded as a sub-layer inside a “System 1” decoder, rather than as a standalone memory system that replaces the macro architecture.

4.2 Compared blocks

We compare the following blocks:

- **Attention block:** multi-head self-attention (Vaswani et al., 2017) + SwiGLU MLP (Shazeer, 2020).
- **LSTM block:** a gated recurrent update applied over the sequence, wrapped with residual connections (Hochreiter and Schmidhuber, 1997).
- **External-memory block:** an NTM/DNC-style controller with differentiable read/write heads (Graves et al., 2014, 2016).
- **VECTUR block:** the VECTUR transition unrolled for T_{\max} steps with learned halting κ .
- **VECSTUR block:** VECTUR with stochastic symbols z .

5 Evaluation Benchmarks

5.1 Reasoning and knowledge

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

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

5.2 Language modeling

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

6 CompGen: Complexity-Stratified Program Generation

6.1 Task format

COMPGEN consists of short Python programs p paired with inputs u and outputs $p(u)$. Each instance is labeled with a target complexity class $(T(n), S(n))$ in terms of input size n (Sipser, 2012). Programs are generated from templates with controlled loop structure, recursion depth, and memory allocation patterns. We view COMPGEN as a utility dataset in the tradition of synthetic algorithmic benchmarks, complementary to the CLRS Algorithmic Reasoning Benchmark (Veličković et al., 2022).

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.

246 6.2 Generalization protocol

247 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.
 248 We report accuracy as a function of n and correlate effective compute (average unroll steps) with
 249 complexity class.

250 7 Randomized Computation Suite

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

258 8 Experimental Setup

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

261 **Blocks and controlled comparisons.** Unless otherwise stated, we run the same experiment for each
 262 block in Section 3 (Attention, LSTM, NTM/DNC, VECTUR, VECSTUR), holding the decoder-only
 263 macro architecture fixed and matching parameter count and training budget as closely as possible.

264 **Experimental protocols (run per block).** We use three complementary training/evaluation proto-
 265 cols:

- 266 1. **Language pretraining → downstream evaluation.** We pretrain on **FineWeb** (general web text),
 267 then evaluate on **GSM8K** (Cobbe et al., 2021), **ARC** (Clark et al., 2018), **HellaSwag** (Zellers
 268 et al., 2019), and **WikiText-103** (Merity et al., 2016). (Table 2.)
- 269 2. **Algorithmic transfer between CLRS and COMPGEN.** (a) **Train on CLRS** (Veličković et al.,
 270 2022) and evaluate on COMPGEN under three regimes: *zero-shot* (no COMPGEN training), *few-*
 271 *shot* (in-context demonstrations at test time), and *fine-tune* (supervised adaptation on COMPGEN
 272 train). (b) **Train on COMPGEN** and evaluate on a held-out COMPGEN split (including out-of-
 273 distribution generalization across input sizes n per Section 5.2). (Figure 3.)
- 274 3. **In-domain CLRS generalization.** We train on CLRS (Veličković et al., 2022) and evaluate
 275 on a held-out CLRS split (standard in-distribution generalization across graphs/sizes/instances).
 276 (Reported alongside other algorithmic results; placeholder in this draft.)

277 **Optimization and budgets.** Within each protocol, we use identical optimizers, learning rate sched-
 278 ules, and token/step budgets across blocks (to isolate architectural effects).

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

Block (model)	Train set	GSM8K (test)	ARC (test)	HellaSwag (test)	WikiText-103 (test)
Attention	FineWeb	Lorem	Ipsum	Dolor	Sit
LSTM	FineWeb	Amet	Consectetur	Adipiscing	Elit
NTM/DNC	FineWeb	Sed	Do	Eiusmod	Tempor
VECTUR	FineWeb	Incididunt	Ut	Labore	Et
VECSTUR	FineWeb	Magna	Aliqua	Ut	Enim

Table 2: **Placeholder (Protocol 1).** Language pretraining on FineWeb, evaluated on downstream benchmarks. Entries are Lorem ipsum placeholders.

Block	Train	Test	Result
Attention	CLRS	COMPGEN (zero-shot)	Lorem ipsum
Attention	CLRS	COMPGEN (few-shot)	Dolor sit
Attention	CLRS	COMPGEN (fine-tune)	Amet consectetur
LSTM	CLRS	COMPGEN (zero-shot)	Adipiscing elit
LSTM	CLRS	COMPGEN (few-shot)	Sed do
LSTM	CLRS	COMPGEN (fine-tune)	Eiusmod tempor
NTM/DNC	CLRS	COMPGEN (zero-shot)	Incididunt ut
NTM/DNC	CLRS	COMPGEN (few-shot)	Labore et
NTM/DNC	CLRS	COMPGEN (fine-tune)	Magna aliqua
VECTUR	CLRS	COMPGEN (zero-shot)	Ut enim
VECTUR	CLRS	COMPGEN (few-shot)	Ad minim
VECTUR	CLRS	COMPGEN (fine-tune)	Veniam quis
VECSTUR	CLRS	COMPGEN (zero-shot)	Nostrud exercitation
VECSTUR	CLRS	COMPGEN (few-shot)	Ullamco laboris
VECSTUR	CLRS	COMPGEN (fine-tune)	Nisi ut

Table 3: **Placeholder (Protocol 2a).** Train on CLRS, test on COMPGEN under zero-shot / few-shot / fine-tune adaptation regimes. Results are placeholders.

282 9 Results (Illustrative Placeholders)

283 **Important note.** The tables below contain **LOREM IPSUM placeholder entries** showing the intended presentation format *and* explicitly recording the train/test split for each experiment protocol.
284 Replace these placeholders with measured metrics.

286 10 Discussion

287 These illustrative results suggest VECTUR provides a useful inductive bias for tasks requiring
288 iterative computation and length extrapolation, while remaining compatible with modern LLM
289 macro architectures. Importantly, the strongest claims in this paper are *not* that differentiable Turing
290 machines are new, but that (i) enforcing strictly local sparse access yields a practical, non-blurring
291 pointer-machine-style block, (ii) treating such a machine as a composable Transformer sub-layer
292 is a strong systems contribution, and (iii) VECSTUR highlights a comparatively underexplored
293 angle: learning to exploit randomness as a computational resource in randomized-algorithm tasks.
294 VECSTUR further improves performance on tasks where randomized strategies are advantageous.

295 11 Limitations and Future Work

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

Block	Train	Test	Result
Attention	COMPGEN (train)	COMPGEN (held-out)	Lorem ipsum
LSTM	COMPGEN (train)	COMPGEN (held-out)	Dolor sit
NTM/DNC	COMPGEN (train)	COMPGEN (held-out)	Amet consectetur
VECTUR	COMPGEN (train)	COMPGEN (held-out)	Adipiscing elit
VECSTUR	COMPGEN (train)	COMPGEN (held-out)	Sed do

Table 4: **Placeholder (Protocol 2b).** Train on COMPGEN, test on held-out COMPGEN (including extrapolation across larger n). Results are placeholders.

Block	Train	Test	Result
Attention	CLRS (train)	CLRS (held-out)	Lorem ipsum
LSTM	CLRS (train)	CLRS (held-out)	Dolor sit
NTM/DNC	CLRS (train)	CLRS (held-out)	Amet consectetur
VECTUR	CLRS (train)	CLRS (held-out)	Adipiscing elit
VECSTUR	CLRS (train)	CLRS (held-out)	Sed do

Table 5: **Placeholder (Protocol 3).** Train on CLRS and evaluate on a held-out CLRS split. Results are placeholders.

300 11.1 Future Work: Mechanistic Interpretability

301 VECTUR is unusually well-suited for mechanistic interpretability (Olah et al., 2020; Elhage et al.,
 302 2021) because its learned dynamics are constrained to resemble an explicit Turing-style transition
 303 system: a finite-dimensional control state Q_t , a tape T_t , and a small number of heads with sparse,
 304 local read/write effects. This structure encourages explanations in terms of *state machines* and
 305 *pointer-based algorithms* (e.g., “scan until condition,” “increment counter,” “copy span,” “simulate
 306 update rule”), rather than opaque global attention patterns.

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

312 Finally, VECTUR may serve as an interpretable *surrogate* for black-box sequence models. Analogous
 313 to knowledge distillation (Hinton et al., 2015; Romero et al., 2015), one can perform *cross-distillation*:
 314 train a VECTUR model to mimic the input–output behavior (and, when available, internal activations)
 315 of an existing architecture, with the goal that the learned tape-and-control dynamics provide a concrete
 316 hypothesis for the black box’s implicit Turing-style computation. Such surrogates could support
 317 “algorithmic guessing”—extracting candidate programs from the VECTUR dynamics—followed by
 318 validation against the teacher via targeted probes and adversarial test cases.

319 Acknowledgments

320 *Placeholder.*

321 References

- 322 Alan M. Turing. On computable numbers, with an application to the Entscheidungsproblem. *Proceedings of the London Mathematical Society*, 1936.
- 324 Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural Computation*, 1997.
- 325 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Łukasz
 326 Kaiser, and Illia Polosukhin. Attention is all you need. In *NeurIPS*, 2017.
- 327 Alex Graves, Greg Wayne, and Ivo Danihelka. Neural Turing machines. *arXiv:1410.5401*, 2014.



Figure 3: **Placeholder.** COMPGEN extrapolation: accuracy vs. input size n , showing how blocks degrade with larger n and how VECTUR modulates effective steps via learned κ .

Block	Train set	Test set	Result
VECTUR	Randomized suite (train)	Freivalds / Miller–Rabin (test)	Lorem ipsum
VECSTUR	Randomized suite (train)	Freivalds / Miller–Rabin (test)	Dolor sit amet

Table 6: **Placeholder.** Randomized computation suite: train/test bookkeeping with placeholder results.

- 328 Alex Graves, Greg Wayne, Malcolm Reynolds, Tim Harley, Ivo Danihelka, Agnieszka Grabska-
 329 Barwińska, Sergio Gómez Colmenarejo, Edward Grefenstette, Tiago Ramalho, John Agapiou,
 330 Adrià Puigdomènech Badia, Karl Moritz Hermann, Yori Zwols, Georg Ostrovski, Adam Cain,
 331 Helen King, Christopher Summerfield, Phil Blunsom, Koray Kavukcuoglu, and Demis Hassabis.
 332 Hybrid computing using a neural network with dynamic external memory. *Nature*, 2016.
- 333 Alex Graves. Adaptive computation time for recurrent neural networks. *arXiv:1603.08983*, 2016.
- 334 Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay
 335 Bashlykov, et al. Llama 2: Open foundation and fine-tuned chat models. *arXiv:2307.09288*, 2023.
- 336 Abhimanyu Dubey et al. The Llama 3 herd of models. *arXiv preprint*, 2024.
- 337 Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukas Kaiser, Matthias
 338 Plappert, et al. Training verifiers to solve math word problems. *arXiv:2110.14168*, 2021.
- 339 Peter Clark, Isaac Cowhey, Oren Etzioni, Tushar Khot, Ashish Sabharwal, Carissa Schoenick, and
 340 Oyyvind Tafjord. Think you have solved question answering? try ARC, the AI2 reasoning challenge.
 341 *arXiv:1803.05457*, 2018.
- 342 Rowan Zellers, Yonatan Bisk, Roy Schwartz, and Yejin Choi. HellaSwag: Can a machine really
 343 finish your sentence? In *ACL*, 2019.
- 344 Stephen Merity, Caiming Xiong, James Bradbury, and Richard Socher. Pointer sentinel mixture
 345 models. *arXiv:1609.07843*, 2016. (Introduces the WikiText-103 benchmark.)
- 346 Raimund Freivalds. Probabilistic machines can use less running time. In *IFIP Congress*, 1977.
- 347 Gary L. Miller. Riemann’s hypothesis and tests for primality. *Journal of Computer and System
 348 Sciences*, 1976.
- 349 Michael O. Rabin. Probabilistic algorithm for testing primality. *Journal of Number Theory*, 1980.
- 350 Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal,
 351 Arvind Neelakantan, et al. Language models are few-shot learners. In *NeurIPS*, 2020.
- 352 Mostafa Dehghani, Stephan Gouws, Oriol Vinyals, Jakob Uszkoreit, and Łukasz Kaiser. Universal
 353 transformers. In *ICLR*, 2019.
- 354 Andrea Banino, Jan Balaguer, Charles Blundell, and Andrew Zisserman. PonderNet: Learning to
 355 ponder. In *NeurIPS*, 2021.

- 356 Ali Behrouz, Meisam Razaviyayn, Peilin Zhong, and Vahab Mirrokni. It’s All Connected: A
357 Journey Through Test-Time Memorization, Attentional Bias, Retention, and Online Optimization.
358 *arXiv:2504.13173*, 2025.
- 359 Nelson Elhage, Sam S. McCandlish, Catherine Olsson, Christopher Henighan, Nicholas Joseph,
360 Ben Mann, Seth Kaplan, et al. A mathematical framework for transformer circuits. *Transformer*
361 *Circuits (Anthropic)*, 2021.
- 362 Łukasz Kaiser and Ilya Sutskever. Neural GPUs learn algorithms. In *ICLR*, 2016.
- 363 Chris Olah, Nick Cammarata, Ludwig Schubert, Gabriel Goh, Michael Petrov, and the OpenAI
364 Clarity team. Zoom in: An introduction to circuits. *Distill*, 2020.
- 365 Ofir Press, Noah A. Smith, and Mike Lewis. Train short, test long: Attention with linear biases
366 enables input length extrapolation. In *ICLR*, 2022.
- 367 Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. Distilling the knowledge in a neural network.
368 *arXiv:1503.02531*, 2015.
- 369 Adriana Romero, Nicolas Ballas, Samira Ebrahimi Kahou, Antoine Chassang, Carlo Gatta, and
370 Yoshua Bengio. FitNets: Hints for thin deep nets. In *ICLR*, 2015.
- 371 Noam Shazeer. GLU variants improve transformer. *arXiv:2002.05202*, 2020.
- 372 Michael Sipser. *Introduction to the Theory of Computation*. Cengage Learning, 3rd edition, 2012.
- 373 Jianlin Su, Yu Lu, Shengfeng Pan, Bo Wen, and Yunfeng Liu. RoFormer: Enhanced transformer with
374 rotary position embedding. *arXiv:2104.09864*, 2021.
- 375 Oriol Vinyals, Meire Fortunato, and Navdeep Jaitly. Pointer networks. In *NeurIPS*, 2015.
- 376 Petar Veličković, Adrià Puigdomènech Badia, David Budden, Razvan Pascanu, Andrea Banino, Misha
377 Daswani, Raia Hadsell, and Charles Blundell. The CLRS Algorithmic Reasoning Benchmark.
378 *arXiv preprint arXiv:2205.15659*, 2022.
- 379 Arnuv Tandon, Karan Dalal, Xinhao Li, Daniel Koceja, Marcel Rød, Sam Buchanan, Xiaolong Wang,
380 Jure Leskovec, Sanmi Koyejo, Tatsunori Hashimoto, Carlos Guestrin, Jed McCaleb, Yejin Choi,
381 and Yu Sun. End-to-End Test-Time Training for Long Context. *arXiv:2512.23675*, 2025.