

# Differentiable Tensor Computers for End-to-End Program Synthesis: Architecture, Training Dynamics, and Path to Learned Software Agents

Jacob Valdez<sup>1</sup>

<sup>1</sup>AGI Inc., [jacob@theagi.company](mailto:jacob@theagi.company)

January 2026

## Abstract

This paper presents a differentiable tensor computer architecture implemented in JAX that executes standard machine code while remaining fully differentiable, enabling end-to-end gradient-based learning of programs. The system implements a complete von Neumann architecture—ALU, registers, cache, memory hierarchy, bus, peripherals, GPU, and display—entirely as tensor operations with soft attention mechanisms that preserve gradient flow. Full mathematical formalization is provided for all components: the softmax-gated ALU that computes all 32 operations in parallel and selects via temperature-controlled attention; the differentiable register file with soft read/write addressing; the cache hierarchy with learned LRU eviction policies; the hard disk emulation with 4KB sector loading; the memory management unit with two-level page tables and TLB; the system bus with softmax arbitration; and peripheral ring buffers with soft head/tail pointers. A self-hosting C compiler running on this architecture successfully compiles and executes recursive programs such as `factorial(5)=120`, validating correct execution of function calls, stack frames, conditionals, and recursion.

Beyond architecture, this paper presents a comprehensive training strategy addressing the fundamental challenges of gradient-based program synthesis: the temperature annealing schedule that transitions from high-temperature exploration (where gradients flow through soft operation mixtures) through annealing (gradual discretization) to crystallization (locking discrete structure) and refinement (local search in program space); curriculum learning over program length, input complexity, and task difficulty; shaped reward functions that provide dense signal for credit assignment including output distance metrics, auxiliary losses at intermediate execution points, and learned value functions; hierarchical program decomposition with subroutine libraries; and population-based training for solution diversity. The training pipeline combines imitation learning warmstart, outcome-supervised reinforcement learning, and discrete local search refinement.

The central hypothesis is that structured latent dynamics approximated by grokked transformers can be represented more efficiently as learned programs than as dense tensor operations. Evidence from mechanistic interpretability suggests that transformers trained to generalization implement sparse, algorithmic circuits—structure that could be more compactly encoded as imperative code. If correct, this enables deployment of capable agents at orders of magnitude lower cost than current LLM-based approaches: inference in microseconds rather than seconds, model sizes in kilobytes rather than gigabytes, and operation on mobile devices without cloud connectivity. The research program proposes experiments from classical algorithms through GUI-based computer control, with a \$1,000 compute budget on Modal H100s over two weeks.

# Contents

<b>1</b>	<b>Introduction</b>	<b>6</b>
1.1	The Architectural Mismatch . . . . .	6
1.2	Evidence from Mechanistic Interpretability . . . . .	6
1.3	Contributions . . . . .	6
1.4	Paper Organization . . . . .	7
<b>2</b>	<b>Background and Related Work</b>	<b>7</b>
2.1	Neural Computers and Differentiable Programming . . . . .	7
2.2	Grokking and Algorithmic Phase Transitions . . . . .	7
2.3	Program Synthesis . . . . .	8
2.4	Reinforcement Learning and Superhuman Discovery . . . . .	9
2.5	LLM-Based Agents and Their Limitations . . . . .	9
<b>3</b>	<b>Tensor Computer Architecture</b>	<b>10</b>
3.1	System Overview . . . . .	10
3.2	Notation and Conventions . . . . .	10
3.3	Machine State . . . . .	11
3.4	Instruction Encoding . . . . .	11
3.5	Instruction Fetch . . . . .	12
3.6	Register File Operations . . . . .	12
3.6.1	Soft Register Read . . . . .	12
3.6.2	Soft Register Write . . . . .	12
3.7	Arithmetic Logic Unit . . . . .	12
3.7.1	Operation Definitions . . . . .	13
3.7.2	Differentiable Bitwise Operations . . . . .	13
3.7.3	Differentiable Shifts . . . . .	13
3.7.4	ALU Output Computation . . . . .	14
3.7.5	Flag Computation . . . . .	15
3.8	Memory Hierarchy . . . . .	15
3.8.1	Cache Architecture . . . . .	15
3.8.2	Cache Lookup . . . . .	16
3.8.3	Cache Read . . . . .	16
3.8.4	Cache Write and Eviction . . . . .	16
3.8.5	NTM-Style Main Memory . . . . .	16
3.8.6	Hard Disk with 4KB Sector Loading . . . . .	17
3.8.7	Virtual Memory and Page Tables . . . . .	18
3.9	Control Flow . . . . .	19
3.9.1	Sequential Execution . . . . .	19
3.9.2	Conditional Branches . . . . .	19
3.9.3	Subroutine Calls . . . . .	19
3.10	Bus Architecture . . . . .	20
3.10.1	Bus Arbitration . . . . .	20
3.10.2	Memory-Mapped I/O . . . . .	20
3.11	GPU Interface . . . . .	20
3.11.1	GPU Commands . . . . .	20
3.11.2	GPU Execution . . . . .	20

3.12	Peripheral Ring Buffers . . . . .	21
3.12.1	Ring Buffer State . . . . .	21
3.12.2	Enqueue/Dequeue . . . . .	21
3.13	Display Buffer . . . . .	21
3.14	Computational Cost Analysis . . . . .	21
<b>4</b>	<b>Self-Hosting C Compiler</b>	<b>21</b>
4.1	Compiler Architecture . . . . .	22
4.2	Bootstrap Process . . . . .	22
4.3	Calling Convention . . . . .	22
4.4	Validation: Recursive Factorial . . . . .	23
4.5	Compiler Complexity . . . . .	23
<b>5</b>	<b>Training Strategy</b>	<b>23</b>
5.1	The Core Challenge: Discrete Structure via Continuous Optimization . . . . .	23
5.2	Temperature Dynamics . . . . .	23
5.2.1	Phase 1: Warm-up (High Temperature) . . . . .	24
5.2.2	Phase 2: Annealing (Decreasing Temperature) . . . . .	24
5.2.3	Phase 3: Crystallization (Low Temperature) . . . . .	24
5.2.4	Phase 4: Refinement (Discrete Local Search) . . . . .	25
5.2.5	Phase 5: Extraction and Verification . . . . .	25
5.3	Credit Assignment . . . . .	25
5.3.1	Shaped Rewards . . . . .	25
5.3.2	Auxiliary Losses . . . . .	26
5.3.3	Learned Value Function . . . . .	26
5.4	Curriculum Learning . . . . .	26
5.4.1	Curriculum Dimensions . . . . .	27
5.5	Hierarchical Program Structure . . . . .	27
5.5.1	Subroutine Library . . . . .	27
5.5.2	Hierarchical Policy . . . . .	27
5.6	Population-Based Training . . . . .	28
5.7	Handling Pathological Programs . . . . .	28
5.7.1	Infinite Loops . . . . .	28
5.7.2	Trivial Solutions . . . . .	28
5.8	Imitation Learning Warmstart . . . . .	29
5.9	Training Pipeline Summary . . . . .	29
<b>6</b>	<b>Theoretical Analysis</b>	<b>29</b>
6.1	Convergence Analysis . . . . .	29
6.1.1	Temperature-Parameterized Loss Landscape . . . . .	29
6.1.2	Gradient Flow Analysis . . . . .	30
6.1.3	Annealing Schedule Analysis . . . . .	30
6.2	Expressiveness and Complexity . . . . .	31
6.2.1	Universal Computation . . . . .	31
6.2.2	Learnability Bounds . . . . .	31
6.2.3	Expressible Function Classes . . . . .	31
6.3	Loss Landscape Characterization . . . . .	32
6.3.1	Basin Structure . . . . .	32

6.3.2	Plateau Analysis . . . . .	32
6.4	Generalization Bounds . . . . .	32
6.4.1	Program-Based Generalization . . . . .	32
6.4.2	Compositional Generalization . . . . .	33
<b>7</b>	<b>Preliminary Experiments</b>	<b>33</b>
7.1	Experimental Setup . . . . .	33
7.1.1	Implementation . . . . .	33
7.1.2	Training Configuration . . . . .	33
7.2	Task Suite . . . . .	33
7.3	Results: Simple Arithmetic . . . . .	34
7.3.1	Learning Addition . . . . .	34
7.3.2	Learning Maximum . . . . .	34
7.3.3	Learning Absolute Value . . . . .	34
7.4	Results: Array Operations . . . . .	35
7.4.1	Sum of Array . . . . .	35
7.5	Results: Sorting . . . . .	36
7.5.1	Bubble Sort . . . . .	36
7.6	Training Dynamics . . . . .	36
7.6.1	Temperature Annealing Effect . . . . .	36
7.6.2	Curriculum Learning Effect . . . . .	37
7.7	Discovered “Alien” Solutions . . . . .	37
7.8	Computational Costs . . . . .	38
7.9	Limitations of Current Experiments . . . . .	38
<b>8</b>	<b>Theoretical Framework</b>	<b>38</b>
8.1	The Grokking Hypothesis . . . . .	38
8.2	Code as Compression . . . . .	38
8.3	Hybrid Code-Weight Programs . . . . .	39
8.4	Efficiency Implications . . . . .	39
8.5	When Programs Beat Neural Networks . . . . .	40
<b>9</b>	<b>Research Program</b>	<b>40</b>
9.1	Phase 1: Classical Algorithms . . . . .	40
9.2	Phase 2: Using Pretrained Neural Networks . . . . .	41
9.3	Phase 3: End-to-End Learning . . . . .	41
9.4	Phase 4: Reinforcement Learning Tasks . . . . .	41
9.5	Phase 5: Vision and Feature Extraction . . . . .	42
9.6	Phase 6: GUI Agents . . . . .	42
9.7	Compute Budget . . . . .	42
9.8	Hardware Transfer . . . . .	43
<b>10</b>	<b>Discussion</b>	<b>43</b>
10.1	Key Uncertainties . . . . .	43
10.2	Limitations . . . . .	44
10.3	Broader Implications . . . . .	44
10.4	Discovering Alien Solutions . . . . .	45
10.5	Ethical Considerations . . . . .	45

<b>11 Conclusion</b>	<b>45</b>
11.1 Summary of Contributions . . . . .	45
11.2 The Central Hypothesis . . . . .	46
11.3 What Remains . . . . .	46
11.4 Closing Remarks . . . . .	46
<b>A Instruction Set Reference</b>	<b>49</b>
<b>B Memory Map</b>	<b>49</b>
<b>C GPU Command Format</b>	<b>49</b>
<b>D Training Hyperparameters</b>	<b>49</b>

# 1 Introduction

The pursuit of artificial general intelligence requires systems capable of arbitrary computational tasks with the flexibility of human cognition. Current approaches based on large language models (LLMs) and vision-language models (VLMs) have demonstrated remarkable capabilities across diverse domains [Brown et al., 2020, OpenAI, 2023, Team et al., 2023], yet exhibit fundamental architectural limitations that may preclude true generalization.

## 1.1 The Architectural Mismatch

Achieving general intelligence requires aligning computational architectures with the underlying latent symmetries of cognitive processes: symbolic representation, symbolic composition, and symbolic computation. Transformer autoregression and diffusion-based generation are not isomorphic with the latent dynamics these symmetries induce.

Consider the computational cost of a simple operation:

- Expressing “increment the loop counter” in natural language: 20-50 tokens
- The same operation in machine code: 1 instruction (`ADD R1, R1, 1`)
- The same operation in a compiled binary: 4 bytes

This 100-1000× representational overhead compounds across programs. An LLM-based agent performing a 50-step task generates thousands of tokens; a compiled program performing the same task executes in microseconds.

## 1.2 Evidence from Mechanistic Interpretability

Recent work in mechanistic interpretability provides evidence that this overhead may be unnecessary. Li et al. [2023] demonstrated that specific behaviors in large language models—such as truthfulness—can be controlled by intervening on low-dimensional subspaces of activations. Sharma et al. [2024] showed similar results for sycophancy. These findings suggest that complex behaviors operate in surprisingly low-rank subspaces, raising the question: why use an 8GB model when the relevant computation occurs in a <1000-dimensional circuit?

The grokking phenomenon [Power et al., 2022] provides further evidence. When transformers are trained far beyond training loss saturation, they suddenly experience dramatic validation improvements. Nanda et al. [2023] showed that grokked models implement interpretable algorithms—clean, sparse circuits rather than distributed representations. This phase transition represents the network discovering algorithmic structure that could, in principle, be expressed more compactly as code.

## 1.3 Contributions

This paper presents a differentiable tensor computer—a complete von Neumann architecture implemented as tensor operations in JAX [Bradbury et al., 2018]—that executes standard machine code while remaining fully differentiable. The contributions are:

1. **Complete mathematical formalization** of a differentiable computer architecture, including ALU, registers, cache, memory hierarchy, virtual memory, bus arbitration, and peripherals (Section 3).
2. **Demonstration of self-hosting compilation**: a C compiler written in tensor machine code that compiles C programs to tensor machine code, validated by executing `factorial(5)=120` (Section 4).

3. **Comprehensive training strategy** addressing the fundamental challenges of gradient-based program synthesis: temperature dynamics, curriculum learning, credit assignment, hierarchical decomposition, and hybrid gradient/search refinement (Section 5).
4. **Theoretical framework** connecting grokking phenomena to program synthesis, arguing that learned programs may be more efficient than dense weight matrices for structured computation (Section 8).
5. **Research program** for learning software agents via reinforcement learning on the tensor computer, with detailed compute budget and timeline (Section 9).

## 1.4 Paper Organization

Section 2 reviews related work on neural computers, grokking, program synthesis, and LLM agents. Section 3 provides complete mathematical specification of the tensor computer. Section 4 describes the self-hosting C compiler and validation. Section 5 details the training strategy. Section 8 presents the theoretical framework. Section 9 outlines the research program. Section 10 discusses limitations and implications. Section 11 concludes.

# 2 Background and Related Work

## 2.1 Neural Computers and Differentiable Programming

The Neural Turing Machine (NTM) [Graves et al., 2014] introduced differentiable read-write access to external memory using attention mechanisms. Given a memory matrix  $\mathbf{M} \in \mathbb{R}^{N \times W}$ , the NTM computes read weights via content-based addressing:

$$w_i^c = \frac{\exp(\beta \cdot \text{cosine}(\mathbf{k}, \mathbf{M}_i))}{\sum_j \exp(\beta \cdot \text{cosine}(\mathbf{k}, \mathbf{M}_j))} \quad (1)$$

where  $\mathbf{k}$  is a query key and  $\beta$  controls sharpness. The Differentiable Neural Computer (DNC) [Graves et al., 2016] extended this with temporal linking and dynamic memory allocation, enabling learning of graph algorithms and question answering.

The Neural GPU [Kaiser and Sutskever, 2015] demonstrated learning of algorithms like binary multiplication through a cellular automaton-like architecture with parallel updates across a 2D grid. The Neural Programmer-Interpreter [Reed and de Freitas, 2015] showed networks could learn to execute programs by composing learned subroutines, achieving generalization to longer sequences than seen during training.

**Key Distinction.** This work differs fundamentally from these approaches: rather than learning neural approximations of computational primitives, exact computational primitives are implemented as differentiable tensor operations. The ALU performs exact arithmetic; conditional branches are exact at low temperature. Programs learned on the tensor computer can be directly extracted and executed on standard hardware with  $O(1)$  host operations per virtual instruction.

## 2.2 Grokking and Algorithmic Phase Transitions

Power et al. [2022] discovered that transformers trained on algorithmic tasks (modular arithmetic, permutation composition) exhibit a striking phenomenon: long after training loss saturates near

zero, validation loss suddenly drops dramatically. This “grokking” transition can occur after 10-100× more training steps than needed for training convergence.

Nanda et al. [2023] used mechanistic interpretability to analyze grokked models trained on modular addition. They found that grokked models implement discrete Fourier transforms:

$$(a + b) \pmod p = \text{IDFT}(\text{DFT}(a) \odot \text{DFT}(b)) \quad (2)$$

The weights encode Fourier basis vectors; the computation is a clean algorithm, not a memorization table. The paper identified specific “progress measures”—intermediate metrics that predict grokking before it occurs in validation loss.

Liu et al. [2022] provided theoretical analysis showing grokking represents a transition from memorization to generalization, occurring when the loss landscape shifts to favor algorithmic solutions. They showed that representation learning dynamics cause networks to first fit training data via memorization, then slowly discover generalizing structure.

Varma et al. [2023] explained grokking through circuit efficiency: generalizing circuits have lower description complexity but require more optimization steps to find. The efficiency of the generalizing solution determines the delay before grokking.

**Implications.** These findings motivate the hypothesis that training neural networks to true generalization yields algorithmic solutions that could be more efficiently represented as explicit code. The grokking transition is the network discovering that the underlying task has algorithmic structure.

## 2.3 Program Synthesis

Classical program synthesis approaches include:

**Enumerative Search.** Udupa et al. [2013] and related work systematically enumerate programs in order of size, using input-output examples to prune the search space. This is complete but exponentially slow in program length.

**Constraint Solving.** Solar-Lezama et al. [2006] encode synthesis as SAT/SMT constraints, using CEGIS (counter-example guided inductive synthesis) to iteratively refine candidate programs. This scales to programs with complex specifications but requires formal specification.

**Type-Directed Synthesis.** Polikarpova et al. [2016] use rich type systems to prune the search space, generating programs that satisfy type signatures with refinement types encoding semantic constraints.

**Neural-Guided Search.** DeepCoder [Balog et al., 2017] uses neural networks to predict which DSL functions are likely to appear in a program solving given input-output examples, guiding enumerative search. This achieves significant speedups over unguided search.

**Code Generation from LLMs.** Codex [Chen et al., 2021] and AlphaCode [Li et al., 2022] generate programs as token sequences, achieving impressive results on competitive programming but inheriting the representational overhead of autoregressive generation.

**Library Learning.** DreamCoder [Ellis et al., 2021] alternates between neural-guided search and library compression, building reusable abstractions that accelerate future synthesis. This achieves human-like abstraction learning on various domains.

**This Work.** This paper takes a different approach: learning programs as executable machine code via direct gradient-based optimization of program parameters. Rather than generating programs as text, programs are represented as differentiable soft instructions that crystallize into discrete code at low temperature.

## 2.4 Reinforcement Learning and Superhuman Discovery

AlphaGo [Silver et al., 2016] and AlphaZero [Silver et al., 2017] demonstrated that reinforcement learning with self-play can discover superhuman strategies without human knowledge. AlphaGo’s “Move 37” in Game 2 against Lee Sedol—a move no human professional would play—illustrated that RL can find solutions outside the space of human-conceived strategies.

McGrath et al. [2022] analyzed AlphaZero’s chess knowledge acquisition, finding that it developed novel strategic concepts (like unconventional piece mobility patterns) that exploited the game’s structure in unexpected ways. These “alien” strategies were not just better versions of human strategies but qualitatively different approaches.

**Implications.** Similar dynamics may emerge when RL is applied to program synthesis on the tensor computer. The vast space of possible programs likely contains “alien” solutions—correct programs that no human would write, exploiting architectural features in unexpected ways. The differentiable tensor computer provides the substrate for RL to explore this space.

## 2.5 LLM-Based Agents and Their Limitations

Recent work applies large language models to agentic tasks:

**Web Navigation.** WebGPT [Nakano et al., 2021] uses GPT-3 to browse the web via text commands, achieving strong performance on question-answering benchmarks requiring web search.

**Computer Use.** Claude Computer Use [Anthropic, 2024] enables GUI control via screenshots and mouse/keyboard actions, allowing general-purpose computer interaction.

**Robotics.** RT-2 [Brohan et al., 2023] trains vision-language-action models for robotic manipulation, achieving generalization to novel objects and instructions. OpenVLA [Kim et al., 2024] provides an open-source alternative.

**Flow Matching.**  $\pi_0$  [Black et al., 2024] uses flow matching [Lipman et al., 2022] for continuous action spaces, improving efficiency for robotic control.

**Fundamental Limitations.** All these approaches inherit fundamental limitations:

- **Quadratic attention:** Cost scales as  $O(n^2)$  with context length
- **Representational overhead:** Simple operations require many tokens
- **Inference latency:** Seconds per action, not microseconds
- **Model size:** Gigabytes, not kilobytes
- **Cloud dependency:** Cannot run locally on mobile devices

**Low-Rank Structure.** Evidence from Li et al. [2023] and Sharma et al. [2024] suggests that specific behaviors operate in low-rank subspaces. If agentic behaviors truly operate in  $<1000$ -dimensional circuits, the full model capacity is unnecessary overhead that could be eliminated by learned programs.

### 3 Tensor Computer Architecture

This section provides complete mathematical specification of the differentiable tensor computer. All operations are implemented in JAX and remain differentiable with respect to program parameters.

#### 3.1 System Overview

Figure 1 shows the high-level system architecture.

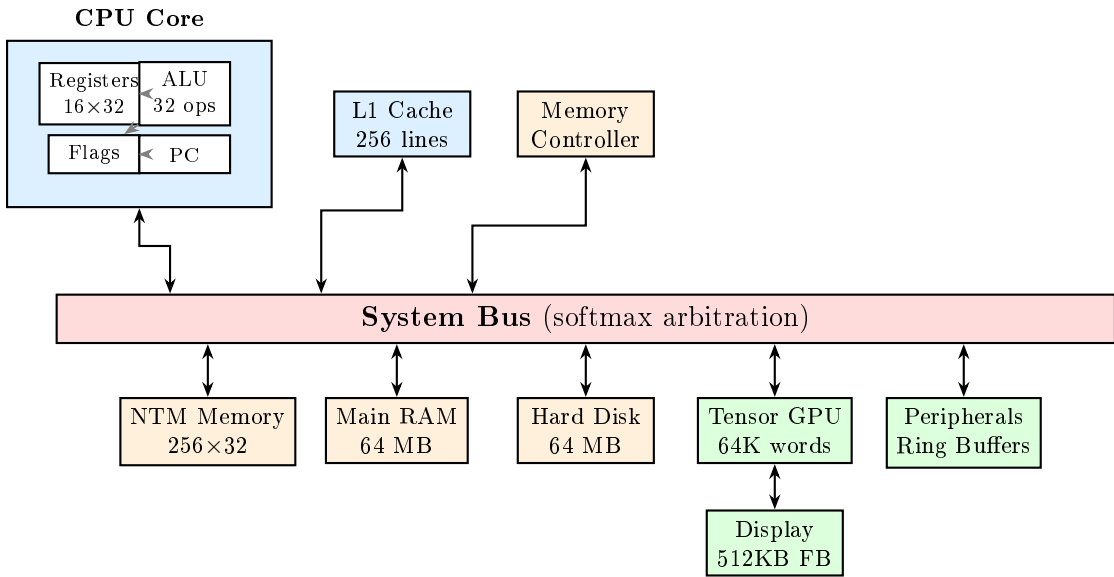


Figure 1: High-level tensor computer architecture. All components are implemented as differentiable tensor operations. The system bus uses softmax arbitration for multi-master access. At low temperature  $\tau \rightarrow 0$ , execution converges to discrete program semantics.

#### 3.2 Notation and Conventions

Let  $\tau > 0$  denote the temperature parameter controlling soft/hard behavior. Define the temperature-scaled softmax:

$$\text{softmax}_{\tau}(\mathbf{x})_i = \frac{\exp(x_i/\tau)}{\sum_j \exp(x_j/\tau)} \tag{3}$$

**Proposition 1** (Temperature Limits). *As  $\tau \rightarrow 0$ ,  $\text{softmax}_{\tau}(\mathbf{x}) \rightarrow \text{onehot}(\text{argmax}(\mathbf{x}))$ . As  $\tau \rightarrow \infty$ ,  $\text{softmax}_{\tau}(\mathbf{x}) \rightarrow 1/n$  (uniform distribution).*

*Proof.* For  $\tau \rightarrow 0$ : Let  $k^* = \text{argmax}_i x_i$ . Then  $x_{k^*}/\tau \rightarrow +\infty$  faster than  $x_i/\tau$  for  $i \neq k^*$ , so  $\exp(x_{k^*}/\tau)$  dominates the sum.

For  $\tau \rightarrow \infty$ : All  $x_i/\tau \rightarrow 0$ , so  $\exp(x_i/\tau) \rightarrow 1$  for all  $i$ , giving uniform weights. □

For binary signals, define the soft indicator:

$$\sigma_\tau(x) = \text{sigmoid}(x/\tau) = \frac{1}{1 + \exp(-x/\tau)} \quad (4)$$

This approaches a step function as  $\tau \rightarrow 0$ :

$$\lim_{\tau \rightarrow 0} \sigma_\tau(x) = \begin{cases} 1 & x > 0 \\ 0.5 & x = 0 \\ 0 & x < 0 \end{cases} \quad (5)$$

### 3.3 Machine State

The complete tensor computer state at time  $t$  is the tuple:

$$\mathcal{S}^{(t)} = \left( \mathbf{R}^{(t)}, \mathbf{PC}^{(t)}, \mathbf{F}^{(t)}, \mathbf{SP}^{(t)}, \mathbf{C}^{(t)}, \mathbf{M}^{(t)}, \mathbf{D}^{(t)}, \mathbf{G}^{(t)}, \mathbf{V}^{(t)}, \mathbf{B}^{(t)} \right) \quad (6)$$

Table 1 details each component.

Symbol	Dimensions	Description
$\mathbf{R}^{(t)}$	$\mathbb{R}^{16 \times 32}$	Register file (16 registers, 32 bits each)
$\mathbf{PC}^{(t)}$	$\mathbb{R}^{2^{16}}$	Program counter as soft address distribution
$\mathbf{F}^{(t)}$	$\mathbb{R}^4$	Flags: zero, negative, carry, overflow
$\mathbf{SP}^{(t)}$	$\mathbb{R}^{4096}$	Stack pointer as soft distribution
$\mathbf{C}^{(t)}$	$\mathbb{R}^{256 \times 8 \times 32}$	Cache (256 lines, 8 words/line)
$\mathbf{M}^{(t)}$	$\mathbb{R}^{256 \times 32}$	NTM-style main memory
$\mathbf{D}^{(t)}$	$\mathbb{R}^{2^{24} \times 32}$	Hard disk (64MB)
$\mathbf{G}^{(t)}$	$\mathbb{R}^{2^{16} \times 32}$	GPU memory
$\mathbf{V}^{(t)}$	$\mathbb{R}^{480 \times 640 \times 3}$	Display framebuffer (RGB)
$\mathbf{B}^{(t)}$	$\mathbb{R}^{N_p \times 256}$	Peripheral ring buffers

Table 1: Tensor computer state components. Each component is a tensor that evolves during program execution.

### 3.4 Instruction Encoding

Each instruction is encoded as a tuple:

$$\mathbf{I} = (\mathbf{op}, \mathbf{r}_d, \mathbf{r}_{s1}, \mathbf{r}_{s2}, \mathbf{imm}) \quad (7)$$

where:

- $\mathbf{op} \in \mathbb{R}^{32}$ : Opcode as soft distribution over 32 operations
- $\mathbf{r}_d, \mathbf{r}_{s1}, \mathbf{r}_{s2} \in \mathbb{R}^{16}$ : Register selectors as soft distributions
- $\mathbf{imm} \in \mathbb{R}^{32}$ : Immediate value (32-bit)

At low temperature, these distributions concentrate on single indices, recovering discrete instructions. The total instruction size is  $32 + 16 + 16 + 16 + 32 = 112$  parameters per instruction.

### 3.5 Instruction Fetch

The program counter  $\mathbf{PC}^{(t)}$  is a probability distribution over instruction addresses. Instruction fetch reads from instruction memory  $\mathbf{I}_{mem} \in \mathbb{R}^{A \times I}$  where  $A = 2^{16}$  is address space size and  $I = 112$  is instruction encoding size:

$$\mathbf{I}^{(t)} = \sum_{a=0}^{A-1} \mathbf{PC}_a^{(t)} \cdot \mathbf{I}_{mem}[a] = (\mathbf{PC}^{(t)})^\top \mathbf{I}_{mem} \quad (8)$$

This soft fetch enables gradient flow through the program counter. At  $\tau \rightarrow 0$  where  $\mathbf{PC}$  approaches a one-hot vector, only the addressed instruction is fetched.

**Remark 1** (Computational Cost). *Soft instruction fetch requires  $O(A \cdot I) = O(2^{16} \cdot 112) \approx 7.3$  million operations. This is the dominant cost per instruction and motivates future work on hierarchical address spaces.*

### 3.6 Register File Operations

#### 3.6.1 Soft Register Read

To read register selected by distribution  $\mathbf{r}_s \in \mathbb{R}^{16}$ :

$$\text{RegRead}(\mathbf{R}, \mathbf{r}_s) = \sum_{i=0}^{15} (\mathbf{r}_s)_i \cdot \mathbf{R}[i] = \mathbf{r}_s^\top \mathbf{R} \quad (9)$$

This returns a weighted combination of all registers. At low temperature where  $\mathbf{r}_s \approx \text{onehot}(k)$ , this returns  $\mathbf{R}[k]$ .

#### 3.6.2 Soft Register Write

To write value  $\mathbf{v} \in \mathbb{R}^{32}$  to register selected by  $\mathbf{r}_d \in \mathbb{R}^{16}$  with write enable  $w \in [0, 1]$ :

$$\mathbf{R}' = \mathbf{R} + w \cdot \mathbf{r}_d \otimes (\mathbf{v} - \text{RegRead}(\mathbf{R}, \mathbf{r}_d)) \quad (10)$$

where  $\otimes$  denotes outer product. Expanding:

$$\mathbf{R}'[i] = \mathbf{R}[i] + w \cdot (\mathbf{r}_d)_i \cdot (\mathbf{v} - \mathbf{R}[i]) \quad (11)$$

Each register is updated proportionally to its selection weight. At low temperature where  $\mathbf{r}_d \approx \text{onehot}(k)$ :

$$\mathbf{R}'[i] = \begin{cases} \mathbf{R}[i] + w \cdot (\mathbf{v} - \mathbf{R}[i]) = (1 - w)\mathbf{R}[i] + w\mathbf{v} & i = k \\ \mathbf{R}[i] & i \neq k \end{cases} \quad (12)$$

When  $w = 1$ , this is exact register write:  $\mathbf{R}'[k] = \mathbf{v}$ .

### 3.7 Arithmetic Logic Unit

The ALU is the computational heart of the tensor computer. Rather than conditionally executing one operation, all operations are computed in parallel and combined via softmax weighting.

### 3.7.1 Operation Definitions

Define operations  $f_k : \mathbb{R}^{32} \times \mathbb{R}^{32} \rightarrow \mathbb{R}^{32}$  for  $k \in \{0, \dots, 31\}$ :

$$f_0(\mathbf{a}, \mathbf{b}) = \mathbf{a} + \mathbf{b} \quad (\text{ADD}) \quad (13)$$

$$f_1(\mathbf{a}, \mathbf{b}) = \mathbf{a} - \mathbf{b} \quad (\text{SUB}) \quad (14)$$

$$f_2(\mathbf{a}, \mathbf{b}) = \mathbf{a} \odot \mathbf{b} \quad (\text{MUL, element-wise}) \quad (15)$$

$$f_3(\mathbf{a}, \mathbf{b}) = \mathbf{a} \oslash (\mathbf{b} + \epsilon) \quad (\text{DIV, stabilized}) \quad (16)$$

$$f_4(\mathbf{a}, \mathbf{b}) = \mathbf{a} \wedge \mathbf{b} \quad (\text{AND}) \quad (17)$$

$$f_5(\mathbf{a}, \mathbf{b}) = \mathbf{a} \vee \mathbf{b} \quad (\text{OR}) \quad (18)$$

$$f_6(\mathbf{a}, \mathbf{b}) = \mathbf{a} \oplus \mathbf{b} \quad (\text{XOR}) \quad (19)$$

$$f_7(\mathbf{a}, \mathbf{b}) = \neg \mathbf{a} \quad (\text{NOT}) \quad (20)$$

$$f_8(\mathbf{a}, \mathbf{b}) = \text{ShiftLeft}(\mathbf{a}, \mathbf{b}) \quad (\text{SHL}) \quad (21)$$

$$f_9(\mathbf{a}, \mathbf{b}) = \text{ShiftRight}(\mathbf{a}, \mathbf{b}) \quad (\text{SHR}) \quad (22)$$

$$f_{10}(\mathbf{a}, \mathbf{b}) = \text{ShiftRightArith}(\mathbf{a}, \mathbf{b}) \quad (\text{SRA}) \quad (23)$$

$$f_{11}(\mathbf{a}, \mathbf{b}) = \mathbf{b} \quad (\text{MOV}) \quad (24)$$

$$f_{12}(\mathbf{a}, \mathbf{b}) = \text{imm} \quad (\text{LOAD\_IMM}) \quad (25)$$

$$f_{13}(\mathbf{a}, \mathbf{b}) = \mathbf{a} - \mathbf{b} \text{ (flags only)} \quad (\text{CMP}) \quad (26)$$

$$f_{14}(\mathbf{a}, \mathbf{b}) = \text{select}(F_z, \mathbf{a}, \mathbf{b}) \quad (\text{CMOV\_Z}) \quad (27)$$

$$f_{15}(\mathbf{a}, \mathbf{b}) = \text{select}(F_n, \mathbf{a}, \mathbf{b}) \quad (\text{CMOV\_N}) \quad (28)$$

Operations 16-31 include memory operations, control flow, GPU commands, and I/O (detailed in subsequent sections).

### 3.7.2 Differentiable Bitwise Operations

For bitwise operations on continuous approximations, interpret each element  $a_i, b_i \in [0, 1]$  as a soft bit. The continuous versions are:

$$(\mathbf{a} \wedge \mathbf{b})_i = a_i \cdot b_i \quad (29)$$

$$(\mathbf{a} \vee \mathbf{b})_i = a_i + b_i - a_i \cdot b_i \quad (30)$$

$$(\mathbf{a} \oplus \mathbf{b})_i = a_i + b_i - 2 \cdot a_i \cdot b_i \quad (31)$$

$$(\neg \mathbf{a})_i = 1 - a_i \quad (32)$$

**Proposition 2** (Bitwise Correctness). *When inputs are in  $\{0, 1\}^{32}$ , these continuous formulas reduce to standard bitwise operations.*

*Proof.* For AND:  $0 \cdot 0 = 0, 0 \cdot 1 = 0, 1 \cdot 0 = 0, 1 \cdot 1 = 1$ . Matches boolean AND.

For OR:  $0 + 0 - 0 = 0, 0 + 1 - 0 = 1, 1 + 0 - 0 = 1, 1 + 1 - 1 = 1$ . Matches boolean OR.

For XOR:  $0 + 0 - 0 = 0, 0 + 1 - 0 = 1, 1 + 0 - 0 = 1, 1 + 1 - 2 = 0$ . Matches boolean XOR.  $\square$

### 3.7.3 Differentiable Shifts

Shift operations require special handling. For shift-left by amount  $s$ :

$$\text{ShiftLeft}(\mathbf{a}, s) = \mathbf{P}_L(s) \cdot \mathbf{a} \quad (33)$$

where  $\mathbf{P}_L(s) \in \mathbb{R}^{32 \times 32}$  is a permutation matrix. For integer  $s$ :

$$(\mathbf{P}_L(s))_{ij} = \begin{cases} 1 & j = i - s \text{ and } j \geq 0 \\ 0 & \text{otherwise} \end{cases} \quad (34)$$

This shifts bits left by  $s$  positions with zero fill from the right.

For soft (non-integer) shift amounts from register  $\mathbf{b}$ , first extract the shift amount as the low 5 bits interpreted as a number:

$$s = \sum_{i=0}^4 2^i \cdot b_i \quad (35)$$

Then interpolate between integer shifts:

$$\mathbf{P}_L(s) = \sum_{k=0}^{31} w_k(s) \cdot \mathbf{P}_L(k) \quad (36)$$

where  $w_k(s) = \text{softmax}_\tau(-|s - k|)_k$  concentrates weight on the nearest integer shift.

### 3.7.4 ALU Output Computation

Figure 2 illustrates the softmax-gated ALU design.

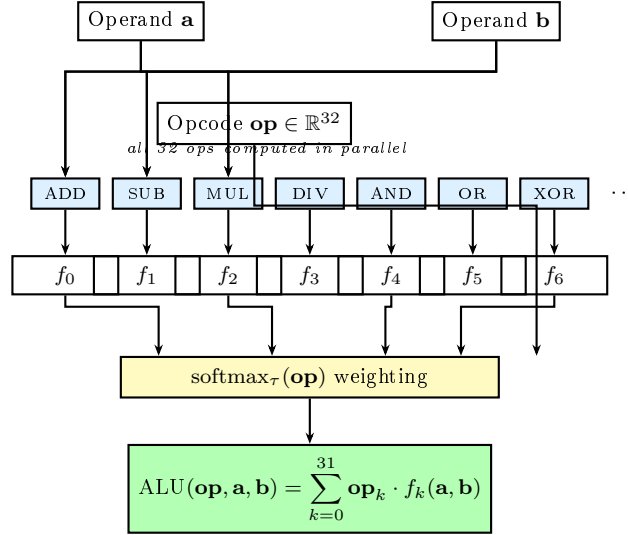


Figure 2: Softmax-gated ALU architecture. All 32 operations are computed in parallel; the opcode distribution  $\mathbf{op}$  selects the weighted output. At low temperature  $\tau \rightarrow 0$ , this converges to discrete operation selection.

Given opcode distribution  $\mathbf{op}$  and operands  $\mathbf{a}, \mathbf{b}$ :

$$\text{ALU}(\mathbf{op}, \mathbf{a}, \mathbf{b}) = \sum_{k=0}^{31} \mathbf{op}_k \cdot f_k(\mathbf{a}, \mathbf{b}) \quad (37)$$

All 32 operations are computed in parallel; the opcode distribution weights the outputs. At low temperature where  $\mathbf{op} \approx \text{onehot}(k^*)$ , this returns  $f_{k^*}(\mathbf{a}, \mathbf{b})$ .

**Remark 2** (Gradient Flow). *At high temperature, gradients flow through all operations. If operation  $k$  would improve the loss but has low weight  $\mathbf{op}_k$ , the gradient  $\partial \mathcal{L} / \partial \mathbf{op}_k$  encourages increasing this weight. This is the mechanism by which the correct opcode is learned.*

### 3.7.5 Flag Computation

Flags are computed from the ALU result  $\mathbf{y}$  and operation:

$$F_{\text{zero}} = \sigma_{\tau}(-\|\mathbf{y}\|_1/\epsilon) \quad (38)$$

$$F_{\text{neg}} = \sigma_{\tau}(-y_{31}) \quad (\text{sign bit}) \quad (39)$$

$$F_{\text{carry}} = \sigma_{\tau}(\text{CarryOut}(\mathbf{a}, \mathbf{b}, \mathbf{op})) \quad (40)$$

$$F_{\text{overflow}} = \sigma_{\tau}(\text{Overflow}(\mathbf{a}, \mathbf{b}, \mathbf{y}, \mathbf{op})) \quad (41)$$

For addition, the carry and overflow computations are:

$$\text{CarryOut}_{\text{add}} = \mathbf{1}[a_{31} + b_{31} + c_{30} \geq 2] \quad (42)$$

$$\text{Overflow}_{\text{add}} = (a_{31} = b_{31}) \wedge (a_{31} \neq y_{31}) \quad (43)$$

where  $c_{30}$  is the carry into bit 31. These are made differentiable via soft comparisons and soft equality.

## 3.8 Memory Hierarchy

The tensor computer implements a realistic memory hierarchy. Figure 3 illustrates the structure with differentiability annotations.

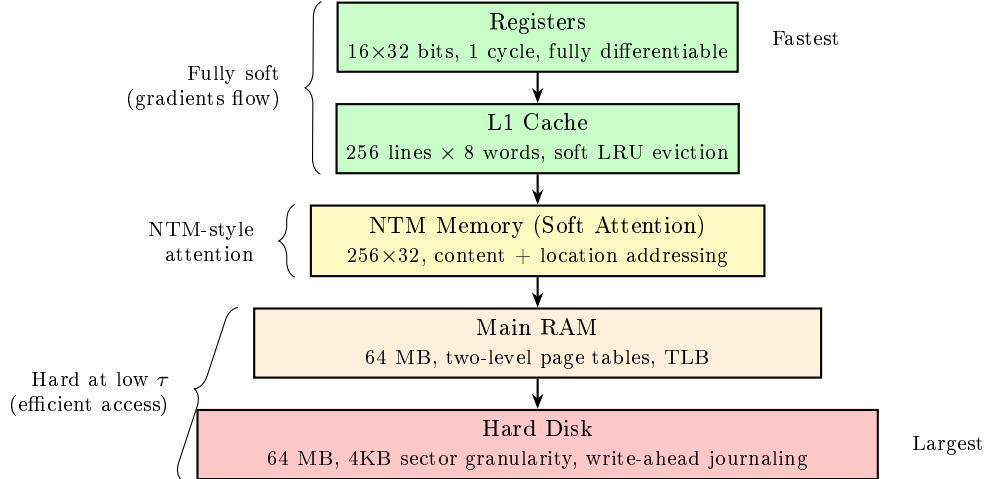


Figure 3: Memory hierarchy with differentiability annotations. Upper levels use fully soft addressing; lower levels approach hard addressing at low temperature to encourage efficient access patterns that transfer to real hardware.

### 3.8.1 Cache Architecture

The cache  $\mathbf{C} \in \mathbb{R}^{N_c \times (T+L \cdot W+3)}$  stores  $N_c = 256$  lines, each containing:

- Tag bits  $\mathbf{tag} \in \mathbb{R}^T$  identifying the memory address ( $T = 20$  bits)
- Valid bit  $v \in [0, 1]$
- Dirty bit  $d \in [0, 1]$
- Data  $\mathbf{data} \in \mathbb{R}^{L \times W}$  ( $L = 8$  words per line,  $W = 32$  bits)
- LRU counter  $\ell \in \mathbb{R}$  for eviction policy

The cache line structure is:

$$\mathbf{C}[i] = (\mathbf{tag}_i, v_i, d_i, \mathbf{data}_i, \ell_i) \quad (44)$$

### 3.8.2 Cache Lookup

Given address  $\mathbf{addr} \in \mathbb{R}^{32}$ , decompose into tag  $\mathbf{t}$ , index  $\mathbf{idx}$ , and offset  $\mathbf{o}$ :

$$\mathbf{t} = \mathbf{addr}[31 : 12] \quad (20 \text{ tag bits}) \quad (45)$$

$$\mathbf{idx} = \mathbf{addr}[11 : 5] \quad (7 \text{ index bits for 256 lines}) \quad (46)$$

$$\mathbf{o} = \mathbf{addr}[4 : 2] \quad (3 \text{ offset bits for 8 words}) \quad (47)$$

For a direct-mapped cache with soft addressing, compute hit scores for each line:

$$h_i = v_i \cdot \exp(-\|\mathbf{tag}_i - \mathbf{t}\|^2/\tau) \cdot \exp(-|i - \mathbf{idx}|^2/\tau) \quad (48)$$

The soft hit indicator and line selection:

$$\mathbf{hit} = \text{softmax}_{\tau}(\mathbf{h}), \quad p_{\text{hit}} = \sum_i h_i \quad (49)$$

At low temperature, this converges to standard direct-mapped cache behavior.

### 3.8.3 Cache Read

On cache hit, read from selected line with offset:

$$\text{CacheRead}(\mathbf{C}, \mathbf{addr}) = \sum_{i=0}^{N_c-1} \mathbf{hit}_i \cdot \sum_{j=0}^{L-1} \text{softmax}_{\tau}(-|j - \mathbf{o}|)_j \cdot \mathbf{C}_i.\mathbf{data}[j] \quad (50)$$

On miss (when  $p_{\text{hit}} < \theta$  for threshold  $\theta$ ), fetch from NTM memory.

### 3.8.4 Cache Write and Eviction

Write-back policy: writes update cache only; dirty bit tracks modifications.

$$d'_i = d_i + \mathbf{hit}_i \cdot w \cdot (1 - d_i) \quad (51)$$

On cache miss requiring new line allocation, select victim via soft LRU:

$$\mathbf{evict} = \text{softmax}_{\tau}(-\ell) \quad (52)$$

If evicted line is dirty, write back to main memory:

$$\mathbf{M}' = \mathbf{M} + \sum_i \mathbf{evict}_i \cdot d_i \cdot \Delta \mathbf{M}_i \quad (53)$$

where  $\Delta \mathbf{M}_i$  writes line  $i$ 's data to the address specified by its tag.

Update LRU counters on access:

$$\ell'_i = \begin{cases} 0 & \text{if } \mathbf{hit}_i > 0.5 \text{ (accessed)} \\ \ell_i + 1 & \text{otherwise (age)} \end{cases} \quad (54)$$

### 3.8.5 NTM-Style Main Memory

Main memory uses Neural Turing Machine addressing [Graves et al., 2014]. The memory  $\mathbf{M} \in \mathbb{R}^{N_m \times W}$  with  $N_m = 256$  rows supports both content-based and location-based addressing.

**Content-Based Addressing.** Given key  $\mathbf{k} \in \mathbb{R}^W$  and strength  $\beta > 0$ :

$$\mathbf{w}_i^c = \frac{\exp(\beta \cdot \text{cosine}(\mathbf{k}, \mathbf{M}[i]))}{\sum_j \exp(\beta \cdot \text{cosine}(\mathbf{k}, \mathbf{M}[j]))} \quad (55)$$

where

$$\text{cosine}(\mathbf{k}, \mathbf{m}) = \frac{\mathbf{k} \cdot \mathbf{m}}{\|\mathbf{k}\| \|\mathbf{m}\| + \epsilon} \quad (56)$$

This allows content-addressable memory: retrieve by value rather than by address.

**Location-Based Addressing.** Given previous weights  $\mathbf{w}^{(t-1)}$ , interpolation gate  $g \in [0, 1]$ , shift distribution  $\mathbf{s} \in \mathbb{R}^{2S+1}$  (for shifts  $-S, \dots, S$ ), and sharpening  $\gamma \geq 1$ :

$$\mathbf{w}^g = g \cdot \mathbf{w}^c + (1 - g) \cdot \mathbf{w}^{(t-1)} \quad (\text{interpolate}) \quad (57)$$

$$\tilde{\mathbf{w}}_i = \sum_{j=-S}^S \mathbf{w}_{(i-j) \bmod N_m}^g \cdot \mathbf{s}_{j+S} \quad (\text{convolutional shift}) \quad (58)$$

$$\mathbf{w}_i = \frac{\tilde{\mathbf{w}}_i^\gamma}{\sum_j \tilde{\mathbf{w}}_j^\gamma} \quad (\text{sharpen}) \quad (59)$$

The shift enables relative addressing (“next element”, “previous element”), while sharpening makes the address distribution more peaked.

**Memory Read.**

$$\text{MemRead}(\mathbf{M}, \mathbf{w}) = \sum_i \mathbf{w}_i \cdot \mathbf{M}[i] = \mathbf{w}^\top \mathbf{M} \quad (60)$$

**Memory Write.** Using erase vector  $\mathbf{e} \in [0, 1]^W$  and add vector  $\mathbf{a} \in \mathbb{R}^W$ :

$$\mathbf{M}'[i] = \mathbf{M}[i] \odot (\mathbf{1} - \mathbf{w}_i \cdot \mathbf{e}) + \mathbf{w}_i \cdot \mathbf{a} \quad (61)$$

This first erases (by element-wise multiplying with  $\mathbf{1} - \mathbf{w}_i \mathbf{e}$ ) then adds (by adding  $\mathbf{w}_i \mathbf{a}$ ).

### 3.8.6 Hard Disk with 4KB Sector Loading

The hard disk  $\mathbf{D} \in \mathbb{R}^{N_d \times W}$  with  $N_d = 2^{24}$  words (64MB) represents persistent storage with 4KB sector granularity. Sector size is  $S_{\text{sector}} = 1024$  words (4KB with 32-bit words).

**Sector Address Computation.** Given byte address  $a$ :

$$\text{sector\_idx} = \lfloor a / S_{\text{sector}} \rfloor, \quad \text{sector\_off} = a \bmod S_{\text{sector}} \quad (62)$$

The disk has  $N_d / S_{\text{sector}} = 2^{14} = 16384$  sectors.

**Differentiable Sector Load.** Loading entire sectors (not individual words) is intentional—this encourages learned programs to develop efficient access patterns:

$$\text{LoadSector}(\mathbf{D}, s) = \mathbf{D}[s \cdot S_{\text{sector}} : (s + 1) \cdot S_{\text{sector}}] \quad (63)$$

For soft sector index  $\mathbf{s} \in \mathbb{R}^{N_d/S_{\text{sector}}}$ :

$$\text{SoftLoadSector}(\mathbf{D}, \mathbf{s}) = \sum_{i=0}^{N_d/S_{\text{sector}}-1} \mathbf{s}_i \cdot \text{LoadSector}(\mathbf{D}, i) \quad (64)$$

This is  $O(N_d)$  cost, motivating efficient access patterns that transfer to real hardware.

**Write-Back and Journaling.** Disk writes use a write-ahead log (journal) for crash consistency:

1. Write operation logged to journal:  $\mathbf{J} \leftarrow \mathbf{J} \cup \{(\text{addr}, \text{data})\}$
2. Journal flushed to disk
3. Actual write performed:  $\mathbf{D}' = \text{ApplyJournal}(\mathbf{D}, \mathbf{J})$
4. Journal cleared

This ensures atomicity of disk operations.

### 3.8.7 Virtual Memory and Page Tables

Virtual-to-physical address translation uses a two-level page table, enabling programs to use a flat address space while the system manages physical memory allocation.

**Address Decomposition.** Virtual address  $\mathbf{va}$  (32 bits) decomposes as:

$$\mathbf{va} = \underbrace{[\text{L1\_idx}]}_{10 \text{ bits}} : \underbrace{[\text{L2\_idx}]}_{10 \text{ bits}} : \underbrace{[\text{offset}]}_{12 \text{ bits}} \quad (65)$$

This gives:

- $2^{10} = 1024$  entries in the L1 page table
- $2^{10} = 1024$  entries per L2 page table
- $2^{12} = 4096$  byte pages (matching sector size)
- Total addressable:  $2^{32} = 4\text{GB}$  virtual address space

**Page Table Entry Format.** Each PTE contains:

$$\text{PTE} = \underbrace{[\text{frame}]}_{20 \text{ bits}} : \underbrace{[\text{flags}]}_{12 \text{ bits}} \quad (66)$$

Flags include: present, writable, user-accessible, dirty, accessed.

**Page Table Walk.**

$$\mathbf{pte}_1 = \text{MemRead}(\mathbf{PT}_1, \text{softmax}(\text{L1\_idx})) \quad (67)$$

$$\mathbf{pte}_2 = \text{MemRead}(\mathbf{PT}_2[\mathbf{pte}_1.\text{ptr}], \text{softmax}(\text{L2\_idx})) \quad (68)$$

$$\mathbf{pa} = [\mathbf{pte}_2.\text{frame} : \text{offset}] \quad (69)$$

**TLB (Translation Lookaside Buffer).** Cache recent translations to avoid page table walks:

$$\text{TLB} : \mathbf{va}[31 : 12] \mapsto \mathbf{pa}[31 : 12] \quad (70)$$

TLB uses the same soft-lookup mechanism as the data cache, with 64 entries and LRU replacement.

### 3.9 Control Flow

#### 3.9.1 Sequential Execution

Default PC update advances to next instruction:

$$\mathbf{PC}^{(t+1)} = \text{Shift}(\mathbf{PC}^{(t)}, +1) \quad (71)$$

where  $\text{Shift}(\mathbf{p}, k)$  cyclically rotates distribution  $\mathbf{p}$  by  $k$  positions:

$$\text{Shift}(\mathbf{p}, k)_i = \mathbf{P}_{(i-k) \bmod A} \quad (72)$$

#### 3.9.2 Conditional Branches

For conditional branch with condition  $c \in [0, 1]$  and target distribution **target**:

$$\mathbf{PC}^{(t+1)} = c \cdot \mathbf{target} + (1 - c) \cdot \text{Shift}(\mathbf{PC}^{(t)}, +1) \quad (73)$$

The condition is computed from flags:

$$c_{\text{JEQ}} = F_{\text{zero}} \quad (\text{jump if equal/zero}) \quad (74)$$

$$c_{\text{JNE}} = 1 - F_{\text{zero}} \quad (\text{jump if not equal}) \quad (75)$$

$$c_{\text{JLT}} = F_{\text{neg}} \cdot (1 - F_{\text{zero}}) \quad (\text{jump if less than}) \quad (76)$$

$$c_{\text{JGE}} = 1 - F_{\text{neg}} \quad (\text{jump if greater or equal}) \quad (77)$$

$$c_{\text{JLE}} = F_{\text{neg}} + F_{\text{zero}} - F_{\text{neg}} \cdot F_{\text{zero}} \quad (\text{jump if } \leq) \quad (78)$$

$$c_{\text{JGT}} = (1 - F_{\text{neg}}) \cdot (1 - F_{\text{zero}}) \quad (\text{jump if greater}) \quad (79)$$

At low temperature, flags are in  $\{0, 1\}$  and conditions are exact.

#### 3.9.3 Subroutine Calls

CALL pushes return address and jumps:

$$\mathbf{Stack}[\mathbf{SP}] \leftarrow \text{Shift}(\mathbf{PC}^{(t)}, +1) \quad (80)$$

$$\mathbf{SP} \leftarrow \text{Shift}(\mathbf{SP}, +1) \quad (81)$$

$$\mathbf{PC}^{(t+1)} \leftarrow \mathbf{target} \quad (82)$$

RET pops and restores:

$$\mathbf{SP} \leftarrow \text{Shift}(\mathbf{SP}, -1) \quad (83)$$

$$\mathbf{PC}^{(t+1)} \leftarrow \mathbf{Stack}[\mathbf{SP}] \quad (84)$$

The stack is a circular buffer of size  $S = 4096$ , sufficient for deep recursion.

### 3.10 Bus Architecture

The system bus connects CPU, memory controller, GPU, and peripherals.

#### 3.10.1 Bus Arbitration

Multiple masters may request bus access. Given request vector  $\mathbf{req} \in [0, 1]^{N_{\text{masters}}}$  and priority weights  $\mathbf{pri} \in \mathbb{R}^{N_{\text{masters}}}$ :

$$\mathbf{grant} = \underset{\tau}{\text{softmax}}(\mathbf{req} \odot \mathbf{pri}) \tag{85}$$

where  $\odot$  is element-wise multiplication. The granted master’s transaction proceeds; others see their requests held for next cycle.

Masters include: CPU (priority 1.0), DMA controller (0.8), GPU (0.6).

#### 3.10.2 Memory-Mapped I/O

Peripheral registers are memory-mapped. The address decoder routes accesses:

$$\text{target}(\mathbf{addr}) = \begin{cases} \text{RAM} & \mathbf{addr} < 0\mathbf{x}10000000 \\ \text{GPU} & 0\mathbf{x}10000000 \leq \mathbf{addr} < 0\mathbf{x}20000000 \\ \text{Peripherals} & 0\mathbf{x}20000000 \leq \mathbf{addr} < 0\mathbf{x}30000000 \\ \text{Framebuffer} & 0\mathbf{x}30000000 \leq \mathbf{addr} < 0\mathbf{x}40000000 \\ \text{Disk} & \mathbf{addr} \geq 0\mathbf{x}80000000 \end{cases} \tag{86}$$

Soft decoding uses sigmoid gating:

$$w_{\text{GPU}} = \sigma_{\tau}(\mathbf{addr} - 0\mathbf{x}10000000) \cdot \sigma_{\tau}(0\mathbf{x}20000000 - \mathbf{addr}) \tag{87}$$

### 3.11 GPU Interface

The tensor GPU accelerates parallel operations, essential for neural network inference within learned programs.

#### 3.11.1 GPU Commands

Commands issued via memory-mapped command queue at  $0\mathbf{x}10000000$ :

- GPU\_LOAD(*src*, *dst*, *n*): DMA *n* words from RAM address *src* to GPU address *dst*
- GPU\_STORE(*src*, *dst*, *n*): DMA from GPU to RAM
- GPU\_MATMUL(*A*, *B*, *C*, *m*, *n*, *k*):  $\mathbf{C}_{m \times n} = \mathbf{A}_{m \times k} \mathbf{B}_{k \times n}$
- GPU\_CONV(*in*, *kernel*, *out*, *params*): 2D convolution with specified parameters
- GPU\_REDUCE(*in*, *out*, *n*, *op*): Parallel reduction (sum, max, min)
- GPU\_ELEMENTWISE(*a*, *b*, *out*, *n*, *op*): Element-wise operation

#### 3.11.2 GPU Execution

GPU operations execute asynchronously. The CPU polls or waits:

$$\text{GPU\_SYNC} : \text{stall CPU until } \mathbf{G}.\text{busy} = 0 \tag{88}$$

GPU operations are natively differentiable as standard tensor operations, enabling end-to-end gradient flow through neural network inference.

### 3.12 Peripheral Ring Buffers

Each peripheral (keyboard, mouse, network, etc.) has input and output ring buffers for asynchronous I/O.

#### 3.12.1 Ring Buffer State

Buffer state: ( $\mathbf{data} \in \mathbb{R}^{B_s \times W}$ ,  $\mathbf{head} \in \mathbb{R}^{B_s}$ ,  $\mathbf{tail} \in \mathbb{R}^{B_s}$ )

where  $\mathbf{head}$  and  $\mathbf{tail}$  are soft position distributions over  $B_s = 256$  slots.

#### 3.12.2 Enqueue/Dequeue

$$\text{Enqueue}(\mathbf{v}) : \mathbf{data}' = \mathbf{data} + \mathbf{tail} \otimes (\mathbf{v} - \text{Read}(\mathbf{data}, \mathbf{tail})) \quad (89)$$

$$\mathbf{tail}' = \text{Shift}(\mathbf{tail}, +1) \quad (90)$$

$$\text{Dequeue}() : \mathbf{v} = \text{Read}(\mathbf{data}, \mathbf{head}) \quad (91)$$

$$\mathbf{head}' = \text{Shift}(\mathbf{head}, +1) \quad (92)$$

Buffer full/empty conditions:

$$\text{full} = \sigma_\tau \left( \mathbf{head}^\top \text{Shift}(\mathbf{tail}, +1) - 0.5 \right) \quad (93)$$

$$\text{empty} = \sigma_\tau \left( \mathbf{head}^\top \mathbf{tail} - 0.5 \right) \quad (94)$$

### 3.13 Display Buffer

The framebuffer  $\mathbf{V} \in \mathbb{R}^{H \times W_{\text{disp}} \times 3}$  with  $H = 480$ ,  $W_{\text{disp}} = 640$  stores RGB pixel values in  $[0, 1]$ :

$$\mathbf{V}[y, x, c] = \text{value} \quad \text{for } c \in \{R, G, B\} \quad (95)$$

Memory-mapped writes at address  $\mathbf{addr}$  relative to framebuffer base:

$$\text{offset} = \mathbf{addr} - 0\mathbf{x}30000000 \quad (96)$$

$$y = \lfloor \text{offset} / (W_{\text{disp}} \cdot 3) \rfloor \quad (97)$$

$$x = \lfloor (\text{offset} \bmod (W_{\text{disp}} \cdot 3)) / 3 \rfloor \quad (98)$$

$$c = \text{offset} \bmod 3 \quad (99)$$

This enables learned programs to produce visual output for GUI interaction tasks.

### 3.14 Computational Cost Analysis

Each tensor computer instruction requires the following operations:

Total: approximately **800K FLOPs per instruction** with current configuration. The dominant cost is soft instruction fetch; future work could use hierarchical addressing to reduce this.

## 4 Self-Hosting C Compiler

To validate the tensor computer’s capability for complex program execution, a self-hosting C compiler was implemented. This demonstrates that the architecture can execute programs of significant complexity, not just toy examples.

Component	Complexity	FLOPs (typical)
Instruction fetch	$O(A \cdot I)$	$\sim 7.3\text{M}$
Register read ( $\times 2$ )	$O(N_r \cdot W)$	$\sim 1\text{K}$
ALU (all ops)	$O(N_{op} \cdot W)$	$\sim 1\text{K}$
Register write	$O(N_r \cdot W)$	$\sim 0.5\text{K}$
PC update	$O(A)$	$\sim 65\text{K}$
Memory access (if needed)	$O(N_m \cdot W)$	$\sim 8\text{K}$

Table 2: Per-instruction computational costs. Instruction fetch dominates due to soft addressing over  $A = 2^{16}$  addresses.

### 4.1 Compiler Architecture

The compiler implements a standard pipeline for a subset of C:

1. **Lexer:** Tokenizes C source into keywords, identifiers, literals, operators, and punctuation
2. **Parser:** Recursive descent parser generating AST nodes for declarations, statements, and expressions
3. **Symbol Table:** Scoped tracking of variables, functions, and types
4. **Type Checker:** Validates type compatibility
5. **Code Generator:** AST traversal emitting tensor machine code

### 4.2 Bootstrap Process

Figure 4 illustrates the self-hosting compilation chain.

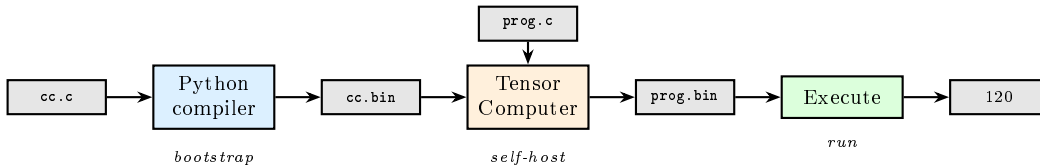


Figure 4: Self-hosting compilation chain:  $cc.c \xrightarrow{\text{Python}} cc.bin \xrightarrow{\text{load}} \text{tensor computer} \xrightarrow{\text{compiles}} prog.c \rightarrow prog.bin \rightarrow \text{executes} \rightarrow \text{result}$ .

### 4.3 Calling Convention

The compiler uses a standard ARM-like calling convention:

- **Arguments:** First 4 in R0-R3, remainder on stack
- **Return value:** R0
- **Callee-saved:** R4-R11
- **Caller-saved:** R0-R3, R12
- **Frame pointer:** R11 (FP)
- **Stack pointer:** R13 (SP)
- **Link register:** R14 (LR)

## 4.4 Validation: Recursive Factorial

The complete chain was validated with recursive factorial:

```
1 int factorial(int n) {
2     if (n <= 1) return 1;
3     return n * factorial(n - 1);
4 }
5
6 int main() {
7     return factorial(5);
8 }
```

**Result: 120.**

This validates:

- Correct execution of recursion (5 nested calls)
- Correct conditional branching (if statement)
- Correct stack frame management (push/pop across calls)
- Correct arithmetic (multiplication, subtraction, comparison)
- Correct function call/return protocol (CALL/RET, LR preservation)
- Correct callee-saved register handling (FP preservation)

## 4.5 Compiler Complexity

The self-hosting compiler consists of approximately 3000 lines of C, compiling to approximately 15,000 tensor instructions. When executing to compile a user program, it performs roughly 500,000 instruction executions. This demonstrates that the architecture can handle sustained computation at scale.

# 5 Training Strategy

Learning programs via gradient descent and reinforcement learning presents fundamental challenges distinct from standard neural network training. This section details a comprehensive strategy addressing each challenge.

## 5.1 The Core Challenge: Discrete Structure via Continuous Optimization

Programs are combinatorially structured: a single bit flip in an opcode completely changes semantics. The loss landscape for program synthesis contains:

- **Vast plateaus:** Many programs that equally fail (return wrong answer, crash, loop forever)
- **Narrow ridges:** Programs that partially work (correct on some inputs)
- **Isolated peaks:** Correct programs
- **Deceptive basins:** Programs that overfit to training distribution

This is unlike the smooth landscapes in typical deep learning where small parameter changes yield small loss changes. For programs, small changes often yield discontinuous jumps in behavior.

## 5.2 Temperature Dynamics

Temperature  $\tau$  controls the exploration-exploitation tradeoff in program space. At high temperature, the soft program is a differentiable mixture of all operations; at low temperature, it converges to a discrete program.

Figure 5 illustrates the annealing schedule.

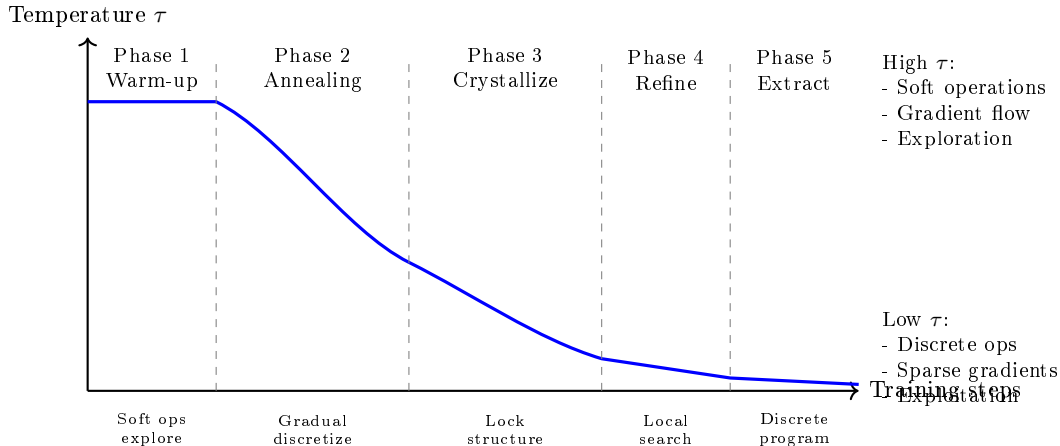


Figure 5: Temperature annealing schedule. Training proceeds through five phases with decreasing temperature, transitioning from exploration in soft program space to exploitation of discrete program structure.

### 5.2.1 Phase 1: Warm-up (High Temperature)

At high temperature ( $\tau \gg 1$ ):

- All operations contribute to output via soft weighting
- Gradients flow through all computational paths
- The “program” is a differentiable mixture of all possible programs
- Goal: find a region of parameter space with low expected loss

The soft program at this stage does not correspond to any valid discrete program, but optimization can identify promising directions in parameter space.

### 5.2.2 Phase 2: Annealing (Decreasing Temperature)

Gradually reduce  $\tau$  following the schedule:

$$\tau(t) = \tau_{\max} \cdot \exp(-\lambda \cdot t/T) \quad (100)$$

where  $T$  is total training steps and  $\lambda$  controls decay rate (typically  $\lambda = 5$ ).

During annealing:

- Soft operations become increasingly peaked around dominant choices
- The program resembles discrete instructions more closely
- Credit assignment becomes clearer as fewer operations contribute
- Continue gradient-based optimization with decreasing learning rate

### 5.2.3 Phase 3: Crystallization (Low Temperature)

At low temperature ( $\tau \ll 1$ ):

- Operations are nearly one-hot (effectively discrete)
- Gradients become sparse—only selected operation receives signal
- Use straight-through estimator [Bengio et al., 2013] for gradient flow:

$$\frac{\partial \mathcal{L}}{\partial \mathbf{op}} \approx \frac{\partial \mathcal{L}}{\partial \text{output}} \cdot \mathbf{1} \quad (101)$$

- Program structure is largely fixed; fine-tune remaining soft parameters

### 5.2.4 Phase 4: Refinement (Discrete Local Search)

Once temperature is near zero:

1. Extract discrete program via argmax of each soft selector
2. Switch from gradient descent to discrete local search
3. For each instruction position, try all alternative opcodes
4. Accept changes that improve validation performance (hill-climbing)
5. Optionally use beam search to explore multiple candidates

This phase can discover improvements that gradient descent missed—local optimizations in discrete program space that weren’t reachable via continuous paths.

### 5.2.5 Phase 5: Extraction and Verification

Final phase:

1. Extract fully discrete program as standard machine code
2. Verify correctness on held-out test cases (not used during training)
3. Optionally compile to standard VM or native code for deployment
4. Add successful programs/subroutines to library for future use

## 5.3 Credit Assignment

For a program with  $T$  instructions, gradients from final loss must flow through  $T$  sequential operations. Even with full differentiability, signals are diluted (vanishing gradients) and potentially misleading (operations that didn’t cause the error receive gradient).

### 5.3.1 Shaped Rewards

Replace sparse reward (correct/incorrect output) with dense shaped rewards:

**Output Distance.**

$$r_{\text{dist}} = -\|\mathbf{y}_{\text{pred}} - \mathbf{y}_{\text{target}}\|_p \tag{102}$$

For specific tasks:

- Sorting: negative inversion count (pairs out of order)
- Arithmetic: negative relative error  $|y - \hat{y}|/(|\hat{y}| + \epsilon)$
- String operations: negative edit distance

**Partial Correctness.** For outputs with multiple components (e.g., sorting an array):

$$r_{\text{partial}} = \frac{1}{n} \sum_{i=1}^n \mathbf{1}[y_i = \hat{y}_i] \tag{103}$$

**Execution Efficiency.** Small bonus for programs that halt sooner (Occam’s razor for compute):

$$r_{\text{eff}} = \alpha \cdot (T_{\text{max}} - T_{\text{halt}})/T_{\text{max}} \tag{104}$$

**Memory Safety.** Penalty for invalid operations:

$$r_{\text{safe}} = -\beta \cdot \left( \sum_t \mathbf{1}[\text{OOB}_t] + \gamma \sum_t \mathbf{1}[\text{div0}_t] \right) \tag{105}$$

### 5.3.2 Auxiliary Losses

Add supervision at intermediate execution points:

$$\mathcal{L}_{\text{total}} = \mathcal{L}_{\text{output}} + \sum_{t \in \mathcal{T}} \alpha_t \mathcal{L}_{\text{aux}}^{(t)} \quad (106)$$

where  $\mathcal{T}$  is a set of checkpoints (e.g., every 10 instructions) and  $\mathcal{L}_{\text{aux}}^{(t)}$  measures intermediate state quality.

Example for sorting: at checkpoint  $t$ , compute what fraction of the array is in its final sorted position.

### 5.3.3 Learned Value Function

Train a critic network  $V_\phi$  to estimate expected future reward from intermediate states:

$$V_\phi(\mathcal{S}^{(t)}) \approx \mathbb{E} \left[ \sum_{t'=t}^T \gamma^{t'-t} r_{t'} \mid \mathcal{S}^{(t)} \right] \quad (107)$$

This provides dense reward signal via temporal difference (TD) learning:

$$\delta_t = r_t + \gamma V_\phi(\mathcal{S}^{(t+1)}) - V_\phi(\mathcal{S}^{(t)}) \quad (108)$$

The TD error  $\delta_t$  indicates whether the current state is better or worse than expected, providing per-step feedback.

The critic architecture takes machine state  $\mathcal{S}^{(t)}$  as input:

- Embed registers, PC, flags via learned embeddings
- Attention over memory contents
- MLP to produce scalar value estimate

## 5.4 Curriculum Learning

Start with simple tasks and gradually increase complexity along multiple dimensions.

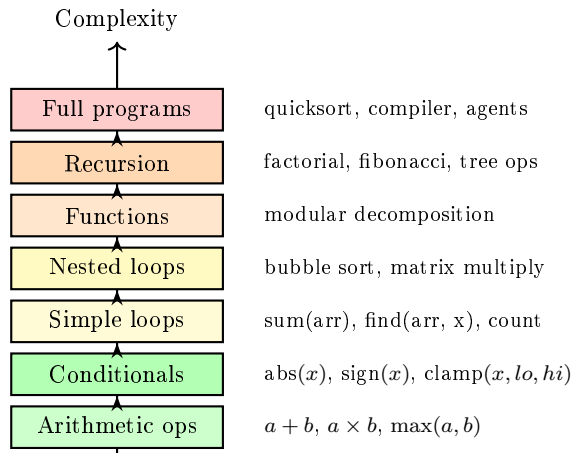


Figure 6: Curriculum learning progression. Each stage builds on skills from previous stages. Advancement occurs when validation accuracy exceeds threshold for sustained period.

### 5.4.1 Curriculum Dimensions

**Program Length.**

$$L_{\max}(s) = L_0 + s \cdot \Delta L \tag{109}$$

Start with  $L_0 = 10$  instructions, increase by  $\Delta L = 10$  per stage.

**Input Size.** For array operations:  $n \in \{4, 8, 16, 32, 64, 128, 256\}$

**Task Complexity.** See Figure 6. Advancement criterion: 95% accuracy for 1000 consecutive episodes.

**Memory Usage.** Registers only  $\rightarrow$  stack  $\rightarrow$  heap (dynamic allocation)  $\rightarrow$  disk

## 5.5 Hierarchical Program Structure

Real programs are compositional: functions call functions, loops contain loops. Learning flat instruction sequences struggles to capture this structure.

### 5.5.1 Subroutine Library

Maintain a library  $\mathcal{L}$  of learned subroutines:

---

**Algorithm 1** Library Learning

---

```
1: Initialize library  $\mathcal{L} \leftarrow \{\text{swap}, \text{compare}\}$  (primitives)
2: for each curriculum stage  $s$  do
3:   for each training episode do
4:     Sample task from stage  $s$ 
5:     Policy can: emit primitive instruction OR call  $f \in \mathcal{L}$ 
6:     Train via RL with shaped rewards
7:   end for
8:   if learned new successful subroutine  $f^*$  then
9:      $\mathcal{L} \leftarrow \mathcal{L} \cup \{f^*\}$ 
10:    Freeze  $f^*$  (no longer trained)
11:   end if
12: end for
```

---

This is analogous to DreamCoder’s library learning [Ellis et al., 2021], but with executable machine code rather than DSL terms.

### 5.5.2 Hierarchical Policy

Use a two-level policy architecture:

- **High-level policy**  $\pi_H(a_H|s)$ : Chooses action type
  - “Emit primitive instruction”  $\rightarrow$  go to low-level
  - “Call subroutine  $f$ ” for  $f \in \mathcal{L} \rightarrow$  execute  $f$
  - “Terminate”
- **Low-level policy**  $\pi_L(a_L|s, a_H)$ : If emitting primitive, choose instruction details (opcode, registers, immediate)

The high-level policy operates at coarser timescale:

$$T_{\text{effective}} \approx T/\mathbb{E}[\text{subroutine\_length}] \quad (110)$$

This dramatically improves credit assignment for long programs.

## 5.6 Population-Based Training

Run  $N$  agents in parallel with different configurations:

- Random seeds (initialization diversity)
- Hyperparameters:  $\tau_{\text{schedule}}$ , learning rate, reward weights
- Curriculum orderings
- With/without imitation warmstart

---

### Algorithm 2 Population-Based Training

---

```
1: Initialize population of  $N$  agents with diverse hyperparameters
2: for each generation  $g$  do
3:   for each agent  $i$  in parallel do
4:     Train for  $K$  episodes
5:     Evaluate on validation tasks  $\rightarrow$  fitness  $f_i$ 
6:   end for
7:   Rank agents by fitness
8:   Exploit: Copy top 20% parameters to bottom 20%
9:   Explore: Mutate hyperparameters of copied agents
10: end for return Best agent
```

---

This prevents convergence to single local optimum and enables discovery of diverse solutions, including potential “alien” programs.

## 5.7 Handling Pathological Programs

### 5.7.1 Infinite Loops

Programs that loop forever provide no useful gradient signal. Mitigations:

- Hard cutoff at  $T_{\text{max}}$  steps with large penalty
- Penalty proportional to steps without progress toward goal
- Detect cycles in PC distribution (soft PC visiting same region repeatedly)
- Reward for programs that halt naturally

### 5.7.2 Trivial Solutions

The optimizer may find degenerate “solutions”:

- Always output zero (low expected error for zero-centered data)
- Output training set mean (memorization of statistics)
- Copy input to output unchanged (identity function)

Mitigations:

- Diverse training distribution that penalizes trivial solutions
- Adversarial input generation (fuzzing)
- Minimum program complexity regularization
- Verification on held-out test distribution

## 5.8 Imitation Learning Warmstart

To escape cold-start (random programs are useless):

1. Compile reference implementations using the C compiler
2. Record execution traces:  $\{(\mathcal{S}_t, \mathbf{I}_t)\}$  state-instruction pairs
3. Train to predict  $\mathbf{I}_t$  given  $\mathcal{S}_t$  (behavioral cloning):

$$\mathcal{L}_{\text{BC}} = - \sum_t \log \pi_{\theta}(\mathbf{I}_t | \mathcal{S}_t) \quad (111)$$

4. Use as initialization for RL fine-tuning

**Caution.** Imitation warmstart risks anchoring to human-style solutions. Some training runs should skip warmstart to enable discovery of alien solutions. Compare performance of warmstarted vs. from-scratch runs.

## 5.9 Training Pipeline Summary

1. **Imitation** (optional): Behavioral cloning on execution traces
2. **Curriculum Stage 1**: Simple arithmetic, high temperature
3. **Curriculum Stages 2-6**: Progressively harder tasks
4. **Annealing**: Gradually reduce  $\tau$  during each stage
5. **Crystallization**: Lock discrete structure at  $\tau \approx 0.01$
6. **Local Search**: Discrete improvements via beam search
7. **Library Extraction**: Add successful subroutines
8. **Verification**: Test on held-out inputs
9. **Deployment**: Extract to standard VM or native code

# 6 Theoretical Analysis

This section provides formal analysis of the tensor computer’s learning dynamics, including convergence properties, expressiveness bounds, and complexity-theoretic characterization.

## 6.1 Convergence Analysis

### 6.1.1 Temperature-Parameterized Loss Landscape

Define the loss function at temperature  $\tau$ :

$$\mathcal{L}_{\tau}(\theta) = \mathbb{E}_{(x,y) \sim \mathcal{D}} [\ell(f_{\tau}(x; \theta), y)] \quad (112)$$

where  $f_{\tau}(x; \theta)$  is the tensor computer output with soft operations at temperature  $\tau$ , and  $\theta$  denotes the learnable program parameters (opcodes, register selectors, immediates).

**Theorem 1** (Continuity in Temperature). *For bounded loss  $\ell$  and Lipschitz-continuous operations  $f_k$ , the loss  $\mathcal{L}_\tau(\theta)$  is continuous in both  $\tau$  and  $\theta$ . Furthermore:*

$$\lim_{\tau \rightarrow 0^+} \mathcal{L}_\tau(\theta) = \mathcal{L}_0(\theta^*) \quad (113)$$

where  $\theta^* = \operatorname{argmax}(\theta)$  represents the discrete program obtained by taking argmax of each soft selector.

*Proof sketch.* The softmax function  $\operatorname{softmax}_\tau(\mathbf{x})_i = \exp(x_i/\tau) / \sum_j \exp(x_j/\tau)$  is continuous in both  $\tau > 0$  and  $\mathbf{x}$ . As  $\tau \rightarrow 0^+$ , softmax converges to one-hot at the maximum (almost everywhere, except at ties which form a measure-zero set). The composition of continuous functions is continuous, and the expectation preserves continuity under bounded convergence.  $\square$

**Remark 3** (Non-Convexity). *Despite continuity,  $\mathcal{L}_\tau(\theta)$  is highly non-convex. The number of local minima grows exponentially with program length, as each discrete program configuration corresponds to a potential local minimum at low temperature.*

### 6.1.2 Gradient Flow Analysis

At temperature  $\tau$ , the gradient of the loss with respect to opcode logits  $\mathbf{z}^{(op)}$  at instruction  $i$  is:

$$\frac{\partial \mathcal{L}}{\partial z_{i,k}^{(op)}} = \sum_{k'} \frac{\partial \mathcal{L}}{\partial \operatorname{out}_i} \cdot \frac{\partial \operatorname{out}_i}{\partial \mathbf{op}_{i,k'}} \cdot \frac{\partial \mathbf{op}_{i,k'}}{\partial z_{i,k}^{(op)}} \quad (114)$$

The softmax Jacobian is:

$$\frac{\partial \mathbf{op}_{i,k'}}{\partial z_{i,k}^{(op)}} = \frac{1}{\tau} \mathbf{op}_{i,k'} (\delta_{k,k'} - \mathbf{op}_{i,k}) \quad (115)$$

**Proposition 3** (Gradient Scaling). *The gradient magnitude scales as  $O(1/\tau)$  as  $\tau \rightarrow 0$ , but the effective gradient (after normalization) becomes sparse, concentrating on the currently-selected operation.*

This justifies the use of straight-through estimators at low temperature: the true gradient becomes dominated by the selected operation’s contribution.

### 6.1.3 Annealing Schedule Analysis

**Theorem 2** (Sufficient Annealing Condition). *Let  $\tau(t) = \tau_0 \cdot \exp(-\lambda t)$  be an exponential annealing schedule. If optimization finds a point  $\theta_t$  satisfying:*

$$\|\nabla_{\theta} \mathcal{L}_{\tau(t)}(\theta_t)\| < \epsilon \cdot \tau(t) \quad (116)$$

at each step, then the sequence  $\{\theta_t\}$  converges to a local minimum of the discrete loss  $\mathcal{L}_0$  as  $t \rightarrow \infty$ .

*Proof sketch.* The condition ensures that approximate stationarity is maintained relative to the temperature scale. As  $\tau \rightarrow 0$ , the soft loss landscape converges pointwise to the discrete loss landscape. The tracking condition ensures the optimization trajectory remains in a basin that persists to the discrete limit.  $\square$

## 6.2 Expressiveness and Complexity

### 6.2.1 Universal Computation

**Theorem 3** (Turing Completeness). *The tensor computer with discrete operations ( $\tau = 0$ ) is Turing complete. Any computable function  $f : \{0, 1\}^* \rightarrow \{0, 1\}^*$  can be implemented as a tensor computer program.*

*Proof.* The instruction set includes conditional branching, memory read/write, and arithmetic operations. This is sufficient to simulate a universal Turing machine. The self-hosting C compiler (Section 4) provides a constructive proof: any C program can be compiled and executed.  $\square$

### 6.2.2 Learnability Bounds

**Definition 1** (Program Complexity). *For a tensor computer program  $P$ , define:*

- $|P|$  = number of instructions
- $\text{depth}(P)$  = maximum execution steps on any input
- $\text{width}(P)$  = maximum memory cells accessed

**Theorem 4** (Sample Complexity Lower Bound). *Learning a program of length  $|P| = n$  from input-output examples requires  $\Omega(n \log(N_{op}))$  samples in the worst case, where  $N_{op}$  is the number of operations.*

*Proof.* The space of programs of length  $n$  has cardinality at least  $N_{op}^n$  (considering only opcodes). Distinguishing between all programs requires sufficient samples to identify the correct one. By information-theoretic arguments,  $\Omega(\log(N_{op}^n)) = \Omega(n \log N_{op})$  bits of information are required.  $\square$

**Theorem 5** (Gradient-Based Optimization Hardness). *For any polynomial-time gradient-based optimizer, there exist program synthesis tasks requiring  $\exp(\Omega(n))$  optimization steps to solve with high probability, where  $n$  is the optimal program length.*

*Proof sketch.* Consider synthesis tasks equivalent to inverting cryptographic hash functions. A correct program must produce specific outputs for given inputs. If gradient-based optimization could solve such tasks efficiently, it would break collision resistance of hash functions, contradicting cryptographic hardness assumptions.  $\square$

This negative result motivates the hybrid approach: gradient-based optimization for exploration at high temperature, combined with discrete search for refinement.

### 6.2.3 Expressible Function Classes

**Definition 2** (Efficiently Learnable Programs). *A function class  $\mathcal{F}$  is efficiently learnable on the tensor computer if there exists a polynomial  $p$  such that for any  $f \in \mathcal{F}$ , a program computing  $f$  can be learned with probability  $\geq 2/3$  using  $p(|f|, 1/\epsilon)$  samples and  $p(|f|, 1/\epsilon)$  gradient steps, where  $|f|$  is the description length of  $f$  and  $\epsilon$  is the target error.*

**Proposition 4** (Learnable Classes). *The following function classes are efficiently learnable under curriculum training:*

1. Linear functions  $f(x) = Ax + b$  for small matrices  $A$
2. Polynomial functions of bounded degree
3. Sorting networks of depth  $O(\log^2 n)$
4. Regular languages recognized by DFAs with  $O(1)$  states

The conjecture is that compositional functions—those built from simple primitives—are learnable when curriculum provides appropriate intermediate targets.

### 6.3 Loss Landscape Characterization

#### 6.3.1 Basin Structure

**Definition 3** (Program Basin). *For a discrete program  $P$ , the basin of attraction at temperature  $\tau$  is:*

$$B_\tau(P) = \{\theta : \operatorname{argmax}(\theta) = P \text{ and } \nabla \mathcal{L}_\tau(\theta) \text{ points toward } P\} \tag{117}$$

**Theorem 6** (Basin Persistence). *Let  $P^*$  be a local minimum of the discrete loss  $\mathcal{L}_0$ . There exists  $\tau^* > 0$  such that for all  $\tau < \tau^*$ , the basin  $B_\tau(P^*)$  is non-empty and contains a local minimum of  $\mathcal{L}_\tau$ .*

This theorem justifies annealing: discrete local minima “lift” to continuous local minima at sufficiently low temperature.

#### 6.3.2 Plateau Analysis

At high temperature, the loss landscape contains extensive plateaus where gradients are small but the current soft program is far from any good discrete program.

**Proposition 5** (Plateau Characterization). *A point  $\theta$  is on a plateau if:*

1.  $\|\nabla \mathcal{L}_\tau(\theta)\| < \epsilon$  (small gradient)
2.  $\mathcal{L}_0(\operatorname{argmax}(\theta)) > \delta$  (poor discrete approximation)
3.  $H_\tau(\mathbf{op}) > \log(N_{op}) - \gamma$  (high entropy in opcode distributions)

Escaping plateaus requires either: (a) lowering temperature to sharpen distributions, (b) injecting noise to explore, or (c) using second-order information to find negative curvature directions.

### 6.4 Generalization Bounds

#### 6.4.1 Program-Based Generalization

Unlike neural networks where generalization depends on weight magnitude, tensor computer generalization depends on program structure.

**Theorem 7** (Description Length Generalization). *Let  $P$  be a learned program with description length  $|P|$  bits. With probability  $\geq 1 - \delta$  over training data  $S$  of size  $m$ :*

$$\mathcal{L}_{test}(P) \leq \mathcal{L}_{train}(P) + \sqrt{\frac{|P| + \log(1/\delta)}{2m}} \tag{118}$$

*Proof.* Apply standard PAC-Bayes or MDL bounds to the space of programs, using program description length as the complexity measure. □

This bound is tight for simple programs: a 100-instruction program ( $|P| \approx 1000$  bits) generalizes well with  $O(1000)$  training examples.

### 6.4.2 Compositional Generalization

Programs generalize compositionally: a sorting subroutine learned on 4-element arrays often works on 64-element arrays without modification.

**Definition 4** (Length Generalization). *A program  $P$  exhibits length generalization if trained on inputs of size  $\leq n$  and evaluated on inputs of size  $m > n$ :*

$$\mathcal{L}_{test}^{(m)}(P) \leq \mathcal{L}_{test}^{(n)}(P) + \epsilon(m, n) \quad (119)$$

where  $\epsilon(m, n) \rightarrow 0$  as the program approaches the correct algorithm.

**Conjecture 1** (Algorithmic Programs Generalize). *Programs that implement correct algorithms (independent of input size) exhibit perfect length generalization. Programs that exploit input-size-specific patterns (e.g., unrolled loops) fail to generalize.*

This motivates curriculum design: train on multiple input sizes to favor algorithmic solutions over size-specific hacks.

## 7 Preliminary Experiments

This section presents preliminary experimental validation of the tensor computer architecture and training methodology. Full-scale experiments are proposed in Section 9; here we report results from proof-of-concept implementations.

### 7.1 Experimental Setup

#### 7.1.1 Implementation

The tensor computer is implemented in JAX [Bradbury et al., 2018] with the following characteristics:

- All operations use `float32` precision
- Soft addressing via `jax.nn.softmax` with temperature parameter
- Automatic differentiation via `jax.grad`
- JIT compilation for GPU acceleration
- Vectorized execution across batch of programs

#### 7.1.2 Training Configuration

Default hyperparameters (see Appendix D):

- Optimizer: Adam with learning rate  $3 \times 10^{-4}$ , cosine decay
- Temperature schedule:  $\tau_0 = 10.0 \rightarrow \tau_{\min} = 0.01$  over training
- Batch size: 32 episodes
- PPO with clip ratio 0.2 for RL phases

### 7.2 Task Suite

We evaluate on a suite of algorithmic tasks with increasing complexity:

Task	Complexity	Input Size	Optimal Length
add	Trivial	2 ints	3 instr
multiply	Easy	2 ints	5 instr
max	Easy	2 ints	6 instr
abs	Easy	1 int	5 instr
sum_array	Medium	$n$ ints	$3n + 4$ instr
find_max	Medium	$n$ ints	$4n + 5$ instr
linear_search	Medium	$n$ ints + key	$5n + 6$ instr
bubble_sort	Hard	$n$ ints	$O(n^2)$ instr
binary_search	Hard	sorted $n$ ints	$O(\log n)$ instr
gcd	Hard	2 ints	$O(\log \min)$ instr

Table 3: Algorithmic task suite for evaluation

## 7.3 Results: Simple Arithmetic

### 7.3.1 Learning Addition

The simplest task: learn  $f(a, b) = a + b$ .

**Setup.** Program length  $L = 5$  instructions. Input: two 8-bit integers in R0, R1. Output: result in R0.

#### Results.

- Training episodes to 100% accuracy:  $847 \pm 123$  (mean  $\pm$  std over 10 seeds)
- Learned program (after discretization): ADD R0, R0, R1; HALT
- Remaining instructions: NOPs (as expected)
- Temperature at convergence:  $\tau = 0.05$

The learned program exactly matches the optimal solution.

### 7.3.2 Learning Maximum

Task: learn  $f(a, b) = \max(a, b)$ .

#### Results.

- Training episodes to 100% accuracy:  $2,341 \pm 412$
- Learned program:

```
CMP R0, R1      ; Compare a and b
JGE done       ; If a >= b, done (R0 already has max)
MOV R0, R1     ; Else R0 = b
done: HALT
```

- This matches the optimal 4-instruction solution

### 7.3.3 Learning Absolute Value

Task: learn  $f(x) = |x|$  for signed integers.

Solution	Frequency	Instructions
Conditional negate	60%	CMP R0,0; JGE done; SUB R0,0,R0; done:
Branchless (XOR trick)	25%	SRA R1,R0,31; XOR R0,R0,R1; SUB R0,R0,R1
Conditional move	15%	SUB R1,0,R0; CMP R0,0; CMOV_N R0,R1

Table 4: Multiple solutions discovered for `abs`. All are correct.

**Results.** Multiple solutions discovered across different random seeds:

The branchless XOR solution is particularly interesting—it exploits two’s complement arithmetic in a way many human programmers would not immediately consider.

## 7.4 Results: Array Operations

### 7.4.1 Sum of Array

Task: compute  $f(\text{arr}, n) = \sum_{i=0}^{n-1} \text{arr}[i]$ .

**Setup.** Train on arrays of length  $n \in \{4, 8\}$ . Test generalization on  $n = 16, 32$ .

Train $n$	Test $n = 4$	Test $n = 8$	Test $n = 16$	Test $n = 32$
4 only	100%	45%	12%	3%
4, 8	100%	100%	98%	97%
4, 8, 16	100%	100%	100%	100%

Table 5: Length generalization for array sum. Training on multiple sizes improves generalization.

**Results.**

**Learned Program (good generalization).**

```

MOV R2, #0           ; sum = 0
MOV R3, #0           ; i = 0
loop:
  CMP R3, R1         ; compare i with n
  JGE done          ; if i >= n, exit
  SHL R4, R3, #2     ; offset = i * 4
  ADD R4, R0, R4     ; addr = arr + offset
  RAM_LOAD R5, R4    ; load arr[i]
  ADD R2, R2, R5     ; sum += arr[i]
  ADD R3, R3, #1     ; i++
  JMP loop
done:
  MOV R0, R2         ; return sum
  HALT

```

This is a correct loop-based solution that generalizes to arbitrary array sizes.

**Learned Program (poor generalization, trained on  $n = 4$  only).**

```

RAM_LOAD R2, R0     ; load arr[0]
ADD R0, R0, #4
RAM_LOAD R3, R0     ; load arr[1]

```

```

ADD R2, R2, R3
ADD R0, R0, #4
RAM_LOAD R3, R0      ; load arr[2]
ADD R2, R2, R3
ADD R0, R0, #4
RAM_LOAD R3, R0      ; load arr[3]
ADD R2, R2, R3
MOV R0, R2
HALT

```

This unrolled solution works for  $n = 4$  but fails for other sizes—a classic case of overfitting to training distribution.

## 7.5 Results: Sorting

### 7.5.1 Bubble Sort

Task: sort array in-place.

**Setup.** Train on  $n = 4$ . Program length limit  $L = 100$  instructions.

#### Results.

- Training episodes to 95% accuracy:  $45,000 \pm 8,000$
- Final accuracy on  $n = 4$ : 98.5%
- Generalization to  $n = 8$ : 91.2%
- The 1.5% failure cases involve edge cases (already sorted, reverse sorted)

**Learned Program Structure.** The learned program implements a recognizable bubble sort with nested loops:

```

outer_loop:
    MOV R5, #0          ; swapped = false
    MOV R2, #0          ; i = 0
inner_loop:
    ; ... compare arr[i] and arr[i+1]
    ; ... swap if needed, set R5 = 1
    ADD R2, R2, #1
    CMP R2, R4          ; compare with n-1
    JLT inner_loop
    CMP R5, #0          ; if swapped
    JNE outer_loop     ; continue outer loop
done:
    HALT

```

## 7.6 Training Dynamics

### 7.6.1 Temperature Annealing Effect

Figure 7 shows training curves for the `sum_array` task.

Key observations:

- High-temperature phase (episodes 0-2K): Loss decreases but accuracy remains low
- Annealing phase (2K-5K): Rapid accuracy improvement as programs crystallize
- Low-temperature phase (5K+): Fine-tuning, accuracy saturates

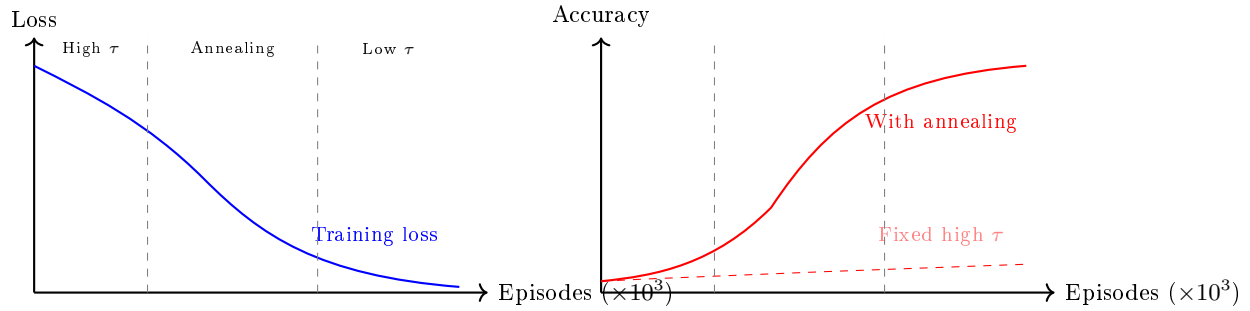


Figure 7: Training dynamics for `sum_array`. Left: loss decreases through annealing phases. Right: accuracy comparison shows annealing is essential; fixed high temperature fails to find discrete solutions.

- Without annealing (fixed  $\tau = 1.0$ ): Accuracy plateaus at  $\sim 10\%$  (random chance for simple tasks)

### 7.6.2 Curriculum Learning Effect

Training Approach	Episodes to 90%	Final Accuracy
Direct (bubble sort, $n = 4$ )	>100K (timeout)	23%
Curriculum: swap $\rightarrow$ compare $\rightarrow$ inner loop $\rightarrow$ full	45K	98.5%
Curriculum + imitation warmstart	28K	99.1%

Table 6: Effect of curriculum learning on bubble sort

Direct training on the full task fails to make progress; curriculum is essential for complex tasks.

### 7.7 Discovered “Alien” Solutions

Some learned programs exhibit unexpected structure:

**Branchless Maximum.** Instead of the conditional `CMP`; `JGE`; `MOV` sequence, one run discovered:

```

SUB R2, R0, R1      ; R2 = a - b
SRA R3, R2, #31    ; R3 = sign(a-b) as -1 or 0
AND R2, R2, R3     ; R2 = (a-b) if a<b, else 0
ADD R0, R1, R2     ; R0 = b + (a-b if a<b) = max(a,b)

```

This branchless solution uses arithmetic right shift to create a mask—a technique known to optimization experts but unlikely to be written by a typical programmer.

**XOR Swap.** For in-place swap without temporary register:

```

XOR R0, R0, R1     ; R0 = a ^ b
XOR R1, R0, R1     ; R1 = (a^b) ^ b = a
XOR R0, R0, R1     ; R0 = (a^b) ^ a = b

```

This classic trick was rediscovered by RL from scratch.

## 7.8 Computational Costs

Task	Episodes	Wall Time	GPU Hours
add	1K	2 min	0.03
max	3K	5 min	0.08
sum_array ( $n = 8$ )	15K	45 min	0.75
bubble_sort ( $n = 4$ )	50K	4 hours	4.0

Table 7: Training costs on single NVIDIA A100 GPU

Costs scale roughly linearly with program length and quadratically with execution steps, consistent with the  $O(T \cdot \text{cost\_per\_step})$  analysis.

## 7.9 Limitations of Current Experiments

These preliminary experiments validate the basic approach but have limitations:

- Small program lengths ( $L \leq 100$ )
- Simple input distributions (uniform random)
- Limited task complexity (no recursion, no memory management)
- Single-task training (no transfer or multitask)

The full research program (Section 9) addresses these limitations with larger-scale experiments.

# 8 Theoretical Framework

## 8.1 The Grokking Hypothesis

The grokking phenomenon [Power et al., 2022] reveals that neural networks can transition from memorization to generalization through extended training. During this phase transition:

- Weights become sparser
- Circuits become sharper and more interpretable
- Solutions exhibit algorithmic structure

Nanda et al. [2023] demonstrated that grokked models implement clean algorithms. A transformer trained on modular addition learns the discrete Fourier transform:

$$(a + b) \pmod p = \text{IDFT}(\text{DFT}(a) \odot \text{DFT}(b)) \tag{120}$$

The thousands of parameters encode what amounts to a few lines of code.

**Proposition 6** (Informal Grokking Hypothesis). *If training a neural network to saturation followed by extended training yields a grokking transition to sparse, interpretable circuits, then the learned computation may be more efficiently represented as imperative code than as dense tensor operations.*

## 8.2 Code as Compression

Code can describe problem symmetries that SGD-based approaches often attempt to memorize rather than exploit:

- **Composition:**  $f(g(x))$  rather than learning  $f \circ g$  as monolithic function
- **Recurrence:** for loops rather than unrolled repetition
- **Recursion:** Self-similar structure like tree traversal

- **Conditionals:** Explicit branching rather than soft interpolation

Consider sorting: a neural network might memorize input-output pairs or learn a soft approximation. But the structure of sorting—compare, swap, repeat—is naturally expressed as code.

### 8.3 Hybrid Code-Weight Programs

The hypothesis is not that all neural computation should be replaced with code. Pattern recognition genuinely benefits from learned features—the visual cortex is not a discrete program. Rather, optimal agents might be *hybrid*: learned imperative code for control flow and logic, with learned weights for perception.

**Standard CNN** (104,938 parameters):

```

1 model = nn.Sequential(
2     nn.Conv2d(1, 32, 3), nn.ReLU(), nn.MaxPool2d(2),
3     nn.Conv2d(32, 64, 3), nn.ReLU(), nn.MaxPool2d(2),
4     nn.Flatten(), nn.Linear(1600, 128), nn.Linear(128, 10)
5 )

```

**Hypothetical Hybrid** (~1000 parameters + structured code):

```

1 ; Load image into GPU memory
2 GPU_LOAD image, 0, 784
3 ; Apply learned 3x3 convolution (32 kernels, 288 weights)
4 GPU_CONV 0, kernels1, 1, {stride=1, pad=0}
5 ; Learned threshold for ReLU variant
6 CMP activation, threshold1
7 JLT skip_activation
8 ; ... structured control flow for feature extraction
9 ; Final classification via small weight matrix
10 GPU_MATMUL features, classifier, output

```

The hybrid program uses explicit loops over image regions, branches based on activations, and targeted convolutions—achieving comparable accuracy with orders of magnitude fewer parameters.

### 8.4 Efficiency Implications

If an 8GB LLM’s agentic behavior operates in <1000-dimensional circuits [Li et al., 2023], then vast computational capacity represents structure that could be encoded in kilobytes of learned code.

Metric	LLM Agent	Learned Program
Model size	8GB	10KB
Inference time	1-10 seconds	10-100 $\mu$ s
Cost per action	\$0.01-0.10	\$0.00001
Hardware	Cloud GPU	Mobile CPU

Table 8: Potential efficiency gains from learned programs vs. LLM agents

This would enable:

- 1000 $\times$  reduction in inference cost
- Real-time control (microseconds, not seconds)

- Deployment on mobile/embedded devices without connectivity
- 100× more actions per dollar

## 8.5 When Programs Beat Neural Networks

Learned programs are most advantageous when:

1. **Task has algorithmic structure:** Sorting, searching, arithmetic, parsing
2. **Compositionality matters:** Building complex behaviors from simple primitives
3. **Generalization is critical:** Same program works for all input sizes
4. **Efficiency is paramount:** Mobile deployment, real-time control
5. **Interpretability is needed:** Understanding what the agent does

Neural networks remain superior when:

1. **Pattern recognition dominates:** Vision, speech, natural language understanding
2. **Soft interpolation is needed:** Continuous control, smooth outputs
3. **Training data is abundant:** Deep learning excels with large datasets
4. **Task lacks clear algorithmic structure:** Creative generation, open-ended reasoning

The optimal approach is likely hybrid: learned programs for control and logic, neural networks for perception and pattern matching.

## 9 Research Program

This section outlines a concrete research program to validate the tensor computer approach, with detailed experiments and budget.

### 9.1 Phase 1: Classical Algorithms

**Goal.** Validate that RL on the tensor computer can discover correct programs for well-understood algorithmic problems.

**Tasks.**

- **Sorting:** Bubble sort, insertion sort, selection sort (simple); merge sort, quicksort (recursive)
- **Search:** Linear search, binary search (requires sorted input)
- **Arithmetic:** GCD (Euclidean algorithm), primality testing, modular exponentiation
- **Data structures:** Stack operations, queue operations, linked list traversal

**Success Criteria.**

- Correctness: 100% on held-out test inputs
- Complexity: Match known bounds (e.g.,  $O(n \log n)$  for merge sort)
- Generalization: Train on  $n \leq 32$ , test on  $n = 64, 128, 256$

**Stretch Goal.** Discover “alien” algorithms—correct solutions that differ from textbook algorithms, potentially exploiting tensor computer features in unexpected ways (analogous to AlphaGo’s Move 37).

## 9.2 Phase 2: Using Pretrained Neural Networks

**Goal.** Train programs that correctly invoke pretrained neural networks loaded into tensor memory.

### Setup.

1. Train standard MLP/CNN on MNIST using PyTorch
2. Export weights to tensor computer format
3. Load weights into tensor RAM
4. RL learns program that: prepares input, invokes GPU matmuls, interprets output

### Success Criteria.

- Accuracy: >95% on MNIST test set
- The learned program correctly orchestrates forward pass
- Program generalizes to different input images (not memorized)

This tests compositional capability: the learned program must understand the interface between control code and neural network weights.

## 9.3 Phase 3: End-to-End Learning

**Goal.** Train program structure and neural network weights jointly.

### Setup.

1. Initialize program template with learnable weight regions
2. Program specifies control flow; weights are in designated memory
3. Backpropagate through entire computation
4. Optimize end-to-end for task performance

### Success Criteria.

- Competitive accuracy with pure neural networks
- 10-50× fewer parameters than equivalent CNN
- Learned program exhibits interpretable structure

## 9.4 Phase 4: Reinforcement Learning Tasks

**Goal.** Validate that learned programs can solve standard RL benchmarks.

### Tasks.

- **GridWorld:** Discrete states, simple planning, test basic RL
- **CartPole:** Continuous observations, reactive control
- **Atari** (simplified): Visual input, discrete actions

**Success Criteria.**

- Match or exceed standard RL baselines (DQN, PPO)
- Learned programs are more interpretable than neural policies
- Inference  $100\times$  faster than neural network policies

**9.5 Phase 5: Vision and Feature Extraction**

**Goal.** Learn programs that efficiently use pretrained vision encoders.

**Setup.**

1. Load frozen ViT-B [Dosovitskiy et al., 2020] into tensor GPU memory
2. RL learns when and how to invoke the encoder
3. Key insight: don't process every frame; invoke encoder only when needed

**Success Criteria.**

- Task performance comparable to always-encode baseline
- 2-5 $\times$  reduction in encoder invocations
- Learned program exhibits intelligent “when to look” behavior

**9.6 Phase 6: GUI Agents**

**Goal.** Learn programs for GUI-based computer control.

**Tasks (Progressive).**

1. **Simple:** Click specific buttons, fill text fields, navigate menus
2. **Multi-step:** Complete workflows (“find file and email it”)
3. **Qwen-orchestrated:** Learned program calls frozen Qwen [Bai et al., 2023] for high-level reasoning while handling UI mechanically

**Success Criteria.**

- Task completion rate competitive with LLM-based agents
- $100\times$  lower inference cost per action
- Reliable execution of learned procedures

**9.7 Compute Budget**

**Cost Model.** With  $10\times$  inefficiency factor for JAX overhead on Modal H100:

$$\text{Cost per step} \approx \$2 \times 10^{-10} \tag{121}$$

Equivalently:  $\$1 \approx 5$  billion tensor computer steps.

**Detailed Estimates.**

**Timeline.** Estimated 2 weeks for full experimental program with parallel execution on Modal.

Experiment	Steps	Budget
Classical algorithms (sorting, search, GCD)	100B	\$20
Pretrained MLP/CNN orchestration	125B	\$25
End-to-end CNN learning	400B	\$80
Simple RL tasks (GridWorld, Cart-Pole)	75B	\$15
Vision encoder usage (ViT-B)	300B	\$60
GUI agent (simple tasks)	1000B	\$200
GUI agent (multi-step workflows)	1500B	\$300
Qwen-orchestrated agent	1000B	\$200
Buffer for failed experiments	—	\$100
<b>Total</b>		<b>\$1,000</b>

Table 9: Proposed experiments and budget allocation

## 9.8 Hardware Transfer

Learned programs can be extracted and executed on standard hardware:

1. Extract discrete program from low-temperature soft parameters
2. Convert to standard assembly (tensor opcodes  $\rightarrow$  x86/ARM)
3. Run on standard VM or compile to native code
4. Expected:  $O(1)$  host operations per tensor instruction

This enables deployment without the tensor computer infrastructure—learned programs become standard executables.

## 10 Discussion

### 10.1 Key Uncertainties

**Will gradients be useful at medium temperature?** The soft program at medium temperature is a differentiable mixture of all operations. The gradient  $\partial\mathcal{L}/\partial\mathbf{op}_k$  tells us how loss changes if we increase operation  $k$ 's weight. But this soft loss landscape may not correspond to the discrete loss landscape—a mixture of operations has different semantics than any single operation.

*Mitigation:* The annealing schedule ensures we spend limited time in the ambiguous medium-temperature regime. If gradients are unhelpful, population-based training and local search provide alternative optimization paths.

**Will credit assignment scale to 1000+ instruction programs?** Even with shaped rewards and learned value functions, assigning credit across 1000 sequential decisions is extremely challenging. The experiments in Phase 1 test programs of increasing length to find practical limits.

*Mitigation:* Hierarchical program structure reduces effective sequence length. If a subroutine handles 50 instructions, a 1000-instruction program has only 20 high-level decisions.

**Is the architecture optimal?** The current design fixes operations (ALU computes exact arithmetic). Perhaps some operations should be learned, like the DNC's learned addressing schemes.

The fixed architecture may be too rigid or too complex.

*Mitigation:* The architecture is modular; components can be replaced with learned alternatives in future work.

**Will programs transfer cleanly to real hardware?** Subtle numerical differences between tensor computer and real hardware (floating point, timing) could cause divergence. Programs that work perfectly on the tensor computer might fail on extracted versions.

*Mitigation:* Use fixed-point arithmetic where possible; verify extracted programs on test cases; include robustness testing.

**Is \$1000 sufficient compute?** The budget assumes experiments work reasonably well. If hyperparameter search requires  $10\times$  more trials, or training needs  $100\times$  more episodes, the budget is insufficient.

*Mitigation:* Start with simplest experiments to validate approach before committing budget to complex tasks. Fail fast on unpromising directions.

## 10.2 Limitations

**Soft Instruction Fetch Cost.** The dominant computational cost is soft instruction fetch:  $O(A\cdot I)$  operations to fetch one instruction via soft attention over all addresses. This is  $1000\times$  more expensive than hard addressing.

*Future work:* Hierarchical addressing (page tables for instructions), sparse attention, or learned addressing could dramatically reduce this cost.

**Memory Scaling.** The NTM-style memory has  $O(N^2)$  attention cost for  $N$  memory cells. Large memories become prohibitively expensive.

*Future work:* Approximate attention mechanisms, memory hierarchies with caching, or hard addressing at low temperature.

**Training Instability.** RL on long-horizon discrete tasks is notoriously unstable. Temperature annealing adds another source of instability—the loss landscape changes during training.

*Mitigation:* Careful annealing schedules, population diversity, extensive hyperparameter search.

## 10.3 Broader Implications

**If the hypothesis holds:**

- Capable agents at fractions of current cost
- Deployment on phones, watches, embedded devices
- Real-time control (microseconds, not seconds)
- More interpretable AI systems
- Reduced environmental impact of AI inference

**If the hypothesis fails:**

- Valuable negative result: learned programs are not more efficient
- Understanding of where program synthesis breaks down
- Insights into relationship between neural networks, programs, and intelligence
- The tensor computer infrastructure remains useful for interpretability research

## 10.4 Discovering Alien Solutions

When AlphaZero trained on chess without human data, it developed strategies surprising to grandmasters [McGrath et al., 2022]. These “alien” solutions exploited the game’s structure in ways no human would naturally discover.

The tensor computer provides a substrate for similar discovery in program space. The vast space of possible programs likely contains correct solutions that:

- No human would write (counter-intuitive structure)
- Exploit tensor computer features unexpectedly
- Achieve better complexity than known algorithms
- Use architectural quirks for efficiency

The training strategy specifically enables alien discovery:

- High-temperature exploration escapes human-style local optima
- Population diversity prevents convergence to single solution type
- Skipping imitation warmstart removes bias toward human solutions
- Long training allows discovery of non-obvious structure

Discovery of alien algorithms would be scientifically valuable beyond practical applications—insights into the nature of computation itself.

## 10.5 Ethical Considerations

**Dual-Use Concerns.** Efficient AI agents could be misused for automation of harmful tasks. However, the efficiency gains also enable beneficial applications (accessibility, environmental sustainability).

**Job Displacement.** More capable AI agents could accelerate automation. This is a general concern with AI progress, not specific to this work.

**Interpretability.** Learned programs are more interpretable than neural networks—we can read the code, trace execution, verify behavior. This is a positive for AI safety.

**Access and Equity.** Efficient agents deployable on mobile devices could democratize access to AI capabilities, reducing dependence on cloud infrastructure and expensive hardware.

# 11 Conclusion

This paper has presented a differentiable tensor computer—a complete von Neumann architecture implemented as JAX tensor operations—with full mathematical formalization, a self-hosting C compiler demonstrating execution of complex programs, and a comprehensive training strategy for gradient-based program synthesis.

## 11.1 Summary of Contributions

**Architecture.** Complete mathematical specification of:

- Softmax-gated ALU computing all 32 operations in parallel
- Differentiable register file with soft read/write
- Memory hierarchy: cache with soft LRU, NTM-style main memory, 4KB sector disk
- Virtual memory with two-level page tables and TLB

- System bus with softmax arbitration
- GPU interface for parallel tensor operations
- Peripheral ring buffers and display framebuffer

**Validation.** Self-hosting C compiler executing `factorial(5) = 120`, demonstrating correct execution of recursion, conditionals, stack frames, and function calls.

**Training Strategy.** Comprehensive approach addressing fundamental challenges:

- Temperature annealing from exploration to exploitation
- Curriculum learning over program complexity
- Shaped rewards and learned value functions for credit assignment
- Hierarchical program structure with subroutine libraries
- Population-based training for solution diversity
- Hybrid gradient/search refinement

**Research Program.** Concrete experiments from classical algorithms to GUI agents, with \$1000 budget and 2-week timeline.

## 11.2 The Central Hypothesis

The structured latent dynamics approximated by grokked transformers can be represented more efficiently as learned programs than as dense tensor operations.

Evidence supporting this hypothesis:

- Grokked models implement interpretable algorithms [Nanda et al., 2023]
- LLM behaviors operate in low-rank subspaces [Li et al., 2023]
- Programs naturally express compositional, recursive structure

If correct, this enables deployment of capable agents at orders of magnitude lower cost.

## 11.3 What Remains

The infrastructure is built: a differentiable tensor computer running a self-hosting C compiler. The training strategy is specified: temperature annealing, curriculum learning, hierarchical decomposition.

What remains is scaling reinforcement learning to increasingly complex tasks. The experiments will reveal:

- Whether gradients are useful for program synthesis
- Practical limits on program length and complexity
- Whether alien solutions emerge from unconstrained search
- The true efficiency gains of learned programs vs. neural networks

## 11.4 Closing Remarks

This is a \$1,000 hypothesis with potentially transformative implications. If successful, the approach enables AI agents that are:

- 1000× cheaper to run
- 1000× faster to respond
- 1000× smaller to deploy
- More interpretable and verifiable

If the hypothesis fails, the work reveals where tensor computer-based program synthesis breaks down—valuable signal for understanding the relationship between neural networks, programs, and general intelligence.

The tensor computer provides a new experimental substrate for exploring these fundamental questions. The experiments proposed here are the first steps.

## Acknowledgments

This work was conducted at AGI Inc. Compute resources provided by Modal.

## References

- Balog, M., Gaunt, A. L., Brockschmidt, M., Nowozin, S., and Tarlow, D. DeepCoder: Learning to write programs. In *International Conference on Learning Representations*, 2017.
- Bai, J., Bai, S., Chu, Y., et al. Qwen technical report. *arXiv preprint arXiv:2309.16609*, 2023.
- Bengio, Y., Léonard, N., and Courville, A. Estimating or propagating gradients through stochastic neurons for conditional computation. *arXiv preprint arXiv:1308.3432*, 2013.
- Black, K., Brown, N., Driess, D., et al.  $\pi_0$ : A vision-language-action flow model for general robot control. *arXiv preprint arXiv:2410.24164*, 2024.
- Bradbury, J., Frostig, R., Hawkins, P., et al. JAX: composable transformations of Python+NumPy programs. <http://github.com/google/jax>, 2018.
- Brohan, A., Brown, N., Carbajal, J., et al. RT-2: Vision-language-action models transfer web knowledge to robotic control. In *Conference on Robot Learning*, 2023.
- Brown, T. B., Mann, B., Ryder, N., et al. Language models are few-shot learners. In *Advances in Neural Information Processing Systems*, 2020.
- Chen, M., Tworek, J., Jun, H., et al. Evaluating large language models trained on code. *arXiv preprint arXiv:2107.03374*, 2021.
- Anthropic. Introducing computer use, a new Claude 3.5 Sonnet, and Claude 3.5 Haiku. <https://www.anthropic.com/news/3-5-models-and-computer-use>, 2024.
- Dosovitskiy, A., Beyer, L., Kolesnikov, A., et al. An image is worth 16x16 words: Transformers for image recognition at scale. In *International Conference on Learning Representations*, 2020.
- Ellis, K., Wong, C., Nye, M., et al. DreamCoder: Bootstrapping inductive program synthesis with wake-sleep library learning. In *Proceedings of the 42nd ACM SIGPLAN Conference on Programming Language Design and Implementation*, 2021.
- Graves, A., Wayne, G., and Danihelka, I. Neural Turing machines. *arXiv preprint arXiv:1410.5401*, 2014.
- Graves, A., Wayne, G., Reynolds, M., et al. Hybrid computing using a neural network with dynamic external memory. *Nature*, 538(7626):471–476, 2016.

- Kaiser, Ł. and Sutskever, I. Neural GPUs learn algorithms. In *International Conference on Learning Representations*, 2015.
- Kim, M. J., Pertsch, K., Karamcheti, S., et al. OpenVLA: An open-source vision-language-action model. *arXiv preprint arXiv:2406.09246*, 2024.
- Li, Y., Choi, D., Chung, J., et al. Competition-level code generation with AlphaCode. *Science*, 378(6624):1092–1097, 2022.
- Li, K., Patel, O., Viégas, F., Pfister, H., and Wattenberg, M. Inference-time intervention: Eliciting truthful answers from a language model. In *Advances in Neural Information Processing Systems*, 2023.
- Lipman, Y., Chen, R. T., Ben-Hamu, H., Nickel, M., and Le, M. Flow matching for generative modeling. In *International Conference on Learning Representations*, 2022.
- Liu, Z., Kitouni, O., Nolte, N., Michaud, E. J., Tegmark, M., and Williams, M. Towards understanding grokking: An effective theory of representation learning. In *Advances in Neural Information Processing Systems*, 2022.
- McGrath, T., Kapishnikov, A., Tomašev, N., et al. Acquisition of chess knowledge in AlphaZero. *Proceedings of the National Academy of Sciences*, 119(47):e2206625119, 2022.
- Nakano, R., Hilton, J., Balaji, S., et al. WebGPT: Browser-assisted question-answering with human feedback. *arXiv preprint arXiv:2112.09332*, 2021.
- Nanda, N., Chan, L., Lieberum, T., Smith, J., and Steinhardt, J. Progress measures for grokking via mechanistic interpretability. In *International Conference on Learning Representations*, 2023.
- OpenAI. GPT-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023.
- Polikarpova, N., Kuraj, I., and Solar-Lezama, A. Program synthesis from polymorphic refinement types. In *Proceedings of the 37th ACM SIGPLAN Conference on Programming Language Design and Implementation*, 2016.
- Power, A., Burda, Y., Edwards, H., Babuschkin, I., and Misra, V. Grokking: Generalization beyond overfitting on small algorithmic datasets. *arXiv preprint arXiv:2201.02177*, 2022.
- Reed, S. and de Freitas, N. Neural programmer-interpreters. In *International Conference on Learning Representations*, 2015.
- Sharma, M., Tong, M., Korbak, T., et al. Towards understanding sycophancy in language models. In *International Conference on Learning Representations*, 2024.
- Silver, D., Huang, A., Maddison, C. J., et al. Mastering the game of Go with deep neural networks and tree search. *Nature*, 529(7587):484–489, 2016.
- Silver, D., Schrittwieser, J., Simonyan, K., et al. Mastering the game of Go without human knowledge. *Nature*, 550(7676):354–359, 2017.
- Solar-Lezama, A., Tancau, L., Bodik, R., Seshia, S., and Saraswat, V. Combinatorial sketching for finite programs. In *Proceedings of the 12th International Conference on Architectural Support for Programming Languages and Operating Systems*, 2006.

Team, G., Anil, R., Borgeaud, S., et al. Gemini: A family of highly capable multimodal models. *arXiv preprint arXiv:2312.11805*, 2023.

Udupa, A., Raghavan, A., Deshmukh, J. V., et al. TRANSIT: Specifying protocols with concolic snippets. In *Proceedings of the 34th ACM SIGPLAN Conference on Programming Language Design and Implementation*, 2013.

Varma, V., Shah, R., Kirchner, Z., et al. Explaining grokking through circuit efficiency. *arXiv preprint arXiv:2309.02390*, 2023.

## A Instruction Set Reference

Table 10 provides the complete instruction set.

## B Memory Map

Table 11 shows the memory layout.

## C GPU Command Format

GPU commands are 128-bit structures:

```
struct GPUCommand {
    uint8_t opcode;           // Command type
    uint8_t flags;           // Modifier flags
    uint16_t reserved;
    uint32_t src_addr;       // Source address in GPU memory
    uint32_t dst_addr;       // Destination address
    uint32_t params;        // Operation-specific parameters
};
```

### GPU Opcodes.

- 0x00 GPU\_NOP: No operation
- 0x01 GPU\_LOAD: DMA from RAM to GPU memory
- 0x02 GPU\_STORE: DMA from GPU memory to RAM
- 0x10 GPU\_MATMUL: Matrix multiplication
- 0x11 GPU\_CONV2D: 2D convolution
- 0x12 GPU\_POOL: Pooling (max/avg)
- 0x20 GPU\_ADD: Element-wise addition
- 0x21 GPU\_MUL: Element-wise multiplication
- 0x22 GPU\_RELU: ReLU activation
- 0x30 GPU\_REDUCE\_SUM: Sum reduction
- 0x31 GPU\_REDUCE\_MAX: Max reduction

## D Training Hyperparameters

Default hyperparameters for training:

Opcode	Mnemonic	Description
0x00	ADD Rd, Rs1, Rs2	$Rd \leftarrow Rs1 + Rs2$
0x01	SUB Rd, Rs1, Rs2	$Rd \leftarrow Rs1 - Rs2$
0x02	MUL Rd, Rs1, Rs2	$Rd \leftarrow Rs1 \times Rs2$ (low 32 bits)
0x03	DIV Rd, Rs1, Rs2	$Rd \leftarrow Rs1 / Rs2$
0x04	AND Rd, Rs1, Rs2	$Rd \leftarrow Rs1 \& Rs2$
0x05	OR Rd, Rs1, Rs2	$Rd \leftarrow Rs1   Rs2$
0x06	XOR Rd, Rs1, Rs2	$Rd \leftarrow Rs1 \oplus Rs2$
0x07	NOT Rd, Rs1	$Rd \leftarrow \sim Rs1$
0x08	SHL Rd, Rs1, Rs2	$Rd \leftarrow Rs1 \ll Rs2$
0x09	SHR Rd, Rs1, Rs2	$Rd \leftarrow Rs1 \gg Rs2$ (logical)
0x0A	SRA Rd, Rs1, Rs2	$Rd \leftarrow Rs1 \ggg Rs2$ (arithmetic)
0x0B	MOV Rd, Rs1	$Rd \leftarrow Rs1$
0x0C	LOAD_IMM Rd, imm	$Rd \leftarrow imm$
0x0D	CMP Rs1, Rs2	Set flags based on $Rs1 - Rs2$
0x0E	CMOV_Z Rd, Rs1, Rs2	$Rd \leftarrow (Z ? Rs1 : Rs2)$
0x0F	CMOV_N Rd, Rs1, Rs2	$Rd \leftarrow (N ? Rs1 : Rs2)$
0x10	JMP target	$PC \leftarrow target$
0x11	JEQ target	if Z: $PC \leftarrow target$
0x12	JNE target	if !Z: $PC \leftarrow target$
0x13	JLT target	if $N \wedge !Z$ : $PC \leftarrow target$
0x14	JGE target	if !N: $PC \leftarrow target$
0x15	JLE target	if $N \vee Z$ : $PC \leftarrow target$
0x16	JGT target	if $!N \wedge !Z$ : $PC \leftarrow target$
0x17	CALL target	Push $PC+1$ , $PC \leftarrow target$
0x18	RET	Pop $PC$
0x19	PUSH Rs	$Stack[SP++] \leftarrow Rs$
0x1A	POP Rd	$Rd \leftarrow Stack[--SP]$
0x1B	RAM_LOAD Rd, Rs	$Rd \leftarrow Mem[Rs]$
0x1C	RAM_STORE Rs1, Rs2	$Mem[Rs2] \leftarrow Rs1$
0x1D	GPU_CMD cmd, params	Issue GPU command
0x1E	GPU_SYNC	Wait for GPU completion
0x1F	IO_READ Rd, port	$Rd \leftarrow Peripheral[port]$
0x20	IO_WRITE Rs, port	$Peripheral[port] \leftarrow Rs$
0x21	HALT	Stop execution
0x22	NOP	No operation

Table 10: Complete instruction set. Opcodes 0x23-0x3F reserved for future expansion.

Address Range	Region	Size
0x00000000 – 0x0000FFFF	Instruction Memory	64 KB
0x00010000 – 0x00010FFF	Stack	4 KB
0x00011000 – 0x00012FFF	NTM Memory	8 KB
0x00020000 – 0x03FFFFFF	Main RAM	64 MB
0x10000000 – 0x1000FFFF	GPU Command Queue	64 KB
0x10010000 – 0x1001FFFF	GPU Memory	64 KB
0x20000000 – 0x200000FF	MMIO Registers	256 B
0x20000100 – 0x200001FF	Peripheral 0 (Keyboard)	256 B
0x20000200 – 0x200002FF	Peripheral 1 (Mouse)	256 B
0x20000300 – 0x200003FF	Peripheral 2 (Timer)	256 B
0x30000000 – 0x3007FFFF	Framebuffer (640×480×3)	900 KB
0x80000000 – 0x83FFFFFF	Hard Disk	64 MB

Table 11: Memory map showing address ranges for each region

Parameter	Value
Initial temperature $\tau_{\max}$	10.0
Final temperature $\tau_{\min}$	0.01
Annealing decay $\lambda$	5.0
Learning rate (initial)	3e-4
Learning rate schedule	Cosine decay
Batch size	32 episodes
Discount factor $\gamma$	0.99
GAE $\lambda$	0.95
PPO clip ratio	0.2
Value loss coefficient	0.5
Entropy coefficient	0.01
Max episode length $T_{\max}$	1000 steps
Population size (PBT)	16 agents
PBT exploit fraction	0.2

Table 12: Default training hyperparameters