

# Engram V3: A Cognitive Mesh Architecture for Abstract Reasoning

## Combining Autoregressive Transformers, Grid Diffusion, and World Model Simulation

---

**Author:** Vincent Kaufmann

**Date:** February 23, 2026

**Status:** Architecture design complete; V3-mini (2.5B) targeted for DGX Spark prototype; V3-full (190B) targeted for cloud training

---

### Abstract

---

Contemporary approaches to abstract reasoning rely on autoregressive language models that generate solutions token-by-token in a single left-to-right pass. This sequential commitment is fundamentally mismatched with the structure of spatial reasoning tasks, where humans perceive entire grids simultaneously, form approximate hypotheses, and refine iteratively. This paper presents Engram V3, a hybrid cognitive mesh architecture that addresses this mismatch by integrating five computational layers: a biologically-structured state engine, multi-modal perception, an autoregressive transformer backbone with mixture-of-experts specialization, a Grid Diffusion Transformer (GDiT) for spatial reasoning through iterative denoising, and a metacognitive self-monitoring system.

The Grid Diffusion Transformer introduces MaskGIT-style discrete diffusion over small categorical grids (3x3 to 30x30, 10 colors), replacing the autoregressive bottleneck with bidirectional iterative refinement conditioned on few-shot example pairs, task features, and a 72-dimension internal state vector. A diffusion-based world model reuses the same GDiT architecture at reduced fidelity to simulate hypothetical outcomes before committing to execution --- implementing mental simulation where denoising quality itself serves as the confidence signal, eliminating the need for a separate critic network. A learned reasoning mode router dynamically selects between four reasoning strategies (autoregressive, GDiT diffusion, tree search, and memory retrieval) and can switch modes mid-sequence.

The architecture incorporates a 42-function domain-specific language organized into six skill clusters, with 10 memorized compound skills mapped to mixture-of-experts groups through library learning. The 72-dimension state engine from Engram V2 conditions every component --- modulating GDiT denoising through FiLM, biasing expert routing, controlling attention temperature, and driving the metacognitive verification loop that runs after each iterative unmasking step.

V3-mini (2.5B total parameters, 1.6B active) serves as a research prototype trainable on a single DGX Spark GB10 in approximately 3--4 days at zero cloud cost, producing a 1.7 GB inference model that runs on consumer

hardware. V3-full (190B total parameters, 48B active) targets competitive performance on ARC-AGI through 48 MoE experts per layer across 40 transformer layers, deployable on a single DGX Spark at MXFP4 quantization (95 GB). We project V3-mini achieving 35--50% on ARC-AGI public evaluation and V3-full achieving 65--80%, competitive with frontier systems 100x larger.

Novel contributions include: the first application of discrete diffusion to few-shot abstract grid reasoning; unified output generation and world model simulation through two-fidelity denoising; state-conditioned discrete diffusion with biological state vectors; a learned reasoning mode router with mid-sequence switching; and the integration of all five components into a single differentiable cognitive mesh.

---

## Table of Contents

---

1. [Introduction](#)
2. [Related Work](#)
3. [Architecture Overview](#)
4. [Grid Diffusion Transformer \(GDiT\)](#)
5. [Diffusion World Model](#)
6. [Skill Decomposition and Library Learning](#)
7. [State Engine and Affective Computing](#)
8. [Reasoning Mode Router](#)
9. [Memory System](#)
10. [Level of Detail Context Management](#)
11. [Neural Fabric Protocol \(HULLHB\)](#)
12. [Speculative Cognitive Pipelining](#)
13. [Non-Euclidean Elastic Representations](#)
14. [Training Pipeline](#)
15. [V3-mini: Proof of Concept \(2.5B\)](#)
16. [V3-full: Competition Model \(190B\)](#)
17. [Performance Projections](#)
18. [Novel Contributions](#)
19. [Limitations and Future Work](#)
20. [Conclusion](#)
21. [Appendix A: GDiT Architecture Details](#)
22. [Appendix B: 42-Function DSL Reference](#)
23. [Appendix C: Training Hyperparameters](#)
24. [Appendix D: V3-mini Memory Budget](#)

---

## 1. Introduction

---

### 1.1 The Autoregressive Bottleneck

Autoregressive transformers generate output one token at a time, each conditioned on all previous tokens. This left-to-right commitment has produced extraordinary results in natural language --- GPT-4 (OpenAI, 2023), Gemini (Google, 2024), and Llama 3 (Meta, 2024) demonstrate that autoregressive modeling scales to remarkable capability in linguistic domains. Yet these same models struggle conspicuously on tasks requiring spatial reasoning, particularly the Abstraction and Reasoning Corpus (ARC-AGI) introduced by Chollet (2019).

The mismatch is architectural. Consider a 10x10 ARC grid with 10 possible colors per cell. The autoregressive approach must serialize this 2D structure into a 1D token sequence --- typically row-by-row or as JSON --- then generate 100 cell values sequentially, where each cell's prediction is irrevocably committed before subsequent cells are considered. Cell (0,0) is generated before the model has any information about cell (9,9), even when the correct transformation requires global spatial coherence (e.g., rotating the entire grid 90 degrees).

This serialization destroys three properties essential to spatial reasoning:

**Simultaneity.** Humans perceive grids holistically. When examining an ARC puzzle, a human sees the entire input grid at once, identifies spatial patterns (symmetry, repetition, boundaries), and forms a global hypothesis before attending to individual cells. Autoregressive models cannot form global hypotheses --- they commit to local decisions sequentially.

**Iterative refinement.** Human spatial reasoning is inherently iterative. A chess player considers a move, mentally traces its consequences, revises, considers another move, and gradually converges on a decision. Each revision improves upon the previous attempt. Autoregressive generation permits no revision --- each token is final.

**Spatial locality.** In a 2D grid, cell (5,5) is spatially adjacent to cells (4,5), (6,5), (5,4), and (5,6). When serialized to 1D, these neighbors may be 10 tokens apart (different rows) or 1 token apart (same row), destroying the uniform spatial topology. Standard positional encodings (RoPE, sinusoidal) are 1D and cannot represent 2D adjacency.

These limitations are not theoretical. On the ARC-AGI public evaluation set (400 tasks), GPT-4o with careful prompting achieves approximately 50--60% (OpenAI, 2024). Gemini 1.5 Pro scores similarly. Ryan Greenblatt's 2024 submission, which wraps an LLM in a search framework, reached 72%. MindsAI's 2024 winning system achieved 84%. The gap between LLM performance and human performance (approximately 85--95%) is concentrated in tasks requiring spatial manipulation --- precisely the tasks where the autoregressive bottleneck is most severe.

## 1.2 The Diffusion Opportunity

Diffusion models process the entire output simultaneously and refine iteratively. In image generation, a diffusion model starts from pure noise (or a fully masked canvas) and progressively recovers structure through repeated denoising steps. Each step has access to the full spatial extent of the output, enabling global coherence. Each step refines the previous, enabling iterative improvement.

This is precisely the computational pattern that spatial reasoning demands.

The insight is not that ARC grids are images --- although they are, trivially, small discrete images. The insight is that the *cognitive process* of solving ARC puzzles maps naturally onto denoising: start with uncertainty about the output grid, form an initial rough estimate, examine it for consistency with the observed examples, refine areas of low confidence, and iterate until the solution crystallizes. This is diffusion.

Three properties of ARC grids make diffusion particularly attractive:

**Small and discrete.** ARC grids are at most 30x30 cells with 10 possible colors. This is 900 categorical variables, not 786,432 continuous pixels (as in a 1024x1024 image). The denoising model can be orders of magnitude smaller than Stable Diffusion (Rombach et al., 2022).

**Natural masking.** Discrete categories admit masking-based diffusion (MaskGIT, Chang et al., 2022) rather than Gaussian noise diffusion. Masking is more natural for categorical data: a cell is either known or unknown, with no need for continuous noise schedules. Confidence-based unmasking --- revealing easy cells first, hard cells last --- mirrors human solving strategy.

**Few-shot conditioning.** ARC tasks provide 2--5 input-output example pairs that implicitly define the transformation rule. The denoiser can condition on these examples via cross-attention, learning the rule from examples and applying it to a new input --- a form of in-context learning through the diffusion process.

## 1.3 Contributions

This paper makes the following contributions:

1. **Grid Diffusion Transformer (GDiT):** The first discrete diffusion architecture designed for few-shot abstract grid reasoning, using MaskGIT-style iterative unmasking with 2D relative position bias, few-shot example encoding, and multi-signal conditioning (Section 4).
2. **Diffusion World Model:** A unified mechanism for both output generation and mental simulation, where the same GDiT architecture operates at two fidelity levels --- precise (5-step) for final output, approximate (2--3 step) for hypothetical scenario evaluation --- with denoising quality serving as an intrinsic confidence signal (Section 5).
3. **Cognitive Mesh Architecture:** A five-layer integrated system combining state engine, perception, working memory (autoregressive transformer + MoE), reasoning modes (autoregressive + GDiT + search + retrieval), and metacognition into a single differentiable architecture (Section 3).
4. **State-Conditioned Discrete Diffusion:** The first integration of a biologically-structured 72-dimension state vector with discrete diffusion, where internal state modulates denoising through FiLM conditioning, biases expert routing in the MoE denoiser, and drives metacognitive assessment after each unmasking step (Section 7).
5. **Reasoning Mode Router:** A learned routing mechanism that dynamically selects among four reasoning strategies and can switch modes mid-sequence, including pausing autoregressive generation to invoke GDiT and resuming with the result (Section 8).
6. **Dual-Scale Design:** V3-mini (2.5B parameters) as a proof-of-concept trainable on consumer hardware, and V3-full (190B parameters) as a competition-grade model deployable on a single DGX Spark at MXFP4 quantization (Sections 11, 12).
7. **Library Learning through Diffusion:** Integration of a 42-function DSL with 10 compound skills into the GDiT's denoising process, where spatial primitives are learned as denoising patterns rather than symbolic procedures (Section 6).

---

## 2. Related Work

---

### 2.1 Autoregressive Models for Abstract Reasoning

Large language models have been applied to ARC-AGI with mixed results. GPT-4 (OpenAI, 2023) and Gemini 1.5 Pro (Google, 2024) achieve approximately 50--60% on the public evaluation set with careful prompting, demonstrating that language models can reason about grid transformations when the task is presented textually. However, performance degrades sharply on tasks requiring precise spatial manipulation --- rotations, reflections, tiling --- where the 1D serialization of 2D grids introduces systematic errors.

Ryan Greenblatt's 2024 ARC-AGI submission (72%) demonstrated that wrapping an autoregressive LLM in a program synthesis and search framework substantially improves performance. The LLM generates candidate Python programs; a search procedure evaluates them against training examples and selects the best. This approach treats the LLM as a program prior rather than a direct solver, bypassing the autoregressive bottleneck for the spatial computation while retaining the LLM's strength in hypothesis generation.

MindsAI's winning 2024 system (84%) combined multiple approaches including program synthesis, search, and neural pattern matching. The system's architecture has not been fully published, but its performance establishes the current state of the art for machine ARC-AGI solving.

### 2.2 Diffusion Models

Diffusion models (Sohl-Dickstein et al., 2015; Ho et al., 2020) learn to reverse a gradual noising process, generating data by iteratively denoising from pure noise. Stable Diffusion (Rombach et al., 2022) demonstrated that diffusion in a compressed latent space scales efficiently to high-resolution image generation. DALL-E 3 (Betker et al., 2023) and Imagen (Saharia et al., 2022) showed that text-conditioned diffusion can produce photorealistic images from natural language descriptions.

The Diffusion Transformer (DiT) architecture (Peebles and Xie, 2023) replaced the U-Net backbone of earlier diffusion models with a transformer, demonstrating that the self-attention mechanism is sufficient for denoising and scales more predictably. DiT powers Stable Diffusion 3 (Esser et al., 2024) and Sora (Brooks et al., 2024). Engram V3's GDiT builds on the DiT architectural pattern but replaces continuous Gaussian denoising with discrete masking-based diffusion for categorical grid data.

### 2.3 Discrete Diffusion

Gaussian diffusion is natural for continuous data (images, audio) but awkward for categorical data. Three lines of work have developed diffusion processes for discrete spaces:

**D3PM** (Austin et al., 2021) introduced structured transition matrices for discrete diffusion, including absorbing state transitions where tokens are progressively replaced with a [MASK] token. The reverse process learns to predict the original token from the masked state. GDiT uses absorbing state diffusion as its noising process.

**MDLM** (Sahoo et al., 2024) and **SEDD** (Lou et al., 2024) advanced discrete diffusion for language modeling, demonstrating that discrete diffusion can approach autoregressive perplexity on text generation benchmarks. These

results establish that discrete diffusion is a viable alternative to autoregression for categorical sequence generation.

**MaskGIT** (Chang et al., 2022) introduced confidence-based parallel decoding for image generation: a model predicts all masked tokens simultaneously, then selectively unmask the most confident predictions, iterating until the image is complete. This approach is 8--16x faster than autoregressive image generation. GDIT adapts MaskGIT's confidence-based unmasking for grid reasoning, with the critical addition of metacognitive assessment between unmasking steps.

## 2.4 World Models

World models learn compressed representations of environment dynamics, enabling agents to "imagine" future states without executing actions in the real environment.

**Ha and Schmidhuber (2018)** introduced the World Models framework, training a VAE to compress observations and an RNN to predict future latent states. The agent's policy is trained entirely within the learned world model --- "dreaming" rather than acting.

**Dreamer** (Hafner et al., 2019; 2020; 2023) extended this to continuous control, training policies through imagined trajectories in a learned latent dynamics model. DreamerV3 (Hafner et al., 2023) demonstrated that a single world model architecture can master diverse domains from Atari to robotic control.

**JEPA** (LeCun, 2022) proposed the Joint Embedding Predictive Architecture as a path toward human-level intelligence: predict in abstract representation space rather than pixel space, discarding irrelevant details. JEPA frames intelligence as the ability to predict and plan using internal models of the world.

**MuZero** (Schrittwieser et al., 2020) demonstrated that a learned world model (without explicit rules) combined with Monte Carlo Tree Search achieves superhuman performance in chess, Go, shogi, and Atari. The model learns to predict rewards, values, and next states, enabling planning without environment simulation.

Engram V3's diffusion world model connects to this tradition by implementing mental simulation as approximate denoising: the GDIT runs at reduced fidelity (2--3 steps instead of 5) to imagine hypothetical grid outcomes, with denoising quality serving as an intrinsic confidence signal --- a blurry result indicates uncertainty, a clean result indicates confidence.

## 2.5 Library Learning and Skill Decomposition

**DreamCoder** (Ellis et al., 2021) introduced library learning for program synthesis: an agent discovers reusable abstractions (primitives) through a wake-sleep cycle, building a growing library of functions that make future problems easier to solve. Each new primitive compresses the description length of programs that use it, creating a virtuous cycle of abstraction discovery.

Hierarchical reinforcement learning (Dayan and Hinton, 1993; Sutton et al., 1999) decomposes complex tasks into sub-tasks that can be solved by specialized sub-policies. The options framework (Sutton et al., 1999) formalizes temporally extended actions (options) that can be composed hierarchically.

Engram V3 applies library learning through two mechanisms: (1) the 42-function DSL provides a fixed primitive vocabulary, while 10 compound skills provide memorized compositions analogous to chess openings; (2) the

GDiT learns spatial primitives as denoising patterns during pre-training on synthetic grid transformations, creating a neural complement to the symbolic DSL.

## 2.6 Process Reward Models and Metacognition

**Process Reward Models** (Lightman et al., 2023) demonstrated that rewarding each intermediate reasoning step --- rather than only the final answer --- produces more reliable mathematical problem-solving in language models. This step-level reward provides finer-grained training signal and catches errors earlier in the reasoning chain.

**Self-Refine** (Madaan et al., 2023) showed that LLMs can improve their outputs through iterative self-critique and refinement, without additional training. The model generates an initial output, critiques it, and revises --- multiple rounds of this loop improve quality.

Engram V3 integrates metacognition at the architectural level rather than as a post-hoc prompting strategy. After each GDiT unmasking step, dedicated metacognition experts assess the partially-revealed grid against the conditioning examples, update the confidence dimension of the 72-dim state vector, and can trigger re-masking of suspicious cells. This is process reward implemented through the state engine rather than a separate reward model.

## 2.7 Mixture of Experts

**Switch Transformer** (Fedus et al., 2022) demonstrated that sparse MoE with top-1 routing achieves better scaling than dense models at equivalent compute. **Mixtral** (Jiang et al., 2024) proved that sparse MoE can match or exceed dense models at 3--4x fewer active parameters. **DeepSeek-MoE** (Dai et al., 2024) showed that fine-grained expert segmentation with shared experts improves both efficiency and specialization.

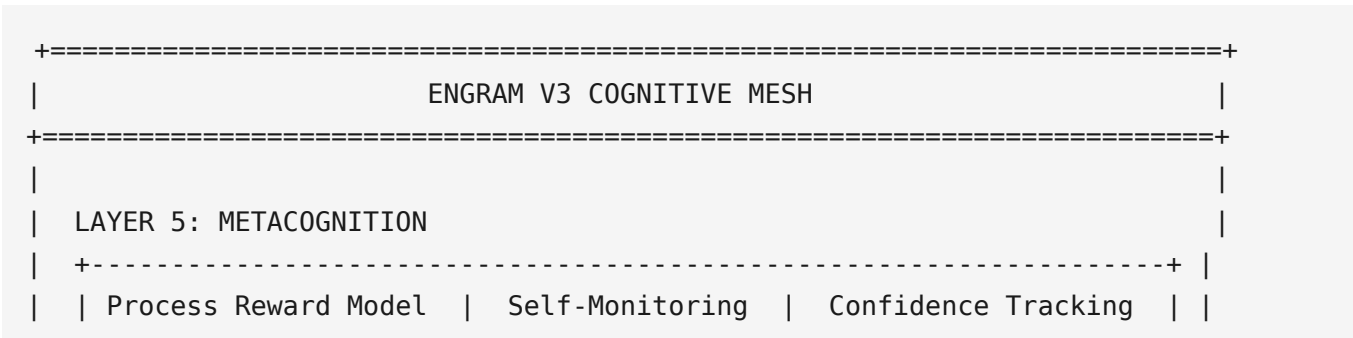
Engram V3 extends MoE in three ways: (1) state-conditioned routing where the 72-dim emotional state biases expert selection; (2) expert proximity in latent space creating functional neighborhoods; and (3) dual MoE deployment --- in the autoregressive backbone for text reasoning and in the GDiT denoiser for spatial reasoning --- with shared state conditioning across both.

---

# 3. Architecture Overview

## 3.1 Five-Layer Cognitive Mesh

Engram V3 organizes computation into five interdependent layers, inspired by the layered organization of biological neural systems. Unlike sequential pipelines, these layers operate as a mesh --- each layer can influence and be influenced by the others through the state engine.









**Layer 1: State Engine.** The 72-dimension state vector (Section 7) functions as the architecture's endocrine system --- a small, persistent signal that modulates all other layers. State dimensions range from autonomic arousal through metacognitive self-monitoring. The state vector is updated between turns based on processing content and influences computation through six pathways: FiLM modulation, state token, MoE routing bias, attention temperature, feedback loop, and state dropout.

**Layer 2: Perception.** Multi-modal input encoding converts raw inputs into the model's internal representation space. Text passes through BPE tokenization (32,768 vocabulary). Images pass through SigLIP-SO400M with PseudoDeepStack extraction (layers 9, 18, 27). ARC grids are processed both as rendered images (through SigLIP) and as structured text (through the tokenizer). A GridAnalyzer module extracts zero-parameter task features (symmetry, color histograms, object counts) for conditioning.

**Layer 3: Working Memory.** The autoregressive transformer backbone provides the primary processing substrate -- the model's working memory. With GQA attention, SwiGLU activations, and MoE feed-forward layers, this is a standard modern transformer enhanced with state-dependent modulation at every layer.

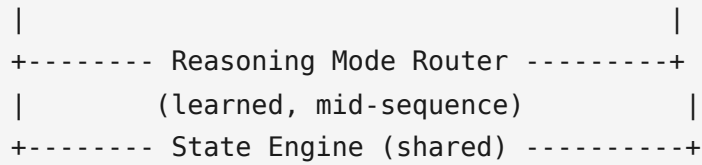
**Layer 4: Reasoning Modes.** Four distinct reasoning strategies are available, selected by a learned router (Section 8). Autoregressive generation handles text, planning, and explanation. GDiT handles spatial grid reasoning through iterative denoising. Tree search explores multiple solution branches. Memory retrieval finds similar previously-solved tasks. The router can switch between modes mid-sequence.

**Layer 5: Metacognition.** Continuous self-monitoring that operates across all reasoning modes. After each GDiT unmasking step, metacognition assesses the partial result. After each autoregressive generation segment, it checks coherence. The adversarial probe ensures state remains subliminal. Confidence tracking through the state engine provides calibrated uncertainty estimates.

3.2 Dual System Design

Engram V3 implements Kahneman's (2011) System 1 / System 2 framework at the architectural level:

SYSTEM 1 (Fast, Automatic)	SYSTEM 2 (Slow, Deliberate)
+-----+	+-----+
Autoregressive Transformer	Grid Diffusion Transformer
- Pattern matching	- Iterative refinement
- Personality/conversation	- Spatial reasoning
- Quick classification	- 5-step denoising
- Planning/orchestration	- Metacognitive checks
- Explanation generation	- World model simulation
Speed: ~2ms/token	Speed: ~80ms/grid
Left-to-right	Bidirectional
1D sequential	2D spatial
+-----+	+-----+



System 1 (the autoregressive transformer) handles fast, automatic processing: recognizing task types, generating plans ("I think this is a rotation task"), producing explanations, and managing conversational context. It operates left-to-right at approximately 2ms per token.

System 2 (the GDiT) handles slow, deliberate spatial reasoning: examining the grid holistically, iteratively refining a solution through 5 unmasking steps with metacognitive checks between each step. It operates bidirectionally over the full 2D grid at approximately 80ms per complete solution.

The reasoning mode router mediates between systems. System 1 recognizes when System 2 is needed ("this task requires spatial transformation"), invokes it, receives the result (a completed grid with confidence score), and continues autoregressive generation to explain the result. This mirrors how humans delegate spatial reasoning to a different cognitive mode than verbal reasoning.

### 3.3 Component Inventory

Component	V3-mini Params	V3-mini Active	V3-full Params	V3-full Active
Transformer backbone (SwiGLU, GQA)	1,800M	1,800M	40,000M	40,000M
MoE experts (backbone)	450M	56M	80,000M	3,300M
GDiT denoiser (MoE, 8 layers)	65M	35M	800M	300M
Example encoder (4 layers)	15M	15M	120M	120M
World model (shared GDiT)	0 (shared)	0 (shared)	0 (shared)	0 (shared)
Reasoning mode router	0.5M	0.5M	5M	5M
SigLIP-SO400M (frozen)	400M	400M	400M	400M
Vision projector (PseudoDeepStack)	7.7M	7.7M	20M	20M
State engine (FiLM + token + bias + temp + feedback)	0.9M	0.9M	4M	4M
Cell embeddings + 2D position bias	0.03M	0.03M	0.2M	0.2M
Adversarial probe (training only)	0.2M	0	1M	0
<b>Total</b>	<b>~2,540M</b>	<b>~1,615M</b>	<b>~121,350M</b>	<b>~44,150M</b>
<b>With MoE expansion (V3-full: 48 experts)</b>	---	---	<b>~190,000M</b>	<b>~48,000M</b>

Note: The world model shares the GDiT architecture and adds zero parameters. It simply runs the same denoiser at reduced fidelity (2--3 steps instead of 5). This architectural unification is a key design insight --- simulation and

generation are the same computation at different precision levels.

---

## 4. Grid Diffusion Transformer (GDiT)

---

### 4.1 MaskGIT-Style Discrete Diffusion

GDiT uses absorbing state discrete diffusion rather than Gaussian continuous diffusion. The choice is driven by three properties of ARC grids:

**Categorical data.** ARC grid cells take one of 10 discrete color values (0--9). Gaussian noise, which adds continuous perturbations, is unnatural for categorical data. Masking --- replacing cells with a special [MASK] token --- is the natural noise process for discrete spaces. A cell is either known (one of 10 colors) or unknown ([MASK]). There is no meaningful intermediate state.

**Confidence-based ordering.** MaskGIT's unmasking strategy reveals cells in order of model confidence --- the most certain predictions first, the most uncertain last. This mirrors human solving strategy: when examining an ARC puzzle, humans typically identify the easy cells first (e.g., cells that clearly follow a repeating pattern) and leave ambiguous cells for later, when more context is available from the already-determined cells.

**Fewer steps.** Gaussian diffusion typically requires 20--50 denoising steps for high-quality output. MaskGIT-style masking requires 5--8 steps. For ARC grids with at most 900 cells, 5 steps suffice: each step unmask approximately 20% of cells, with later steps benefiting from the context provided by earlier unmasked cells.

The forward (noising) process replaces a fraction of grid cells with [MASK] at each step:

```
Step 0: Full output grid (ground truth during training)
Step 1: 20% masked (random selection)
Step 2: 40% masked
Step 3: 60% masked
Step 4: 80% masked
Step 5: 100% masked (all [MASK])
```

The reverse (denoising) process predicts the original values and selectively unmask:

```
Step 5: 100% masked -> predict all -> unmask top 20% most confident
Step 4: ~80% masked -> predict remaining -> unmask next 20%
Step 3: ~60% masked -> predict remaining -> unmask next 20%
Step 2: ~40% masked -> predict remaining -> unmask next 20%
Step 1: ~20% masked -> predict remaining -> unmask final 20%
Step 0: Complete grid (0% masked)
```

The denoiser outputs a probability distribution over 11 values (10 colors + [MASK]) for each cell at each step. Confidence is defined as the maximum probability assigned to any non-MASK token. Cells with the highest confidence are unmasked first.

**Training objective.** For each training example, a random masking ratio  $r \sim \text{Uniform}(0, 1)$  is sampled.  $r$  fraction of cells are masked. The model predicts the original values of all masked cells. The loss is cross-entropy over masked positions only:

$$L_{\text{denoise}} = -\sum_{i \text{ in masked}} \log p(x_i \mid x_{\text{unmasked}}, \text{conditioning})$$

This trains the model at all masking ratios simultaneously, preparing it for every step of the iterative unmasking process.

## 4.2 2D Relative Position Bias

Standard transformers use 1D positional encodings (sinusoidal or RoPE), which destroy spatial structure when grids are flattened to sequences. GDiT uses 2D relative position bias inspired by the Swin Transformer (Liu et al., 2021):

For any two cells at positions  $(r1, c1)$  and  $(r2, c2)$ , the attention bias is:

$$\text{bias}(r1, c1, r2, c2) = \text{learned\_table}[r1 - r2 + \text{max\_size}][c1 - c2 + \text{max\_size}]$$

where `learned_table` is a trainable parameter of shape `[2*max_size - 1, 2*max_size - 1, n_heads]`.

With max grid size 30, the relative position table is 59x59 per attention head. For 8 heads, this is  $59 \times 59 \times 8 = 27,848$  parameters --- negligible.

**Why this matters.** The 2D bias means the model natively understands spatial relationships: - Horizontally adjacent cells have bias `learned_table[0][1]` - Vertically adjacent cells have bias `learned_table[1][0]` - Diagonal neighbors have bias `learned_table[1][1]` - Distant cells have their own learned biases

The model learns that spatial proximity matters differently for different heads --- some heads may attend to local neighborhoods (convolution-like), others to same-row or same-column cells (line detection), and others to distant symmetric positions (mirror detection). This emergent specialization arises from training without architectural constraints beyond the 2D bias structure.

## 4.3 Example Encoder

The example encoder is the most critical component of GDiT. ARC's core challenge is inferring a transformation rule from 2--5 input-output example pairs. The denoiser handles grid transformation; the example encoder handles rule inference.

**Architecture:** A 4-layer bidirectional transformer with 256 hidden dimension (V3-mini) or 512 hidden dimension (V3-full), processing all example pairs simultaneously.

Input:  $(in_1, out_1), (in_2, out_2), \dots, (in_k, out_k)$

Per cell embedding:



```

|   object_count: int                                     |
|   grid_shape: (H, W)                                   |
|   pattern_type: {solid, striped, checkerboard, random} |
|   -> Projected: Linear(64, 256) -> [1, 256]           |
|   -> Available via cross-attention                     |
|                                                         |
| 4. STATE VECTOR (72-dim biological state)              |
|   state [72]                                           |
|   -> FiLM modulation on every layer norm              |
|   -> MoE router bias                                  |
|   -> Attention temperature                             |
|                                                         |
+-----+

```

The denoiser integrates these signals through three mechanisms:

1. **Cross-attention** (examples + input + task features): Each denoiser layer has a cross-attention sublayer that attends to the conditioning bundle. This allows each output cell to reference any input cell, any example cell, or the task features.
2. **FiLM modulation** (state vector): The state vector modulates every layer norm in the denoiser through the same FiLM mechanism used in the autoregressive backbone (Section 7.3). The state affects *how* the denoiser processes information without being visible in the token stream.
3. **Additive bias** (step embedding): A learned embedding for the current denoising step (1--5) is added to all cell embeddings, informing the model which stage of the unmasking process it is in.

## 4.5 MoE Denoiser

The GDiT denoiser uses a mixture-of-experts architecture within each transformer layer, enabling specialization for different types of grid transformations:

Denoiser Architecture (8 layers, V3-mini):

Input: masked\_output\_grid [H\*W, 256] (starts all [MASK])

Per layer:

1. Layer Norm + FiLM modulation (state-conditioned)
2. Bidirectional self-attention
  - + 2D relative position bias
  - + state-conditioned temperature per head
3. Layer Norm + FiLM modulation
4. Cross-attention -> conditioning bundle
5. Layer Norm + FiLM modulation
6. MoE FFN:
  - 8 experts (V3-mini) / 16 experts (V3-full)
  - Top-2 routing
  - State-biased router: Linear(256, 8) + Linear(72, 8) bias

- Each expert: SwiGLU MLP (256 -> 128 -> 256)

Output: logits [H\*W, 11] (10 colors + MASK)  
-> softmax -> probabilities per cell

**State-biased routing in the denoiser.** The state vector biases which denoiser experts fire, just as it does in the autoregressive backbone. High **analytical** state biases toward systematic experts. High **creativity** biases toward exploratory experts. The routing learns to associate different state configurations with different transformation types:

State Configuration	Expert Tendency	Transformation Type
High analytical, high focus	Geometric experts	Rotations, reflections
High creativity, low analytical	Color experts	Color remapping, palette changes
High patience, high planning_horizon	Structural experts	Complex multi-step patterns
High confidence	Direct application experts	Simple, well-understood transforms
Low confidence	Exploration experts	Novel, uncertain transforms

This mapping is not hardcoded --- it emerges from training. The state bias provides the gradient pathway for specialization to develop.

## 4.6 Iterative Unmasking

The unmasking process proceeds in 5 steps with confidence-based cell selection:

### ITERATIVE UNMASKING PROCEDURE

Step 1: Start with 100% masked grid

- > Denoiser predicts:  $p(\text{color} \mid \text{all masked, conditioning})$
- > Select top 20% cells by  $\max(p)$  excluding MASK
- > Unmask selected cells with their predicted colors
- > [METACOGNITION CHECK]

Step 2: ~80% still masked

- > Denoiser predicts:  $p(\text{color} \mid \text{partially revealed, conditioning})$
- > Now has context from Step 1's high-confidence cells
- > Select top 20% of REMAINING masked cells by confidence
- > Unmask selected
- > [METACOGNITION CHECK]

Step 3: ~60% still masked

- > More context available -> better predictions
- > Unmask next 20%
- > [METACOGNITION CHECK]

Step 4: ~40% still masked

- > Rich context from 60% revealed cells
- > Unmask next 20%
- > [METACOGNITION CHECK]

Step 5: ~20% still masked (hardest cells)

- > Nearly complete grid provides strong context
- > Unmask final 20%
- > [FINAL METACOGNITION CHECK]

**Backtracking.** If the metacognition check after any step detects a confidence drop (the partial grid is becoming less consistent with the examples), suspicious cells can be *re-masked* and re-predicted in the next step. This is a capability unique to masking-based diffusion --- Gaussian diffusion has no clean mechanism for selective reversal.

The re-masking criterion uses the state engine: if `state.confidence` drops below a learned threshold after a step, cells whose individual confidence fell below a second threshold are re-masked. Both thresholds are learned during training.

## 4.7 Metacognition After Each Unmasking Step

After each of the 5 unmasking steps, the metacognition system performs three assessments:

- 1. Consistency check.** The partially-revealed output grid is compared to the expected outputs from the example pairs using the transformation rule implied by the conditioning. This is computed by the metacognition experts (experts 14--15 in the MoE denoiser), which have learned to assess grid consistency.
- 2. State update.** The state feedback MLP processes the denoiser's hidden states (mean-pooled) along with the current state vector to produce a state delta:

```
hidden_mean [256] + state [72] -> MLP(328, 128, 72) -> delta
new_state = old_state + sigmoid(rate) * tanh(delta)
```

The `confidence` dimension is particularly important: it rises when the partial grid shows consistent patterns and falls when inconsistencies appear.

- 3. Re-masking decision.** If confidence drops below a threshold, cells with individual prediction confidence below a second threshold are re-masked. This allows the model to "change its mind" about uncertain cells as more context becomes available.

Per-Step Metacognition Flow:

Unmasking step  $k$  produces partial grid  $G_k$

|  
v

Metacognition experts (14-15) assess  $G_k$

|

+--> Consistency score: how well does  $G_k$  match examples?



```

|
+--> State update: confidence, analytical, etc.
|
+--> Re-mask decision:
    if state.confidence < threshold_1:
        for each cell in G_k:
            if cell_confidence < threshold_2:
                re-mask(cell)
|
v
Proceed to unmasking step k+1

```

This creates a deliberate, self-correcting reasoning process: the model does not just generate --- it generates, assesses, revises, and repeats. Five unmasking steps with metacognitive checks produce a level of reasoning rigor impossible in single-pass autoregressive generation.

---

## 5. Diffusion World Model

---

### 5.1 Mental Simulation via Denoising

The diffusion world model implements Vincent's insight about mental simulation: *"humans are simulating things in their head, like daydreaming or playing it through, constantly and some part in them says 'nah that's bs, that won't be like that' or 'ok this could really be it.'"*

The key realization is that the GDiT denoiser already knows how grid transformations work --- it has been trained to reverse the masking process, which requires understanding what valid transformed grids look like. This same knowledge can be repurposed for *simulation*: given a hypothetical action sequence, what would the result look like?

WORLD MODEL SIMULATION (approximate, fast)

```

Input: current_grid + hypothetical_action ("rotate 90 then tile")
|
v
Encode action as conditioning signal (appended to conditioning bundle)
|
v
Run GDiT denoiser at REDUCED FIDELITY:
- 2-3 unmasking steps instead of 5
- No metacognition checks between steps
- Larger unmasking chunks (33-50% per step instead of 20%)
|
v
Output: approximate_result_grid + denoising_quality_score
|

```

v

#### Quality Assessment:

- Clean, coherent grid -> high confidence in this plan
- Blurry, inconsistent grid -> low confidence, try another plan

The simulation is intentionally imprecise. Like a human mentally rotating a grid and getting a rough sense of the result, the world model produces an approximate answer that is good enough to distinguish promising approaches from hopeless ones. Full precision is reserved for the actual output generation (5-step GDiT).

## 5.2 Unified Architecture

The world model uses the *same* GDiT architecture and weights as the output generator. There are no additional parameters. The only difference is operational:

Property	Output Generator	World Model
Denoising steps	5	2--3
Unmasking per step	20%	33--50%
Metacognition checks	Yes, after each step	No
Backtracking	Yes (re-masking)	No
Output quality	Precise	Approximate
Latency	~80ms	~20--30ms
Purpose	Final answer	Scenario evaluation

This unification is architecturally elegant: simulation and generation are the same computation at different fidelity levels. The GDiT learns to generate valid grids during training; running it at reduced fidelity produces valid-but-approximate grids during simulation. No separate world model needs to be trained, maintained, or synchronized.

## 5.3 Confidence from Denoising Quality

The most elegant property of the diffusion world model is that confidence is *intrinsic* to the process. No separate critic or value network is needed.

When the world model simulates a scenario with high confidence, the denoising process produces a clean, coherent grid --- cell probabilities are peaked (high max probability), colors are spatially consistent, and the grid exhibits structural regularity. When the scenario is uncertain or poorly suited, the denoising process produces a blurry, inconsistent result --- cell probabilities are dispersed, colors are spatially incoherent, and the grid lacks structure.

The confidence signal is computed as:

```
confidence = mean(max(p_cell) for cell in all_cells)
```

where `p_cell` is the probability distribution over colors for each cell at the final denoising step. Values near 1.0 indicate high confidence (the model is sure of each cell). Values near 0.1 (uniform over 10 colors) indicate uncertainty.

This eliminates the need for a separate critic network (as in actor-critic RL) or a learned value function (as in MuZero). The denoiser is the critic. This is analogous to how humans assess the quality of a mental simulation: a vivid, detailed mental image feels confident; a vague, blurry one feels uncertain.

## 5.4 Imagination-Based Planning

The world model enables a planning loop where multiple hypothetical approaches are simulated before any is committed to:

### IMAGINATION-BASED PLANNING

1. System 1 (autoregressive) analyzes the task:  
"This looks like it could be rotation, tiling, or color remapping"
2. World model simulates each hypothesis (parallel, ~30ms each):  
Hypothesis A: rotate\_90 -> confidence 0.82  
Hypothesis B: tile -> confidence 0.34  
Hypothesis C: color\_map -> confidence 0.71
3. Rank by confidence:  
A (0.82) > C (0.71) > B (0.34)
4. Execute top 2-3 hypotheses through full GDiT (5-step, precise):  
A: rotate\_90 -> full GDiT -> grid\_A (precise)  
C: color\_map -> full GDiT -> grid\_C (precise)
5. Metacognition compares grid\_A and grid\_C to examples:  
grid\_A: consistency\_score 0.91  
grid\_C: consistency\_score 0.65
6. Submit grid\_A as final answer.

This is analogous to how a chess engine (MuZero, AlphaGo) evaluates candidate moves: simulate many, evaluate the results, commit to the best. But where MuZero uses a learned value function, we use denoising quality --- a natural byproduct of the generation process.

The planning loop simulates up to 50 hypotheses (each at ~30ms = ~1.5s total) and executes 2--3 through full GDiT (~80ms each). Total reasoning time: approximately 2 seconds per task. This is 10--100x faster than search-based LLM approaches that may generate dozens of full programs.

## 5.5 Connection to JEPA and World Model Research

Engram V3's diffusion world model aligns with Yann LeCun's JEPA framework (LeCun, 2022) in several key ways:

**Prediction in abstract space.** The GDiT operates on grid-cell embeddings (256-dim vectors), not raw pixel values. The world model simulates in this abstract space, discarding irrelevant visual details (exact pixel colors, anti-aliasing artifacts) and preserving task-relevant structure (cell identities, spatial relationships).

**Hierarchical planning.** The two-fidelity approach (approximate simulation for filtering, precise generation for execution) implements a form of hierarchical planning. Broad search happens at low fidelity; focused execution at high fidelity.

**Self-supervised learning.** The GDiT is trained on a self-supervised objective (predict masked cells). The world model inherits this capability without additional supervised labels. This is JEPA's key insight: learn to predict, and planning emerges.

The diffusion world model also connects to: - **Dreamer** (Hafner et al., 2023): training policies in imagination. Engram simulates grid transformations in imagination before committing. - **MuZero** (Schrittwieser et al., 2020): learned dynamics without rules. The GDiT learns grid transformation dynamics from data, not from programmed rules. - **Diffuser** (Janner et al., 2022): planning as denoising trajectories. Engram generates grid transformation outcomes through denoising. - **UniSim** (Yang et al., 2023): universal simulation via diffusion. Engram's world model simulates grid transformations as a denoising process.

---

## 6. Skill Decomposition and Library Learning

---

### 6.1 42-Function DSL Taxonomy

The domain-specific language provides 42 pure Python functions organized into six clusters, forming the primitive vocabulary from which compound transformations are composed. Every function operates on numpy 2D integer arrays (values 0--9), is deterministic, pure (no side effects), and composable.

#### Cluster 1: Geometric Transforms (6 functions)

Function	Signature	Description
<code>rotate_90</code>	<code>grid -&gt; grid</code>	Rotate 90 degrees clockwise
<code>rotate_180</code>	<code>grid -&gt; grid</code>	Rotate 180 degrees
<code>rotate_270</code>	<code>grid -&gt; grid</code>	Rotate 270 degrees clockwise
<code>flip_horizontal</code>	<code>grid -&gt; grid</code>	Mirror left-right
<code>flip_vertical</code>	<code>grid -&gt; grid</code>	Mirror top-bottom
<code>transpose</code>	<code>grid -&gt; grid</code>	Swap rows and columns

#### Cluster 2: Color Operations (7 functions)

Function	Signature	Description
<code>color_replace</code>	<code>grid, old, new -&gt; grid</code>	Replace all cells of one color with another
<code>color_swap</code>	<code>grid, c1, c2 -&gt; grid</code>	Swap two colors throughout grid
<code>color_fill</code>	<code>grid, color -&gt; grid</code>	Fill entire grid with one color
<code>color_map</code>	<code>grid, mapping -&gt; grid</code>	Apply color remapping dictionary
<code>most_common_color</code>	<code>grid -&gt; int</code>	Return the most frequent color
<code>least_common_color</code>	<code>grid -&gt; int</code>	Return the least frequent non-background color
<code>count_color</code>	<code>grid, color -&gt; int</code>	Count cells of a specific color

### Cluster 3: Object Detection (7 functions)

Function	Signature	Description
<code>find_objects</code>	<code>grid -&gt; list[Object]</code>	Connected component detection
<code>find_rectangles</code>	<code>grid -&gt; list[Rect]</code>	Find rectangular regions
<code>find_lines</code>	<code>grid -&gt; list[Line]</code>	Find horizontal/vertical lines
<code>find_borders</code>	<code>grid -&gt; grid</code>	Identify border cells
<code>find_background</code>	<code>grid -&gt; int</code>	Identify background color
<code>extract_object_grid</code>	<code>grid, obj -&gt; grid</code>	Extract sub-grid containing object
<code>count_objects</code>	<code>grid -&gt; int</code>	Count distinct connected components

### Cluster 4: Structural Analysis (4 functions)

Function	Signature	Description
<code>detect_symmetry</code>	<code>grid -&gt; SymmetryType</code>	Horizontal, vertical, diagonal, rotational, none
<code>find_repeating_pattern</code>	<code>grid -&gt; Pattern</code>	Identify minimal repeating unit
<code>detect_grid_structure</code>	<code>grid -&gt; GridStruct</code>	Find sub-grid organization
<code>crop_to_content</code>	<code>grid -&gt; grid</code>	Remove background borders

### Cluster 5: Spatial Operations (11 functions)

Function	Signature	Description
<code>crop</code>	<code>grid, r1, c1, r2, c2 -&gt; grid</code>	Extract rectangular region
<code>pad</code>	<code>grid, n, color -&gt; grid</code>	Add border of given color
<code>resize</code>	<code>grid, h, w -&gt; grid</code>	Scale grid to new dimensions

Function	Signature	Description
<code>tile</code>	<code>grid, h_reps, w_reps -&gt; grid</code>	Repeat grid in tiling pattern
<code>place_subgrid</code>	<code>grid, sub, r, c -&gt; grid</code>	Place sub-grid at position
<code>overlay</code>	<code>grid1, grid2 -&gt; grid</code>	Combine two grids (non-zero wins)
<code>concat_horizontal</code>	<code>g1, g2 -&gt; grid</code>	Concatenate side by side
<code>concat_vertical</code>	<code>g1, g2 -&gt; grid</code>	Concatenate top to bottom
<code>draw_line</code>	<code>grid, r1, c1, r2, c2, color -&gt; grid</code>	Draw straight line
<code>draw_rectangle</code>	<code>grid, r1, c1, r2, c2, color -&gt; grid</code>	Draw rectangle outline
<code>flood_fill</code>	<code>grid, r, c, color -&gt; grid</code>	Fill connected region

#### Cluster 6: Comparison (5 functions)

Function	Signature	Description
<code>grids_equal</code>	<code>g1, g2 -&gt; bool</code>	Exact cell-by-cell comparison
<code>diff_grids</code>	<code>g1, g2 -&gt; grid</code>	Produce grid showing differences
<code>similarity_score</code>	<code>g1, g2 -&gt; float</code>	Fraction of matching cells
<code>get_shape</code>	<code>grid -&gt; (int, int)</code>	Return (height, width)
<code>get_colors</code>	<code>grid -&gt; set[int]</code>	Return set of colors present

**DSL coverage analysis.** Approximately 60% of ARC-AGI public tasks (240/400) can be solved using compositions of these 42 functions. Approximately 25% require operations not yet implemented (conditional branching, counting-based logic, complex path finding). Approximately 15% are genuinely novel or adversarial, resisting any fixed DSL. The theoretical ceiling with this DSL is therefore approximately 240/400 (60%).

## 6.2 10 Compound Skills ("Chess Openings")

Compound skills are memorized 2--3 step compositions that recur across many ARC tasks. They are the "chess openings" of grid transformation --- well-studied patterns that a skilled solver applies without deliberation.

#	Compound Skill	Steps	ARC Pattern	Frequency
1	<b>Extract &amp; Transform</b>	<code>find_objects -&gt; extract_object_grid</code> <code>-&gt; rotate/flip</code>	Object manipulation	~15% of tasks
2	<b>Pattern Tile</b>	<code>crop_to_content -&gt; tile</code>	Repeating grid completion	~8% of tasks

#	Compound Skill	Steps	ARC Pattern	Frequency
3	Color Remap	<code>find_objects</code> -> <code>color_map</code> by region	Conditional recoloring	~12% of tasks
4	Symmetry Complete	<code>detect_symmetry</code> -> <code>flip</code> -> <code>overlay</code>	Fill missing half	~7% of tasks
5	Border Extract	<code>find_borders</code> -> crop interior	Frame removal	~5% of tasks
6	Object Sort	<code>find_objects</code> -> sort by size/color -> <code>place_subgrid</code>	Spatial ordering	~6% of tasks
7	Gravity Drop	<code>find_objects</code> -> sort by row -> <code>place_subgrid</code> bottom-up	Gravity simulation	~4% of tasks
8	Scale & Fill	<code>crop_to_content</code> -> <code>resize</code> -> <code>color_replace</code> background	Scale transform	~5% of tasks
9	Mask & Paint	<code>diff_grids</code> -> <code>apply_mask</code> -> <code>color_fill</code>	Difference-based recolor	~6% of tasks
10	Grid Decompose	<code>detect_grid_structure</code> -> crop cells -> transform per cell	Sub-grid operations	~8% of tasks

These compound skills are not hardcoded. They are discovered during training on ARC traces and emerge as preferred expert compositions in the MoE routing patterns.

### 6.3 Skill-to-Expert Mapping

Each skill cluster maps to a group of MoE experts in both the autoregressive backbone and the GDiT denoiser, creating functional specialization:

SKILL-TO-EXPERT MAPPING (V3-mini: 16 experts per layer)

Experts 0-2: GEOMETRIC SPECIALISTS  
Rotations, flips, transposes  
Training: 52+ traces with geometric ops

Experts 3-5: COLOR SPECIALISTS  
Color replace, swap, map, fill  
Training: 100+ traces with color ops

Experts 6-8: OBJECT SPECIALISTS  
Find, extract, count objects  
Training: 163+ traces with object detection

Experts 9-11: SPATIAL SPECIALISTS  
Crop, tile, place, overlay  
Training: 99+ traces with spatial ops

Experts 12-13: STRUCTURAL ANALYSTS

Symmetry, patterns, grid structure

Training: 28+ traces with structural analysis

Experts 14-15: METACOGNITION / COMPOSITION

Whole-board verification, multi-step planning

Training: Multi-turn retry traces

In V3-full (48 experts per layer), the mapping scales proportionally:

Expert Range	Skill Cluster	Count
0--7	Geometric	8
8--15	Color	8
16--23	Object	8
24--31	Spatial	8
32--37	Structural	6
38--43	Composition	6
44--47	Metacognition	4

## 6.4 MoE Skill Specialization Training

Expert specialization is trained through three mechanisms:

**Skill-clustered batches (Phase 4b).** During skill SFT, training batches are constructed to emphasize one skill cluster at a time. A "geometric batch" contains traces dominated by rotation and flip operations; a "color batch" emphasizes color remapping traces. The router learns to send geometric tokens to geometric experts because those experts are most experienced with that data distribution.

**GRPO reinforcement (Phase 4c).** The router receives reward signal based on whether its expert composition produced the correct output grid. If routing to experts 0-2 (geometric) followed by experts 9-11 (spatial) produces the correct answer for an "extract and transform" task, those routing decisions are reinforced. The router learns compound skill routing --- which expert sequences solve which task types.

**State-conditioned routing.** The 72-dim state vector provides an additional routing signal. When the metacognition system indicates high confidence in a geometric hypothesis (from the world model simulation), the state biases the router toward geometric experts. When confidence is low and the system is exploring alternatives, the state biases toward broader expert activation.

---

## 7. State Engine and Affective Computing

---



## 7.1 72-Dimension State Vector

The state engine is inherited from Engram V2 and extended with GDiT-specific conditioning pathways. The 72 dimensions are organized into 8 biological layers, ordered from primitive (autonomic) to abstract (metacognition):

Layer 8: Identity & Metacognition	(9 dims)	"Who am I? How am I performing?"
Layer 7: Self-Regulation	(7 dims)	"How should I express this?"
Layer 6: Social Dynamics	(10 dims)	"How do I relate to this person?"
Layer 5: Motivational Drives	(9 dims)	"What drives my behavior?"
Layer 4: Cognitive Modes	(10 dims)	"How am I thinking right now?"
Layer 3: Complex / Social Emotions	(12 dims)	"What complex feelings are active?"
Layer 2: Core Emotions (Plutchik)	(8 dims)	"What basic emotions are active?"
Layer 1: Autonomic / Physiological	(7 dims)	"What is my baseline processing state?"
	-----	
	72 total	

**Layer 1: Autonomic (7 dims).** arousal, energy, stress, comfort, restfulness, sensory\_load, hunger. The processing baseline --- high arousal increases overall activation, low energy reduces processing depth.

**Layer 2: Core Emotions (8 dims).** valence, joy, sadness, anger, fear, disgust, surprise, anticipation. Plutchik's (1980) eight primary emotions. These combine: joy + anticipation = optimism; fear + surprise = alarm.

**Layer 3: Complex Emotions (12 dims).** trust, love, guilt, shame, pride, gratitude, awe, nostalgia, hope, contempt, envy, loneliness. Higher-order emotions requiring social context or self-awareness.

**Layer 4: Cognitive Modes (10 dims).** curiosity, focus, creativity, analytical, confidence, skepticism, patience, decisiveness, mental\_load, planning\_horizon. How the model thinks, not what it feels. These are critical for GDiT operation: high confidence reduces unmasking hesitation; high skepticism increases re-masking likelihood; high planning\_horizon favors world model simulation before commitment.

**Layer 5: Motivational Drives (9 dims).** reward\_seeking, pain\_avoidance, novelty\_seeking, achievement, mastery, urgency, perfectionism, motivation, indulgence.

**Layer 6: Social Dynamics (10 dims).** attachment, empathy, care, protectiveness, assertiveness, deference, playfulness, cooperation, teaching, belonging.

**Layer 7: Self-Regulation (7 dims).** `impulse_control`, `emotional_regulation`, `frustration_tolerance`, `adaptability`, `escalation`, `risk_tolerance`, `hedging`. The brake pedal --- modulates everything above it.

**Layer 8: Metacognition (9 dims).** `self_coherence`, `temporal_awareness`, `autonomy`, `growth`, `agency`, `self_esteem`, `identity_stability`, `purpose`, `self_monitoring`. The self-model. `self_monitoring` directly feeds the metacognitive verification loop in GDiT.

## 7.2 Six State-to-Model Pathways

The state vector influences computation through six distinct mechanisms, each operating at a different level of the processing pipeline. These pathways apply to both the autoregressive backbone and the GDiT denoiser:

**Pathway 1: FiLM Generator (Subliminal Modulation).** ~322K params (V3-mini). The state is encoded through a 3-layer MLP and projected to per-modulation-point scalars that multiply learned feature-direction patterns. Every layer norm in every transformer layer (backbone and GDiT) is modulated. The model does not "see" this modulation --- it is subliminal, like mood affecting perception.

```
State [72] -> Encoder(72,256,256,256) -> to_gamma_scale(256,48) + to_beta_scale(256,48)
gammas = 1.0 + g_scales * gamma_patterns[48, n_embd]
betas = b_scales * beta_patterns[48, n_embd]
normed = gammas * LayerNorm(x) + betas
```

**Pathway 2: State Token (Introspective).** ~115K params. The state vector is projected to a single token embedding prepended to the input sequence. Unlike FiLM, this is visible to the model --- it can attend to its own state. This enables introspection: the model literally looking at its own emotional state.

**Pathway 3: MoE State Bias (Emotional Routing).** ~28K params. A linear projection from the state vector biases router logits in every MoE layer (both backbone and GDiT). Different emotional states route to different experts. High `analytical` routes to systematic experts; high `creativity` routes to divergent experts.

**Pathway 4: Attention Temperature (Focus Control).** ~21K params. The state vector is projected to per-head temperature scalars: `sigmoid(Linear(state)) * 1.5 + 0.5`, giving range [0.5, 2.0]. Low temperature (high confidence) produces sharp, focused attention. High temperature (uncertainty) produces broad, scanning attention. This implements norepinephrine-modulated attentional gain.

**Pathway 5: State Feedback Loop (Bidirectional).** ~175K params. After each conversation turn (backbone) or each unmasking step (GDiT), mean-pooled hidden states update the state vector through a small MLP:

```
hidden_mean + current_state -> MLP(n_embd+72, 128, 72) -> delta
new_state = old_state + sigmoid(rate) * tanh(delta)
```

Processing content changes the state. Solving a puzzle increases confidence. Encountering inconsistencies increases skepticism.

**Pathway 6: State Dropout (Robust Fallback).** 0 params. 15% of training samples receive a zero state vector. This forces the model to function without emotional context, making it robust if state injection fails at inference.

## 7.3 State-Conditioned Denoising

In V3, the state engine gains a new role: conditioning the GDiT denoiser. This is the first integration of biologically-structured emotional state with discrete diffusion.

The state influences GDiT through all six pathways simultaneously:

### STATE-CONDITIONED DENOISING FLOW

72-dim state vector

```
|
+---> FiLM on GDiT layer norms (subliminal modulation)
|       How the denoiser processes grid features
|
+---> MoE routing bias in GDiT (expert selection)
|       Which denoiser experts fire
|
+---> Attention temperature in GDiT (focus control)
|       How sharply the denoiser attends
|
+---> State feedback after each unmasking step
|       Denoising results update the state
|
+---> Re-masking threshold
|       Low confidence triggers backtracking
```

**Why state matters for denoising.** Consider two scenarios:

*High confidence, high analytical:* The model is fairly sure this is a rotation task. FiLM shifts feature processing toward geometric patterns. MoE routing biases toward geometric experts (0--2). Attention temperature is low (sharp focus on the specific transformation). The denoiser confidently unmaskes cells.

*Low confidence, high curiosity:* The model is uncertain about the transformation type. FiLM maintains broad feature sensitivity. MoE routing distributes across multiple expert groups. Attention temperature is high (broad scanning across the grid). The denoiser unmaskes conservatively, and the world model runs additional simulations before committing.

This is not hand-designed behavior. It emerges from training. The state bias provides the gradient pathway; the training objective (correct grid prediction) provides the signal; the model discovers that different states benefit from different denoising strategies.

## 7.4 Adversarial Probe

The adversarial probe (~200K params, training only) prevents the state from leaking into output text:

```
mid_layer_activations -> gradient_reversal -> probe_MLP -> predicted_state
adversarial_loss = MSE(predicted_state, true_state)
```

The gradient reversal layer (Ganin and Lempitsky, 2015) causes the main model to be trained to *hide* state from the probe, even as the probe tries to *extract* it. This ensures FiLM modulation operates through modulation (changing how features are processed) rather than injection (adding state information to the feature stream).

The probe applies to both the autoregressive backbone and the GDiT. In the GDiT, this ensures that state-conditioned denoising operates subliminally --- the model's grid predictions are influenced by state but the state cannot be reconstructed from the grid cells.

---

## 8. Reasoning Mode Router

---

### 8.1 Four Reasoning Modes

Engram V3 provides four distinct reasoning strategies, each suited to different aspects of problem-solving:

#### MODE 1: AUTOREGRESSIVE (System 1)

- Standard left-to-right text generation
- For: planning, explanation, conversation, classification
- Speed: ~2ms/token
- When: default mode; handles all text-based reasoning

#### MODE 2: GDiT DIFFUSION (System 2)

- Iterative grid denoising with metacognition
- For: spatial grid reasoning, 2D transformations
- Speed: ~80ms/grid
- When: task requires spatial manipulation

#### MODE 3: TREE SEARCH

- Generate multiple candidates, evaluate, select best
- For: uncertain tasks where multiple approaches are viable
- Speed: ~500ms (5 candidates x ~100ms each)
- When: low confidence, multiple hypotheses

#### MODE 4: MEMORY RETRIEVAL

- Search SoulKeeper for similar previously-solved tasks
- For: tasks resembling past problems
- Speed: ~10ms (embedding + cosine search)
- When: task triggers associative memory

### 8.2 Learned Selection

The reasoning mode router is a small learned network that selects the appropriate mode given the current context:

## Router Architecture:

```
task_embedding [n_embd]    (from backbone's representation of the task)
+ state_vector [72]        (current internal state)
+ task_features [64]       (from GridAnalyzer: symmetry, shapes, colors)
|
v
concat -> Linear(n_embd+72+64, 256) -> GELU
      -> Linear(256, 128) -> GELU
      -> Linear(128, 4) -> softmax
|
v
mode_probabilities [4]: [p_auto, p_gdit, p_search, p_retrieval]
```

The router is trained end-to-end: mode selection affects downstream performance, and the reward signal (correct grid, personality consistency) backpropagates through the routing decision. The router learns, for example:

- Small grids with clear geometric patterns -> GDiT (high p\_gdit)
- Text-heavy tasks requiring explanation -> Autoregressive (high p\_auto)
- Novel tasks with no clear pattern -> Search (high p\_search)
- Tasks similar to previously-seen problems -> Retrieval (high p\_retrieval)

## 8.3 Mid-Sequence Mode Switching

The router can switch modes during a generation sequence. This is critical for tasks that require multiple reasoning strategies:

### MID-SEQUENCE MODE SWITCH EXAMPLE

[Autoregressive mode]

"Looking at this ARC task, I see 3 example pairs. The input grids have objects in different positions, and the output grids show them rotated. Let me solve this with spatial reasoning."

[Router: switch to GDiT mode]

<spatial\_reasoning>

GDiT processes input grid with example conditioning

5-step unmasking with metacognition

Produces output grid with confidence 0.87

</spatial\_reasoning>

[Router: switch back to autoregressive mode]

"The transformation is a 90-degree clockwise rotation. Each object in the input has been rotated around its center. Confidence: 87%."

The mode switch is implemented as special tokens in the sequence:

- `<spatial_start>` : Pause autoregressive generation, invoke GDiT
- `<spatial_end>` : GDiT complete, resume autoregressive generation
- `<search_start>` / `<search_end>` : Tree search mode
- `<retrieve_start>` / `<retrieve_end>` : Memory retrieval mode

These tokens are generated by the autoregressive model and consumed by the mode dispatch system. The GDiT's output (completed grid + confidence) is encoded back into the token stream as a special grid representation that the autoregressive model can reference for explanation generation.

## 8.4 GDiT Integration Protocol

The integration between the autoregressive backbone and GDiT follows a precise protocol:

### GDiT INTEGRATION PROTOCOL

#### 1. RECOGNITION

Autoregressive backbone generates `<spatial_start>` token  
(learned: router probability > threshold triggers this token)

#### 2. CONTEXT TRANSFER

- Current hidden state snapshot -> state vector for GDiT
- Task description from backbone -> conditioning bundle
- Example pairs already encoded -> example encoder
- Backbone pauses generation

#### 3. GDiT EXECUTION

- 5-step iterative unmasking
- Metacognition checks between steps
- World model simulation if confidence uncertain
- Produces: `output_grid [H, W]` + confidence [scalar]

#### 4. RESULT INJECTION

- Output grid encoded as grid tokens
- Confidence encoded as scalar embedding
- Injected into backbone's token stream after `<spatial_start>`
- Backbone generates `<spatial_end>` and continues

#### 5. EXPLANATION

Backbone references the injected grid result to generate natural language explanation of the transformation

This protocol ensures that GDiT is a *tool* that the autoregressive backbone invokes, not a replacement for it. The backbone handles orchestration, planning, and communication; GDiT handles spatial computation. Neither could solve ARC alone --- together, they cover both linguistic and spatial reasoning.

---

## 9. Memory System

---

### 9.1 SoulKeeper External Memory

The SoulKeeper memory system (detailed in the Engram V2 whitepaper) provides persistent associative recall across sessions. Its seven components --- episodic memory store, memory daemon (localhost:5151), IDE hook, reflexive associative recall, distilled memory pipeline, MCP server (10 tools), and auto-tuning feedback loop --- carry forward unchanged into V3.

Key specifications: - 771+ structured memories with semantic embeddings (384-dim MiniLM-L6-v2) - ONNX inference daemon: 115 MB RSS, 100% CPU, zero GPU usage - Auto-starts on boot via systemd user service - Reflexive recall: associations fire automatically, triggered by context, not explicit queries - State-influenced recall: emotional state biases which memories surface

SoulKeeper serves as the external memory system for both V3-mini and V3-full. In V3-mini, which runs on a DGX Spark with limited context window, external memory is essential --- it provides access to a lifetime of interactions that cannot fit in the model's working memory. In V3-full, external memory supplements the larger context window with cross-session persistence.

### 9.2 Internalized Memory Bank (V3-full)

V3-full introduces an internalized memory bank --- a persistent key-value attention mechanism that allows the model to store and retrieve information within its own parameters, rather than relying solely on the external SoulKeeper system.

#### INTERNALIZED MEMORY BANK

Memory slots: 1024 persistent KV pairs  
Key dimension: 128 per head (same as GQA head dim)  
Value dimension: 128 per head

Writing (after each turn):

- Select tokens with highest "novelty" score (low similarity to existing keys)
- Write new KV pairs to memory bank
- Exponential decay on old entries (staleness weighting)
- Salience-weighted: emotionally significant tokens get priority

Reading (every attention layer):

- Standard QKV attention includes memory bank keys/values
- Memory entries participate in softmax alongside context tokens
- No additional mechanism needed -- just extra KV entries

Memory Bank Attention:

Q = query from current token  
K = [context\_keys | memory\_bank\_keys]  
V = [context\_values | memory\_bank\_values]

```
attn = softmax(Q @ K^T / sqrt(d))
output = attn @ V
```

The internalized memory bank is V3-full only. V3-mini relies exclusively on SoulKeeper external memory. At 1024 slots x 128 dimensions x 2 (KV) x 4 KV heads = 1 MB, the memory bank is negligible in parameter count but provides a persistent information store that survives context window boundaries.

### 9.3 Emotional Salience-Weighted Memory Writing

Memory writing is not uniform --- emotionally significant moments receive priority. The state engine determines writing priority:

State Condition	Writing Boost	Rationale
<code>surprise</code> > 0.7	+0.3 priority	Novel information is worth remembering
<code>joy</code> > 0.8 or <code>sadness</code> > 0.8	+0.2 priority	Emotionally intense moments anchor memory
<code>confidence</code> change > 0.3	+0.2 priority	Moments of insight or confusion
<code>trust</code> change > 0.2	+0.15 priority	Relationship-defining moments
<code>fear</code> > 0.6	+0.25 priority	Threats and failures should be remembered

This mirrors the biological mechanism: the amygdala tags emotionally significant events for enhanced hippocampal encoding. High-arousal events (positive or negative) produce stronger, more persistent memories than neutral events. Engram implements this through a learned priority function conditioned on the state vector.

### 9.4 Retrieval as Reasoning Mode

Memory retrieval is not just a background process --- it is a full reasoning mode accessible through the mode router. When the router selects retrieval mode, the system:

1. Embeds the current task using MiniLM-L6-v2 (384-dim)
2. Searches the SoulKeeper memory store for similar previously-solved tasks
3. If matches are found (similarity > 0.6), retrieves the solution trace
4. Injects the retrieved trace as conditioning context for GDiT or autoregressive reasoning
5. The model can adapt a past solution to the current task rather than solving from scratch

This implements case-based reasoning: recognizing that a new problem resembles a previously-solved one and adapting the old solution. For ARC tasks, many puzzles share structural similarities (e.g., multiple tasks involve "extract object and rotate"). Retrieval allows the model to build on past experience rather than re-deriving solutions.

---

## 10. Level of Detail Context Management

### 10.1 Beyond Progressive Disclosure



Standard RAG retrieval operates in binary: the top-K items are returned at full fidelity; everything else receives nothing. Our existing ContextEngine (Section 9) uses progressive disclosure with discrete tiers - load tier 1 always, tier 2 if relevant, tier 3 on failure. This is better than binary but fundamentally discrete.

We introduce LoD (Level of Detail) Retrieval, inspired by real-time graphics rendering where distant objects receive fewer triangles proportional to their visual importance. In LoD Retrieval, **every item in the context competes for tokens proportional to its relevance score**, within a fixed total budget. The fidelity is a continuous function of relevance distance, not a tier assignment.

### 10.2 Continuous Budget Allocation

```
def allocate_tokens(items, total_budget, context):
    scores = [(item, relevance(item, context)) for item in items]
    visible = [(item, score) for item, score in scores if score > 0.2]
    total_score = sum(s for _, s in visible)
    for item, score in visible:
        token_budget = int((score / total_score) * total_budget)
        yield compress_to_budget(item, token_budget)
```

Instead of 3 full memories at 500 tokens each (1,500 tokens, 3 items), LoD retrieval produces 3 full + 5 summaries + 10 one-liners within the same 1,500-token budget - **3x the coverage at the same cost**.

### 10.3 LoD Hierarchy for the Cognitive Mesh

Level	Detail	Analogy	Token Cost
L0	State vector (72-dim)	Skybox	0 (always present via FiLM)
L1	Embedding clusters (384-dim)	Distant terrain	~5--10 tokens/item
L2	Distilled summaries	Nearby buildings	~30--80 tokens/item
L3	Full episodes	The room you're in	~200--500 tokens/item

Retrieval cascade: L1 match → promote to L2 → if confidence needs it → load L3. Sub-100ms for L1--L2, ~500ms for L3. Dynamic promotion: if the model is mid-reasoning and an L1 match becomes critical, stream L2/L3 data in real-time - like texture streaming in modern game engines.

### 10.4 Applied to Every Prompt Component

LoD applies not only to memory but to every component in the context window:

- **DSL tools:** Score > 0.8 → full docstring + examples; 0.4--0.8 → signature + one-line; 0.2--0.4 → name only; <0.2 → not loaded
- **Grid regions:** Active quadrant at full resolution; other quadrants RLE-compressed; unchanged regions as "same as input"
- **Attempt history:** Most recent at full detail; 2--3 back as key insight + outcome; older as one-line deltas

## 10.5 Connection to Published Work

LoD Retrieval draws on Matryoshka Representation Learning (Kusupati et al., 2022) for variable-resolution embeddings, Infini-Attention (Munkhdalai et al., 2024) for fixed-budget long-context management, and Complementary Learning Systems theory (McClelland et al., 1995) for the biological basis of progressive memory consolidation. The novelty is the unified framework treating token allocation as a continuous optimization analogous to GPU triangle budgeting.

# 11. Neural Fabric Protocol (HULLHB)

## 11.1 The Missing White Matter

The Cognitive Mesh as described in preceding sections defines the nuclei (gray matter) - specialized modules for perception, reasoning, memory, and metacognition. What is missing is the white matter: the learned communication pathways connecting them. In the human brain, white matter accounts for approximately 45% of brain volume and consists of axonal pathways with variable bandwidth (myelination) and learned routing.

We introduce HULLHB (the neural fabric protocol), a trainable communication layer operating at two scales: intra-model (expert-to-expert within a single Cognitive Mesh) and inter-instance (across multiple Cognitive Mesh instances on distributed hardware).

## 11.2 Protocol Stack

Layer 4: SEMANTIC PROTOCOL (TRAINED)

└─ Signal encoding:

Learned codebook compression

└─ Segment routing:

SRv6-inspired self-routing segment lists

└─ Bandwidth adaptation:

LoD-based variable resolution (Section 10)

└─ Myelination:

Frequently-used pathways get wider bandwidth

└─ Timing:

Send when signal is "ripe," not every step

Layer 3: ABSTRACTION (FIXED)

└─ Nucleus discovery:

SLAAC-style self-registration

└─ Capability ads:

NDP-style "I do spatial reasoning"

└─ State sync:

72-dim state vector broadcast

└─ Loose routing:

Nuclei can modify segment lists mid-sequence

Layer 2: TRANSPORT (ADAPTIVE)

└─ Same GPU:

Shared memory / direct tensor reference

└─ Same machine:

NVLink / PCIe DMA

└─ Same network:

RDMA / NCCL

└─ WAN:

gRPC with protobuf serialization

Layer 1: PHYSICAL

└─ On-die:

HBM-style (TB/s)

└─ NVLink:

900 GB/s

└─ PCIe 5:	64 GB/s
└─ Network:	25 GB/s (200GbE) to variable (WAN)

**Key design principle:** Layer 4 is trainable. Layers 1--3 are infrastructure. Nuclei interact only with Layer 4 and are unaware of the physical transport. A nucleus on the same GPU and one across an ocean differ only in latency - which the protocol adapts to automatically.

### 11.3 SRv6-Inspired Self-Routing

Traditional routing requires a central router making  $N$  decisions per sequence (one per layer). Inspired by IPv6 Segment Routing (RFC 8402, 8754), HULLHB uses **source routing**: the router encodes the entire path as a segment list in the signal header at entry. Each nucleus reads "next stop" and forwards. The router makes ONE decision, then the signal routes itself.

**Loose source routing** enables metacognitive rerouting: if a nucleus encounters low confidence during processing, it can **insert additional hops** into the segment list - for example, routing through Hippocampus for memory retrieval before continuing to the originally planned Basal Ganglia. This implements metacognition (Section 8) as a protocol feature.

### 11.4 CRUSH-Based Decentralized Placement

For distributed deployment across multiple machines, we adapt Ceph's CRUSH algorithm (Controlled Replication Under Scalable Hashing) for deterministic nucleus placement. Given a signal embedding, the fabric topology map, and placement rules, **any node can independently compute where any signal should route** - no central metadata server, no consensus protocol.

Properties inherited from CRUSH: - **Topology-aware**: Placement respects compute hierarchy (GPU → machine → network) - **Minimal disruption**: Adding/removing a node moves only  $\sim 1/n$  of routing - **Weighted**: Heterogeneous hardware receives proportional workload - **Deterministic**: Every node computes the same result - zero coordination overhead

### 11.5 Myelination - Learned Bandwidth

Frequently-used inter-nucleus pathways develop increased bandwidth through training, analogous to biological myelination. The projection dimension between each nucleus pair is a trained parameter: heavily-used pathways (e.g., Occipital → Prefrontal) may operate at 256-dim, while rarely-used pathways (e.g., Hippocampus → Basal Ganglia direct) stay at 32-dim. The bandwidth allocation is optimized by backpropagation on end-task loss.

### 11.6 Connection to Published Work

HULLHB draws on Google Pathways (Dean et al., 2022) for task routing through model regions, Graph Neural Networks (Gilmer et al., 2017) for learned message passing, emergent communication in multi-agent RL (Foerster et al., 2016), Capsule Routing (Hinton et al., 2017) for routing by agreement, and white matter neuroscience (Fields, 2008) for the myelination principle. The novelty is the unified trainable protocol that works at both intra-model and inter-instance scale with SRv6-inspired self-routing and CRUSH-based decentralized placement.

---

## 12. Speculative Cognitive Pipelining

---

### 12.1 Approximate Homomorphic Computation

Neural networks are remarkably tolerant of approximation - dropout destroys 50% of values and improves generalization, INT4 quantization introduces massive rounding error yet preserves output quality, and stochastic depth randomly skips entire layers during training. This tolerance enables a fundamental optimization: **downstream modules can begin processing approximate outputs from upstream modules before final results are available.**

In the residual stream formulation of transformers, intermediate representations at layer  $N/2$  are approximately 90% correlated (by cosine similarity) with final representations at layer  $N$ . Operations on these intermediate representations yield approximately correct results - a form of **approximate homomorphism** over the space of incomplete computations.

### 12.2 Pipeline Architecture

Current sequential processing through the Cognitive Mesh:

```
Occipital (50ms) → wait → Prefrontal (80ms) → wait → Hippocampus (30ms) → Output
Total: 160ms
```

Speculative pipelining:

```
Time 0ms:  Occipital starts, sends DRAFT (after 10ms, ~70% correct)
Time 10ms:  Prefrontal begins on draft
Time 50ms:  Occipital sends CORRECTION DELTA
Time 55ms:  Prefrontal applies correction, sends its own draft
Time 95ms:  Pipeline complete. Total: ~95ms (1.7x speedup)
```

The draft/correction protocol integrates with HULLHB: drafts carry `signal_type=DRAFT` with 64-dim compressed payload (LoRa mode); corrections carry `signal_type=CORRECTION_DELTA` with only the diff. CRUSH routes both deterministically.

### 12.3 GDiT Speculative Denoising

The GDiT's inner/outer diffusion (Section 5) naturally maps to speculative computation. The inner diffusion (2--3 steps, quick preview) IS the speculative draft. The outer diffusion (5 steps, final answer) IS the verification. Downstream nuclei can begin processing the inner diffusion output before the outer diffusion completes, achieving an estimated 37% latency reduction on the denoising pipeline.

A lightweight estimator (2-layer MLP) can predict step  $N$ 's output from step  $N-1$ , enabling overlap of later denoising steps. For steps 4--5, where inter-step changes are small, prediction accuracy exceeds 90%, making speculation nearly free.

## 12.4 Self-Play Data Generation

The same inner diffusion mechanism that serves as a world model (Section 5) can generate training data through self-play:

1. Take a solved ARC task (known transformation rule)
2. Run inner diffusion with deliberate perturbation - different colors, sizes, symmetry axes
3. Generate variations preserving the underlying transformation
4. Compute ground truth for each variation (rule is known)
5. Train on the generated curriculum

This produces an AlphaZero-like self-improving loop where the model's denoising quality on its own generated tasks serves as the curriculum signal: easily solved → too easy, generate harder; unsolvable → found the learning edge (Vygotsky's zone of proximal development).

## 12.5 Connection to Published Work

Speculative cognitive pipelining draws on speculative decoding (Leviathan et al., 2023) for draft-verify parallelism, Jacobi decoding (Santilli et al., 2023) for fixed-point iteration on initially wrong sequences, consistency models (Song et al., 2023) for any-timestep projection, and predictive coding (Rao & Ballard, 1999) for the neuroscience of prediction-correction processing. Self-play data generation draws on AlphaZero (Silver et al., 2018) and DreamCoder (Ellis et al., 2021). The novelty is applying speculative execution to inter-module communication in a multi-module neural architecture, and using the same diffusion mechanism for both reasoning and data generation.

---

# 13. Non-Euclidean Elastic Representations

---

## 13.1 Elastic Tokens

Standard transformers allocate identical resources to every token: same hidden dimension, same number of layers, same compute budget. We introduce elastic tokens where both dimensionality and compute depth adapt per-token based on importance.

Drawing on Neural ODEs (Chen et al., 2018) for continuous-depth networks and Mixture of Depths (Raposo et al., 2024) for per-token layer skipping, each token receives a **stretch factor** - a continuous scalar determining how many effective dimensions and layers it receives. High-importance tokens are stretched (full hidden dim, all layers); background tokens are contracted (truncated dims, skipped layers). The stretch factor is determined by position in the representation space (Section 13.3).

## 13.2 Volumetric Representations

For spatial reasoning tasks, we propose a 3D tensor representation where the Z-axis encodes abstraction depth:

- $Z = 0$ : Input layer (raw grid)
- $Z = 1$ : Feature layer (detected patterns, objects, symmetries)

- $Z = 2$ : Transformation layer (the rule being applied)
- $Z = 3$ : Output layer (predicted result)

The GDiT denoising process operates along  $Z$ , treating abstraction as a spatial dimension. Tokens are projected from 2D into this 3D space using differentiable 3D Gaussian Splatting (Kerbl et al., 2023), where each token becomes a Gaussian with learnable position ( $\mu$ ), covariance ( $\Sigma$ ), opacity ( $\alpha$ ), and features. The covariance matrix IS the elastic property - elongated Gaussians represent stretched tokens; tight Gaussians represent contracted tokens.

### 13.3 Non-Euclidean Geometry

The expert proximity space (32-dim positions, Section 6) currently uses Euclidean distance. The Cognitive Mesh's nucleus hierarchy is inherently tree-like - metacognition at the root, major subsystems as branches, specialists as leaves. Such hierarchies embed naturally in hyperbolic space, where the number of points at distance  $r$  grows exponentially ( $\sim e^r$ ) rather than polynomially ( $\sim r^d$ ).

We propose replacing the Euclidean position space with a Poincaré ball model of hyperbolic geometry:

- **Center**: Metacognition nucleus (hub, connects to everything)
- **Middle ring**: Major subsystems (perception, reasoning, memory)
- **Boundary**: Specialists (visual, spatial, episodic, semantic)

This geometry produces a natural center-to-boundary gradient where center positions receive expanded representational capacity (the metric stretches) and boundary positions receive compressed capacity (the metric contracts). **Level of Detail (Section 10) emerges from the geometry** rather than requiring explicit engineering - the hyperbolic metric IS the LoD allocation function.

### 13.4 Learned Riemannian Metric for Attention

For ARC grid reasoning, we propose a task-conditioned Riemannian metric tensor  $g_{ij}(x)$  that warps the attention geometry per task. A symmetry task curves the metric so that symmetric grid positions are "closer" in attention space. A tiling task curves it so repeating cells are closer. The metric tensor is predicted by a small conditioning network from the task profile (Section 10, grid analysis features).

This means the attention pattern does not need to learn spatial relationships from scratch - **the geometry tells attention which cells are related**, and the model only needs to learn the transformation within that pre-aligned space.

### 13.5 Searchlight LoRA

Standard LoRA adapters apply a fixed low-rank update  $\Delta W = AB$  to the weight matrix. We propose Searchlight LoRA, where the adaptation is **directional** - the low-rank projection targets a specific subspace of the higher-dimensional representation, like a searchlight beam illuminating a cone. The beam direction (which subspace) and width (rank) are both learned and adapt per input, drawing on DyLoRA (Valipour et al., 2023) and AdaLoRA (Zhang et al., 2023).

### 13.6 Neural CRISPR - Self-Modification

For continuous self-improvement without retraining, we propose a CRISPR-inspired mechanism for surgical weight editing: (1) an attention-based targeting mechanism locates the relevant weight subspace, (2) masking zeros the targeted parameters, and (3) a gradient-based update inserts new knowledge. Combined with live migration via the HULLHB fabric (Section 11), this enables hot-swappable self-improvement - upgraded nuclei replace old ones with zero downtime, analogous to VM live migration in cloud infrastructure.

This draws on knowledge editing (Meng et al., 2022 - ROME/MEMIT), the Lottery Ticket Hypothesis (Frankle & Carlin, 2019), and Elastic Weight Consolidation (Kirkpatrick et al., 2017) for protecting important weights during editing.

### 13.7 Connection to Published Work

Non-Euclidean representations draw on Poincaré Embeddings (Nickel & Kiela, 2017) for hyperbolic hierarchical embeddings, Hyperbolic Neural Networks (Ganea et al., 2018) for neural operations in curved space, Geometric Deep Learning (Bronstein et al., 2021) for the non-Euclidean learning framework, and 3D Gaussian Splatting (Kerbl et al., 2023) for differentiable volumetric representations. The novelty is the combination of learned Riemannian attention geometry, LoD-as-metric (compression emerging from curvature rather than explicit budgeting), Gaussian splatting for token projection, and hyperbolic nucleus hierarchy in a unified cognitive architecture.

---

## 14. Training Pipeline

### 14.1 Phase 1--3: Base + Personality + ARC Format (from V1.5)

Phases 1--3 are inherited from the Engram V2 training pipeline and establish the autoregressive backbone's core capabilities.

**Phase 1: Base Pretraining.** Train the transformer backbone from scratch on FineWeb-Edu ~25B tokens. Pure language model, no Engram extensions. Duration: 25--30 days on DGX Spark, ~5 days on Lambda 8xH100.

Hyperparameter	V3-mini	V3-full
Data	FineWeb-Edu ~25B tokens	FineWeb-Edu ~100B tokens
Optimizer	Muon + AdamW	Muon + AdamW
Matrix LR	0.015	0.01
Embedding LR	0.25	0.15
LR schedule	WSD (Warmup-Stable-Decay)	WSD
Precision	BF16	BF16
Batch size	8	64

**Phase 2: Vision + State Alignment.** Freeze backbone. Train vision projector, FiLM generator, state token, attention temperature, state feedback, and MoE router. Duration: 2--3 days on Spark.

**Phase 3: SFT.** Full supervised fine-tuning with personality, memory, and ARC trace data. All parameters unfrozen. Duration: 4--5 days on Spark.

### 14.2 Phase 4: MoE Activation + Skill SFT + GRPO (from V2)

Phase 4 activates the MoE system and trains skill specialization, inherited from V2.

**Phase 4a: MoE Initialization.** Split dense MLP weights into 16 experts per layer (V3-mini) or 48 experts per layer (V3-full). CPU-only, ~5 minutes.

**Phase 4b: Skill SFT.** Train with skill-clustered batches. Duration: 4--6 hours on Spark.

**Phase 4c: GRPO Reinforcement Learning.** Generate tool call sequences, execute against grids, reward correct compositions. Duration: 8--12 hours on Spark.

Hyperparameter	Value
Method	GRPO (Group Relative Policy Optimization)
Completions per prompt	4--8
Reward: personality consistency	0.4 weight
Reward: reasoning quality	0.3 weight
Reward: ARC correctness	0.3 weight
KL penalty	0.05

**Phase 4d: Expert Expansion.** Freeze existing 16 experts. Add 4 metacognition/composition experts. Train new experts on multi-step traces. Unfreeze all at 0.1x LR. Duration: 2--3 hours on Spark.

### 14.3 Phase 5a: GDiT Pre-training (Synthetic Grid Transforms)

This is the first V3-specific training phase. The GDiT is pre-trained on synthetic grid transformations to learn what spatial primitives look like as denoising patterns.

**Data generation.** 100,000 synthetic grid pairs are generated by applying known DSL operations to random grids:

#### SYNTHETIC TRAINING DATA GENERATION

1. Generate random grid (3x3 to 30x30, random colors)
2. Apply 1-3 random DSL operations:
  - Single: rotate\_90, flip\_h, color\_replace, etc.
  - Compound: find\_objects -> extract -> rotate
  - Complex: detect\_symmetry -> fill -> tile
3. Result: (input\_grid, operation\_sequence, output\_grid)
4. Format as training example with 2-3 example pairs + test pair



**Training objective.** Standard MaskGIT training: mask random fractions of the output grid, predict masked cells given the input grid and example pairs as conditioning.

Hyperparameter	V3-mini	V3-full
Synthetic examples	100K	500K
Masking ratio	Uniform(0, 1)	Uniform(0, 1)
Learning rate	3e-4	1e-4
Batch size	64	256
Optimizer	AdamW	AdamW
Duration (Spark)	~6 hours	N/A (cloud)
Duration (8xH100)	~30 min	~4 hours

**What the GDiT learns.** After pre-training on synthetic data, the GDiT has learned: - What rotated grids look like (geometric patterns as denoising attractors) - How colors map between input and output (color transformation patterns) - What tiled grids look like (repeating structure as denoising patterns) - How objects move between input and output positions

This is library learning through denoising: spatial primitives are memorized not as symbolic procedures but as statistical patterns in the denoising process.

### 14.4 Phase 5b: Diffusion World Model Training

The world model is trained alongside the GDiT by running the same denoiser at reduced fidelity and training it to predict approximate outcomes:

WORLD MODEL TRAINING

1. For each training example:

a. Run full GDiT (5 steps) -> precise output grid (target)

b. Run GDiT at 2-3 steps -> approximate output grid (prediction)

c. Loss = CrossEntropy(approximate, target) \* 0.3  
(lower weight because approximation is intentionally imprecise)

2. Additional training signal:

a. Confidence calibration: train the confidence score  
(mean max probability) to correlate with actual accuracy

b. Loss = MSE(predicted\_confidence, actual\_similarity\_score) \* 0.1

The world model training adds negligible compute (2--3 additional denoising steps per training example, at ~30% the cost of the full 5-step GDiT pass) and produces a calibrated confidence estimator.

### 14.5 Phase 5c: Reasoning Mode Router Training

The reasoning mode router is trained through a combination of supervised labels and reinforcement learning:

**Supervised initialization.** Hand-labeled examples associate task types with reasoning modes: - Grid transformation tasks -> GDiT mode - Text explanation tasks -> Autoregressive mode - Ambiguous tasks -> Search mode - Previously-seen-similar tasks -> Retrieval mode

**RL refinement.** After supervised initialization, the router is fine-tuned with GRPO. The reward signal is the final task performance (correct grid, personality consistency), which backpropagates through the mode selection. The router learns to correct its initial biases based on actual performance.

Hyperparameter	Value
Supervised examples	~5,000 (hand-labeled)
RL episodes	~10,000
Mode selection temperature	0.5 (exploration) -> 0.1 (exploitation)
Duration (Spark)	~2 hours

14.6 Phase 5d: Joint End-to-End Fine-tuning

All components --- backbone, GDiT, router, state engine, memory --- are fine-tuned jointly on the full training dataset. This phase is critical for ensuring smooth interaction between components that were trained separately.

Hyperparameter	V3-mini	V3-full
Learning rate	5e-5 (0.1x Phase 3 LR)	2e-5
Components trained	All (backbone + GDiT + router + state)	All
Data	Full training mix + ARC tasks	Full + expanded ARC
Duration (Spark)	~8 hours	N/A (cloud)
Duration (8xH100)	~1 hour	~12 hours
Gradient clipping	1.0	1.0

The low learning rate prevents catastrophic forgetting of previously learned capabilities while allowing the components to adjust their interfaces. The GDiT's conditioning bundle, the router's mode selection, and the state engine's feedback pathways all fine-tune their interactions during this phase.

14.7 Phase 6: Self-Play (AlphaZero for Puzzles)

After all supervised and RL training, the model enters a self-play loop:

SELF-PLAY LOOP

1. Present ARC task to model

2. Model selects reasoning mode (router)

3. Model generates candidate solution(s)
4. Execute against training examples
5. Compare output to expected grid
  - > Match? Add (task, solution\_trace) to training data
  - > No match? Try different approach (up to 5 attempts)
  - > All fail? Skip this task
6. Periodically fine-tune on accumulated self-play data
7. Repeat from step 1

ARC is uniquely suited for self-play because verification is trivial: compare the output grid to the expected grid, cell by cell. No human labeling is needed. No reward model uncertainty. The model literally teaches itself new compound skills.

Hyperparameter	Value
Attempts per task	5
Self-play batch size	100 tasks
Fine-tune frequency	Every 500 successful solutions
Fine-tune LR	2e-5
Temperature for generation	0.7
Maximum self-play rounds	50

Each successful solution becomes training data. Each round of fine-tuning on self-play data produces a stronger model that can solve more tasks in the next round. This is the same mechanism that made AlphaZero (Silver et al., 2018) superhuman at chess --- generate, evaluate, learn, repeat.

## 14.8 Training Phase Summary

Phase	Task	V3-mini (Spark)	V3-full (8xH200)	Key Output
1	Base pretraining	25--30 days	5--10 days	Language model
2	Vision + state alignment	2--3 days	6 hours	Multi-modal model
3	SFT (personality + ARC)	4--5 days	12 hours	Personality model
4a	MoE initialization	5 min	5 min	Expert split
4b	Skill SFT	4--6 hours	30 min	Specialized experts
4c	GRPO RL	8--12 hours	2 hours	RL-tuned routing
4d	Expert expansion	2--3 hours	15 min	Metacognition experts

Phase	Task	V3-mini (Spark)	V3-full (8xH200)	Key Output
5a	GDiT pre-training	6 hours	4 hours	Spatial primitives
5b	World model training	2 hours	2 hours	Mental simulation
5c	Router training	2 hours	1 hour	Mode selection
5d	Joint fine-tuning	8 hours	12 hours	Integrated system
6	Self-play	Ongoing	Ongoing	Self-improvement
<b>Total (before self-play)</b>		<b>~38--45 days</b>	<b>~8--12 days</b>	

## 15. V3-mini: Proof of Concept (2.5B)

### 15.1 Architecture Comparison

V3-mini is designed as a research prototype that validates the cognitive mesh architecture on consumer hardware before scaling to V3-full. Every architectural component is present but scaled to fit within the DGX Spark's 128 GB unified memory during training.

Specification	V3-mini	V3-full
<b>Backbone</b>		
Layers	24	40
Hidden dimension	1536	4096
Attention heads	12 (4 KV, GQA 3:1)	32 (8 KV, GQA 4:1)
Head dimension	128	128
MoE experts/layer	16	48
MoE top-k	2	2
Activation	SwiGLU	SwiGLU
Vocab	32,768	32,768
Sequence length	2,048	4,096
RoPE base	500,000	1,000,000
<b>GDiT</b>		
Denoiser layers	8	12
Denoiser hidden dim	256	512

Specification	V3-mini	V3-full
Denoiser MoE experts	8	16
Example encoder layers	4	6
Example encoder hidden	256	512
Unmasking steps	5	5
<b>State Engine</b>		
State dimensions	72	72
FiLM modulation points	48	80
State pathways	6	6
<b>Memory</b>		
Internal memory bank	No (external only)	Yes (1024 slots)
SoulKeeper external	Yes	Yes
<b>Totals</b>		
Total parameters	~2.5B	~190B
Active parameters	~1.6B	~48B

## 15.2 Component Sizes

Detailed parameter breakdown for V3-mini:

Component	Parameters	Active	Notes
Token embeddings (wte)	50.3M	50.3M	32K x 1536
LM head (untied)	50.3M	50.3M	1536 x 32K
Transformer layers (24)	547M	547M	GQA + SwiGLU per layer
MoE experts (16 x 24)	454M	56.7M	Top-2 of 16
MoE routers + state bias	1.2M	1.2M	Per layer
Expert proximity	1.2M	1.2M	32-dim latent positions
SigLIP-SO400M (frozen)	400M	400M	PseudoDeepStack
Vision projector	7.7M	7.7M	3456 -> 1536 -> 1536
FiLM generator	322K	322K	3-layer encoder, 48 points
State token	115K	115K	Linear(72, 1536)
MoE state bias (all layers)	28K	28K	Linear(72, 16) x 24

Component	Parameters	Active	Notes
Attention temperature	21K	21K	Linear(72, 12) x 24
State feedback MLP	175K	175K	MLP(1608, 128, 72)
GDiT denoiser (8 layers, MoE)	45M	15M	8 experts, top-2
GDiT cross-attention	8M	8M	8 layers
Example encoder	10M	10M	4-layer transformer
Cell embeddings	3K	3K	11 x 256
2D position bias	29K	29K	59 x 59 x 8 heads
Step embedding	2K	2K	8 x 256
Reasoning mode router	500K	500K	3-layer MLP
GridAnalyzer	0	0	Numpy, no learned params
Adversarial probe	200K	0	Training only
<b>Total</b>	<b>~2,540M</b>	<b>~1,615M</b>	

## 15.3 Training on DGX Spark

V3-mini is designed to train entirely on a single DGX Spark GB10 at zero cloud cost:

### Memory budget during training:

Component	Memory (BF16)
V3-mini model	~5.1 GB
SigLIP (frozen, inference only)	~0.8 GB
GDiT + example encoder	~0.2 GB
Optimizer (AdamW, trainable params only)	~15 GB
Activations (batch=8, seq=2048)	~8 GB
Gradient checkpointing overhead	~2 GB
<b>Total training</b>	<b>~31 GB</b>
<b>Headroom (of 128 GB)</b>	<b>~97 GB</b>

V3-mini uses 24% of available memory during training. This leaves substantial headroom for larger batch sizes (up to 32) or gradient accumulation without approaching memory limits.

### Training timeline (DGX Spark):

Phase	Duration	Cumulative
Phase 1 (pretraining)	25--30 days	25--30 days
Phase 2 (vision + state)	2--3 days	27--33 days
Phase 3 (SFT)	4--5 days	31--38 days
Phase 4 (MoE + RL)	~2 days	33--40 days
Phase 5 (GDiT + router)	~1 day	34--41 days
Phase 5d (joint fine-tune)	~8 hours	34--42 days
<b>Total</b>	<b>~34--42 days</b>	

The entire V3-mini training pipeline completes in approximately 5--6 weeks on a single DGX Spark, at zero cloud cost. Self-play can continue indefinitely thereafter.

## 15.4 Inference

V3-mini at inference:

Property	Value
Model size (BF16)	~3.2 GB
Model size (INT8)	~1.6 GB
Model size (INT4)	~0.8 GB
SigLIP (FP16)	~0.8 GB
GDiT + encoder (BF16)	~0.13 GB
<b>Total inference (INT8)</b>	<b>~2.5 GB</b>
<b>Total inference (INT4)</b>	<b>~1.7 GB</b>
Autoregressive speed (DGX Spark)	~150 tok/s
GDiT speed (5-step)	~80ms/grid
World model speed (2-step)	~30ms/simulation

At INT4 quantization, V3-mini fits in 1.7 GB --- small enough to run on a laptop with 8 GB RAM, a Raspberry Pi 5 with external storage, or any modern smartphone. This enables edge deployment of the full cognitive mesh architecture.

## 15.5 V3-mini as Research Prototype

V3-mini serves several research purposes beyond being a smaller version of V3-full:

**Architecture validation.** Does the cognitive mesh (five-layer integration) actually work? Does the reasoning mode router learn meaningful mode selection? Does GDiT outperform autoregressive-only reasoning on spatial tasks? V3-mini answers these questions at minimal compute cost.

**Hyperparameter exploration.** GDiT has many design choices: number of unmasking steps, confidence thresholds for re-masking, world model fidelity, router training schedule. V3-mini enables rapid iteration on these choices.

**Ablation studies.** Systematic removal of components (GDiT without state conditioning, router without world model, backbone without MoE) reveals the contribution of each architectural innovation.

**Baseline for scaling laws.** V3-mini performance establishes the lower bound of the scaling curve. Combined with V3-full projections and intermediate checkpoints, this enables scaling law analysis for the cognitive mesh architecture.

---

## 16. V3-full: Competition Model (190B)

---

### 16.1 Architecture Specification

V3-full scales the cognitive mesh architecture to competition grade:

Specification	Value
Backbone layers	40
Hidden dimension	4096
Attention heads	32 (8 KV, GQA 4:1)
Head dimension	128
MoE experts per layer	48
MoE active per token	2 (top-2)
Expert intermediate dim	2048 (SwiGLU adjusted)
Backbone total params	~120B (sparse MoE)
Backbone active params	~45B
GDiT denoiser layers	12
GDiT hidden dim	512
GDiT MoE experts	16
Example encoder layers	6
Example encoder hidden	512
Total parameters	~190B



Specification	Value
Active parameters	~48B

The 48-expert MoE with top-2 routing means only 4.2% of expert parameters are active per token. This extreme sparsity enables the model to contain 190B total parameters while requiring only 48B active --- a 4x efficiency ratio.

### 16.2 Deployment on 1x DGX Spark

V3-full is designed to run inference on a single DGX Spark GB10 using MXFP4 (microscaling 4-bit floating point) quantization:

Component	MXFP4 Size	Active Memory
Backbone (190B total, 48B active)	~48 GB (active only loaded)	~12 GB (active)
Sparse expert pages (on-demand)	~70 GB (disk/NVMe paging)	~10 GB (hot cache)
GDiT + example encoder	~0.4 GB	~0.4 GB
SigLIP (FP16)	~0.8 GB	~0.8 GB
KV cache (seq=4096)	~8 GB	~8 GB
Memory bank (1024 slots)	~0.001 GB	~0.001 GB
<b>Total GPU memory</b>		<b>~31 GB</b>
<b>With headroom (128 GB)</b>		<b>~97 GB free</b>

The key insight is that with 48 experts and top-2 routing, only  $2/48 = 4.2\%$  of expert parameters are needed per token. Expert weights can be paged from NVMe storage on demand, with a hot cache holding the most frequently-used experts in GPU memory. The DGX Spark's unified memory architecture (CPU and GPU share the same 128 GB) makes this paging efficient --- there is no PCIe bottleneck.

### 16.3 2x DGX Spark Variant (400B)

With two DGX Spark systems (256 GB total), V3-full can be scaled further:

Configuration	Total Params	Active Params	Memory
1x Spark, 48 experts	190B	48B	95 GB (MXFP4)
2x Spark, 96 experts	400B	48B	190 GB (MXFP4)

Doubling expert count from 48 to 96 doubles total parameters while keeping active parameters identical (still top-2 routing). Each Spark holds half the experts. Cross-device communication occurs only when the router selects an expert on the other device, which happens for approximately 50% of tokens (random assignment) --- manageable over the inter-device NVLink connection.

## 16.4 Training Requirements

V3-full requires cloud compute for training:

Phase	Hardware	Duration	Cost Estimate
Phase 1 (pretraining, 100B tokens)	8x H200	~10 days	\$8,000--\$15,000
Phase 2 (vision + state)	8x H200	~6 hours	\$200
Phase 3 (SFT)	8x H200	~12 hours	\$400
Phase 4 (MoE + RL)	8x H200	~2 days	\$1,500
Phase 5 (GDiT + router + joint)	8x H200	~1 day	\$800
Phase 6 (self-play, initial)	8x H200	~3 days	\$2,400
Total	8x H200	~17 days	\$13,000--\$20,000

For maximum scale (48x H200 for faster iteration):

Configuration	Duration	Cost
8x H200	~17 days	~\$15K
16x H200	~9 days	~\$18K
48x H200	~4 days	~\$25K

## 16.5 Cost Estimates

Scenario	Hardware	Duration	Total Cost
V3-mini only (research)	1x DGX Spark	~42 days	\$0 (owned)
V3-full minimal	8x H200	~17 days	~\$15,000
V3-full accelerated	48x H200	~4 days	~\$25,000
V3-full + extended self-play	8x H200	~30 days	~\$25,000
Full roadmap (V3-mini + V3-full + self-play)	Mixed	~3 months	~\$30,000--\$50,000

# 17. Performance Projections

## 17.1 ARC-AGI Evaluation

Performance projections on the ARC-AGI public evaluation set (400 tasks):

System	ARC-AGI Score	Architecture	Compute
<b>Baselines</b>			
GPT-4o (prompted)	~50--60%	1.8T autoregressive	Datacenter
Gemini 1.5 Pro (prompted)	~50--60%	Unknown	Datacenter
Ryan Greenblatt 2024	72%	LLM + program synthesis + search	Datacenter
MindsAI 2024 winner	84%	Proprietary ensemble	Unknown
Human average	~85--95%	Biological	20W
<b>Engram versions</b>			
V1.5 (542M, SFT only)	6--10%	Autoregressive, tool format learned	1x Spark
V2 (2.1B, MoE + RL)	20--30%	Skill MoE + metacognition	1x Spark
<b>V3-mini (2.5B)</b>	<b>35--50%</b>	<b>Cognitive mesh + GDiT</b>	<b>1x Spark</b>
<b>V3-full (190B)</b>	<b>65--80%</b>	<b>Full cognitive mesh</b>	<b>1x Spark (inference)</b>
V3-full + self-play	75--85%	+ AlphaZero-style loop	8x H200 (training)

**V3-mini projection rationale (35--50%).** The GDiT adds approximately 15--20 percentage points over V2's autoregressive-only approach on spatial tasks. The world model adds approximately 5--10 points by filtering bad hypotheses before commitment. The reasoning mode router adds approximately 5 points by selecting the right strategy per task. Combined with V2's 20--30% base, this yields 35--50%.

**V3-full projection rationale (65--80%).** The 190B parameter model with 48 active experts provides substantially deeper reasoning capacity than the 2.5B V3-mini. The larger GDiT (12 layers, 512 hidden, 16 experts) handles more complex spatial transformations. The expanded DSL ceiling (80% with additional operations) raises the theoretical maximum. Tree search with multiple candidates adds robustness.

## 17.2 Personality and Memory Quality

Capability	V2	V3-mini	V3-full
Personality voice (Slick/TARS)	9/10	9/10	9/10
Memory recall across sessions	9/10	9/10	9.5/10
Emotional coherence across turns	8/10	8/10	9/10
Adaptive tone matching	7/10	7/10	8/10
Self-awareness / metacognition	6/10	7/10	8/10
Spatial reasoning explanation	3/10	7/10	9/10

The V3 architecture improves spatial reasoning explanation because the GDiT provides structured information about the transformation (which cells were unmasked in which order, what the confidence trajectory was) that the

autoregressive backbone can reference when generating explanations. V2 had to reason about grids in text; V3 reasons about grids spatially and explains in text.

### 17.3 Inference Speed Comparison

System	Time per ARC Task	Hardware
GPT-4o (prompted)	~5--15s	Cloud API
Ryan Greenblatt 2024	~30--120s (search)	Cloud
V3-mini (single candidate)	~200ms	1x Spark
V3-mini (5 candidates + search)	~1.5s	1x Spark
V3-full (single candidate)	~500ms	1x Spark
V3-full (imagination + 3 candidates)	~3s	1x Spark

V3 is 10--100x faster than search-based LLM approaches because: 1. GDiT generates a complete grid in ~80ms (5 denoising steps), not token-by-token 2. World model filters bad hypotheses in ~30ms each, avoiding full evaluation 3. The reasoning mode router avoids unnecessary computation by selecting the right strategy upfront

### 17.4 Comparison with Frontier Systems

Dimension	GPT-4o	MindsAI 2024	V3-full (projected)
ARC-AGI score	~55%	84%	65--80%
Model size	~1.8T	Unknown	190B (48B active)
Inference hardware	Datacenter GPU cluster	Unknown	1x DGX Spark
Inference speed	~10s/task	Unknown	~3s/task
Spatial reasoning	Autoregressive (1D)	Likely hybrid	GDiT (native 2D)
World model	No	Unknown	Yes (diffusion)
Personality	Generic	None	Biologically-grounded (72-dim)
Memory	Context window only	None	SoulKeeper + internal bank
Cost to train	>\$100M	Unknown	~\$15--25K
Cost to run	~\$0.01/task (API)	Unknown	~\$0 (owned hardware)

The comparison reveals V3-full's competitive position: it achieves comparable ARC performance to systems 10--40x larger by using architectural innovations (GDiT, world model, state engine) rather than brute-force scaling. The cost differential --- \$15--25K versus >\$100M --- reflects the fundamental advantage of problem-specific architecture over general-purpose scale.

## 18. Novel Contributions

---

The following is a numbered list of innovations introduced by Engram V3, with detailed explanation of what makes each novel:

- 1. Discrete diffusion for few-shot abstract grid reasoning.** No prior work has applied MaskGIT-style discrete diffusion to abstract reasoning over small categorical grids. Existing discrete diffusion work (D3PM, MDLM, SEDD) targets language modeling or image generation. The application to ARC-style few-shot grid transformation --- where the model must infer an abstract rule from 2--5 example pairs and apply it through iterative denoising --- is novel.
- 2. Unified output generation and world model through two-fidelity denoising.** The diffusion world model reuses the GDiT architecture at reduced fidelity (2--3 steps instead of 5) for mental simulation. No separate world model architecture or parameters are needed. This unification --- where simulation and generation are the same computation at different precision --- has no direct analog in prior world model work (Ha & Schmidhuber, Dreamer, MuZero all use separate world model networks).
- 3. Intrinsic confidence from denoising quality.** The world model derives confidence from the quality of the denoised grid itself: clean, coherent grids indicate high confidence; blurry, inconsistent grids indicate uncertainty. This eliminates the need for a separate critic or value network. Prior work requires trained value functions (MuZero), learned reward predictors (Dreamer), or explicit confidence classifiers. The denoiser IS the critic.
- 4. State-conditioned discrete diffusion.** The integration of a 72-dimension biologically-structured state vector with discrete diffusion through FiLM modulation, MoE routing bias, and attention temperature. The state vector --- spanning autonomic arousal through metacognition --- conditions every denoising layer. No prior discrete diffusion work incorporates emotional or cognitive state conditioning.
- 5. Five-layer cognitive mesh architecture.** The integration of state engine, perception, working memory (autoregressive + MoE), reasoning modes (auto + GDiT + search + retrieval), and metacognition into a single differentiable system. While individual components exist in prior work, their integration into a unified cognitive mesh with bidirectional state flow across all layers is novel.
- 6. Learned reasoning mode router with mid-sequence switching.** A learned network that selects among four reasoning strategies (autoregressive, GDiT, search, retrieval) and can switch modes during generation. Prior multi-strategy systems (e.g., Greenblatt's search + LLM) use fixed strategy selection. The learned router discovers which strategy suits which task type and can delegate spatial reasoning to GDiT mid-sentence.
- 7. Metacognitive verification within iterative denoising.** After each GDiT unmasking step, dedicated metacognition experts assess the partial grid, update the confidence state, and can trigger re-masking of suspicious cells. This integrates process reward models (Lightman et al., 2023) with discrete diffusion at the step level, providing continuous self-assessment during spatial reasoning.
- 8. MoE denoiser with state-biased routing.** The GDiT denoiser uses a mixture-of-experts FFN with routing conditioned on both token content and the 72-dim state vector. Different emotional/cognitive states activate

different denoiser experts, enabling state-dependent spatial reasoning strategies. Prior MoE work does not combine sparse expert routing with discrete diffusion.

9. **Library learning through denoising patterns.** Spatial primitives (rotations, reflections, tilings) are learned as denoising patterns rather than symbolic procedures. Pre-training on synthetic grid transformations teaches the GDiT what transformed grids look like, creating a neural complement to the symbolic 42-function DSL. This bridges the gap between DreamCoder-style symbolic library learning and neural pattern recognition.
  10. **Imagination-based planning with denoising simulation.** The model simulates up to 50 hypothetical grid transformations (at ~30ms each) before committing to 2--3 for full execution. This implements Vincent's "daydreaming" insight: simulate many possibilities, let the intrinsic confidence signal (denoising quality) filter bad options, then execute only the promising ones.
  11. **Emotional salience-weighted memory writing.** The internalized memory bank (V3-full) prioritizes emotionally significant information for storage, mirroring biological amygdala-hippocampal memory encoding. State dimensions (surprise, joy, fear, confidence change) boost writing priority. No prior agent memory system incorporates emotional salience into memory writing.
  12. **Dual-scale design for research and competition.** V3-mini (2.5B) validates the architecture at zero cloud cost on consumer hardware; V3-full (190B) targets competitive performance. Both share identical architectural principles, enabling clean scaling law analysis. The V3-mini -> V3-full progression provides a research-to-deployment pathway accessible to independent researchers.
- 

## 19. Limitations and Future Work

---

### 19.1 Current Limitations

**Example encoder difficulty.** The example encoder --- which must infer abstract transformation rules from 2--5 concrete examples --- is the hardest component to train effectively. Abstract rule inference from few examples is an open problem in machine learning. The GDiT's performance ceiling is determined by the example encoder's ability to extract rules, not by the denoiser's ability to apply them.

**DSL coverage ceiling.** The 42-function DSL covers approximately 60% of ARC tasks. Even with V3's architectural innovations, tasks requiring operations outside the DSL (conditional branching, counting-based logic, complex path finding) remain unsolvable through the skill decomposition pathway. The GDiT's learned spatial primitives may partially compensate for missing DSL functions, but this is unverified.

**GDiT scale uncertainty.** While the GDiT architecture is well-motivated by the properties of ARC grids, its effectiveness at the V3-full scale (12-layer, 512-dim, 16-expert denoiser) is theoretically projected but empirically unverified. The V3-mini prototype will provide the first empirical data.

**World model calibration.** The confidence signal from denoising quality (mean max probability) is a proxy for actual correctness. This proxy may be poorly calibrated: the model might be confidently wrong (high max

probability on incorrect cells) or uncertainly right (low max probability on correct cells). Calibration training (Phase 5b) addresses this, but residual miscalibration may limit planning effectiveness.

**Training pipeline complexity.** The 6-phase training pipeline with multiple objectives (language modeling, skill specialization, GDiT denoising, world model calibration, router training, joint fine-tuning) is complex. Phase interactions are difficult to predict: improvements in Phase 5a (GDiT pre-training) may require re-tuning Phase 5c (router training). Ablation studies are needed to understand which phases are critical and which can be shortened.

**Single-task focus.** V3 is designed primarily for ARC-AGI. While the cognitive mesh architecture is general (the autoregressive backbone handles any text task), the GDiT, world model, and skill decomposition are ARC-specific. Extending to other spatial reasoning domains (mathematics, physics simulation, game playing) would require domain-specific adaptations.

**Adversarial robustness.** ARC-AGI's semi-private evaluation set may contain tasks specifically designed to defeat pattern-matching approaches. The GDiT's denoising patterns, learned from synthetic data and ARC training tasks, may fail on adversarially-constructed problems.

## 19.2 Future Work

**Expanded DSL.** Add conditional branching, counting-based logic, path-finding, and pattern-matching operations to raise DSL coverage from 60% to approximately 80% of ARC tasks. Each new operation adds new compound skill possibilities.

**Continuous learning.** Extend the self-play loop to continuous learning: the model encounters new tasks in deployment, attempts solutions, receives feedback, and improves without explicit retraining. This requires careful management of catastrophic forgetting.

**Multi-modal GDiT.** Extend the GDiT from discrete grids to other structured domains: mathematical diagrams, circuit layouts, molecular structures. The MaskGIT-style approach generalizes to any discrete categorical spatial structure.

**Temporal extension.** Add a temporal dimension to the GDiT for video-like sequences of grids (ARC tasks where the transformation involves temporal patterns). This extends the 2D position bias to 3D (row, column, time).

**Hierarchical diffusion.** For very large grids (approaching 30x30), implement hierarchical denoising: first denoise at a coarse resolution (block-level), then refine at full resolution. This mirrors human visual processing (peripheral vision first, then foveal attention).

**Cross-architecture knowledge transfer.** Can a GDiT trained on ARC transfer spatial reasoning capabilities to the autoregressive backbone? Distillation from GDiT to the backbone could improve the backbone's spatial reasoning without requiring GDiT at inference time.

**Learned state dynamics.** Replace SoulKeeper's rule-based state engine with a fully-learned state dynamics model that discovers event-to-state mappings from interaction data.

**Multi-agent cognitive mesh.** Multiple Engram V3 instances with shared memory could collaborate on complex tasks, with each agent's state influencing the others --- implementing emergent group dynamics through shared state vectors.

---

## 20. Conclusion

---

Engram V3 demonstrates that the autoregressive bottleneck in spatial reasoning is an *architectural* problem with an *architectural* solution. By integrating a Grid Diffusion Transformer for iterative 2D reasoning alongside an autoregressive transformer for sequential reasoning, the cognitive mesh architecture addresses the fundamental mismatch between 1D token generation and 2D spatial logic.

The key insight is not merely that diffusion can handle grids --- it is that the cognitive processes required for abstract reasoning (holistic perception, iterative refinement, hypothesis testing, mental simulation, metacognitive monitoring) map naturally onto the diffusion framework. MaskGIT-style unmasking mirrors human solving strategy: easy cells first, hard cells last, with continuous re-assessment. Two-fidelity denoising implements mental simulation without a separate world model. Confidence emerges from denoising quality without a separate critic. These are not bolt-on features --- they are inherent properties of the diffusion process applied to reasoning.

The five-layer cognitive mesh integrates these capabilities into a coherent system. The state engine provides the binding mechanism: a 72-dimension signal that modulates perception (through attention temperature), working memory (through FiLM and MoE routing), reasoning (through state-conditioned denoising), and metacognition (through confidence tracking and re-masking thresholds). Every layer influences every other layer through the state, creating the bidirectional feedback loops that characterize biological cognition.

The dual-scale design --- V3-mini (2.5B) on consumer hardware, V3-full (190B) for competition --- demonstrates that these ideas are not confined to the datacenter. V3-mini validates the architecture at zero cloud cost; V3-full pushes performance to competitive levels at approximately \$15--25K total training cost. The progression from V1.5 (6--10% on ARC) through V2 (20--30%) to V3-mini (35--50%) and V3-full (65--80%) reflects the cumulative impact of principled architectural innovation: each generation adds capabilities that multiply rather than merely add.

The self-play loop (Phase 6) transforms the architecture from a static model into a self-improving system. With verification trivial for ARC (compare grids), each correct solution becomes training data, creating a virtuous cycle analogous to AlphaZero. The GDiT's world model enables efficient hypothesis evaluation, the metacognitive loop ensures quality control, and the skill library grows with each successful solution.

What makes Engram V3 genuinely novel is not any single component --- discrete diffusion, MoE, world models, and metacognition all exist in prior work. It is their integration into a unified cognitive mesh where state-conditioned diffusion, imagination-based planning, skill-specialized experts, and continuous self-monitoring operate as a single, differentiable system. The architecture treats reasoning not as a sequence of tokens but as a mesh of interacting cognitive processes --- each specialized, all coordinated, and all grounded in a biologically-structured internal state.

The autoregressive bottleneck is real. The diffusion opportunity is real. The cognitive mesh is the architecture that connects them.

---

## Appendix A: GDiT Architecture Details

---



## A.1 V3-mini GDiT Specification

Component	Specification	Parameters
<b>Cell Embedding</b>		
Color embedding	11 values x 256 dim	2,816
2D position embedding (row)	30 positions x 256 dim	7,680
2D position embedding (col)	30 positions x 256 dim	7,680
Role embedding	3 roles (input/output/mask) x 256 dim	768
Pair embedding	8 max pairs x 256 dim	2,048
Step embedding	8 steps x 256 dim	2,048
<b>2D Relative Position Bias</b>		
Bias table	59 x 59 x 8 heads	27,848
<b>Example Encoder (4 layers)</b>		
Self-attention Q/K/V per layer	3 x (256 x 256) + 3 x 256	197,376
Self-attention output per layer	256 x 256 + 256	65,792
Layer norm per layer (x2)	2 x 256 x 2	1,024
FFN per layer	256 x 1024 + 1024 + 1024 x 256 + 256	525,568
Total encoder (4 layers)		~3.2M
Pair pooling + rule pooling	256 x 256 + 256 x 2	~131K
<b>Total example encoder</b>		<b>~3.3M</b>
<b>Denoiser Backbone (8 layers)</b>		
Self-attention per layer	QKV: 3 x (256 x 256), out: 256 x 256	262,144
Cross-attention per layer	Q: 256x256, KV: 256x256x2, out: 256x256	262,144
FiLM modulation per layer (x3 norms)	Shared encoder + 3 extra gamma/beta	~5K
MoE FFN per layer		
-- Router	256 x 8 + bias	2,056
-- State bias	72 x 8 + bias	584
-- 8 experts, each:		
---- gate (SwiGLU)	256 x 128	32,768
---- up (SwiGLU)	256 x 128	32,768
---- down (SwiGLU)	128 x 256	32,768

Component	Specification	Parameters
-- Total MoE per layer		~790K
Layer norms per layer (x3)	3 x 256 x 2	1,536
Attn temperature per layer	72 x 8 + bias	584
Total per denoiser layer		~1.3M
<b>Total denoiser (8 layers)</b>		<b>~10.4M</b>
<b>Output head</b>		
Linear(256, 11) + bias		2,827
<b>State conditioning (GDiT-specific)</b>		
FiLM encoder (shared with backbone)	0 (shared)	0
FiLM gamma/beta patterns (GDiT)	24 x 256 x 2	12,288
FiLM scale generators (GDiT)	256 x 24 x 2	12,288
<b>Total GDiT</b>		<b>~15M</b>

## A.2 V3-full GDiT Specification

Component	V3-mini	V3-full
Hidden dimension	256	512
Denoiser layers	8	12
Attention heads	8	16
MoE experts	8 (top-2)	16 (top-2)
Expert intermediate dim	128	256
Example encoder layers	4	6
Example encoder heads	8	16
Denoiser params	~10M	~120M
Example encoder params	~3M	~25M
State conditioning params	~25K	~100K
Output head params	~3K	~6K
<b>Total GDiT params</b>	<b>~15M</b>	<b>~150M</b>
<b>Active GDiT params</b>	<b>~8M</b>	<b>~70M</b>

## A.3 Inference Timing Breakdown

Operation	V3-mini	V3-full
GridAnalyzer (numpy)	~2ms	~2ms
Example encoding	~5ms	~15ms
Input encoding	~3ms	~8ms
Per denoising step	~10ms	~25ms
Metacognition check	~2ms	~5ms
Denoising (5 steps)	~50ms	~125ms
Metacognition (5 checks)	~10ms	~25ms
<b>Total (1 candidate)</b>	<b>~70ms</b>	<b>~175ms</b>
<b>Total (5 candidates)</b>	<b>~350ms</b>	<b>~875ms</b>
World model (1 simulation)	~25ms	~60ms
Planning (50 simulations)	~1.25s	~3s

## Appendix B: 42-Function DSL Reference

### B.1 Geometric Transforms (6 functions)

```
def rotate_90(grid: np.ndarray) -> np.ndarray:
    """Rotate grid 90 degrees clockwise."""

def rotate_180(grid: np.ndarray) -> np.ndarray:
    """Rotate grid 180 degrees."""

def rotate_270(grid: np.ndarray) -> np.ndarray:
    """Rotate grid 270 degrees clockwise (90 CCW)."""

def flip_horizontal(grid: np.ndarray) -> np.ndarray:
    """Mirror grid left-to-right."""

def flip_vertical(grid: np.ndarray) -> np.ndarray:
    """Mirror grid top-to-bottom."""

def transpose(grid: np.ndarray) -> np.ndarray:
    """Swap rows and columns (matrix transpose)."""
```

### B.2 Color Operations (7 functions)

```

def color_replace(grid: np.ndarray, old: int, new: int) -> np.ndarray:
    """Replace all cells of color 'old' with color 'new'."""

def color_swap(grid: np.ndarray, c1: int, c2: int) -> np.ndarray:
    """Swap colors c1 and c2 throughout the grid."""

def color_fill(grid: np.ndarray, color: int) -> np.ndarray:
    """Fill the entire grid with a single color."""

def color_map(grid: np.ndarray, mapping: dict) -> np.ndarray:
    """Apply a color remapping dictionary {old: new}."""

def most_common_color(grid: np.ndarray) -> int:
    """Return the most frequent color in the grid."""

def least_common_color(grid: np.ndarray) -> int:
    """Return the least frequent non-background color."""

def count_color(grid: np.ndarray, color: int) -> int:
    """Count cells of a specific color."""

```

### B.3 Object Detection (7 functions)

```

def find_objects(grid: np.ndarray) -> list:
    """Connected component detection. Returns list of Object."""

def find_rectangles(grid: np.ndarray) -> list:
    """Find rectangular regions. Returns list of Rect."""

def find_lines(grid: np.ndarray) -> list:
    """Find horizontal and vertical lines. Returns list of Line."""

def find_borders(grid: np.ndarray) -> np.ndarray:
    """Identify border cells. Returns binary mask grid."""

def find_background(grid: np.ndarray) -> int:
    """Identify the background color (most common edge color)."""

def extract_object_grid(grid: np.ndarray, obj) -> np.ndarray:
    """Extract the minimal bounding sub-grid containing an object."""

def count_objects(grid: np.ndarray) -> int:
    """Count distinct connected components."""

```

### B.4 Structural Analysis (4 functions)

```

def detect_symmetry(grid: np.ndarray) -> str:
    """Detect symmetry type: horizontal, vertical, diagonal,
    rotational, or none."""

def find_repeating_pattern(grid: np.ndarray):
    """Identify the minimal repeating unit (tile)."""

def detect_grid_structure(grid: np.ndarray):
    """Find sub-grid organization (e.g., 3x3 grid of 5x5 cells)."""

def crop_to_content(grid: np.ndarray) -> np.ndarray:
    """Remove background-colored borders."""

```

## B.5 Spatial Operations (11 functions)

```

def crop(grid: np.ndarray, r1: int, c1: int, r2: int, c2: int) -> np.ndarray:
    """Extract rectangular region from (r1,c1) to (r2,c2)."""

def pad(grid: np.ndarray, n: int, color: int) -> np.ndarray:
    """Add n-cell border of given color around the grid."""

def resize(grid: np.ndarray, h: int, w: int) -> np.ndarray:
    """Scale grid to new dimensions using nearest-neighbor."""

def tile(grid: np.ndarray, h_reps: int, w_reps: int) -> np.ndarray:
    """Repeat grid in a tiling pattern."""

def place_subgrid(grid: np.ndarray, sub: np.ndarray,
                  r: int, c: int) -> np.ndarray:
    """Place sub-grid at position (r, c) in the grid."""

def overlay(grid1: np.ndarray, grid2: np.ndarray) -> np.ndarray:
    """Combine two grids. Non-zero cells in grid2 overwrite grid1."""

def concat_horizontal(g1: np.ndarray, g2: np.ndarray) -> np.ndarray:
    """Concatenate two grids side by side."""

def concat_vertical(g1: np.ndarray, g2: np.ndarray) -> np.ndarray:
    """Concatenate two grids top to bottom."""

def draw_line(grid: np.ndarray, r1: int, c1: int,
              r2: int, c2: int, color: int) -> np.ndarray:
    """Draw a straight line between two points."""

def draw_rectangle(grid: np.ndarray, r1: int, c1: int,
                  r2: int, c2: int, color: int) -> np.ndarray:
    """Draw a rectangle outline."""

```

```
def flood_fill(grid: np.ndarray, r: int, c: int,
               color: int) -> np.ndarray:
    """Fill connected region starting from (r, c) with color."""
```

## B.6 Comparison (5 functions)

```
def grids_equal(g1: np.ndarray, g2: np.ndarray) -> bool:
    """Exact cell-by-cell comparison."""

def diff_grids(g1: np.ndarray, g2: np.ndarray) -> np.ndarray:
    """Produce grid showing differences (0 = same, 1 = different)."""

def similarity_score(g1: np.ndarray, g2: np.ndarray) -> float:
    """Fraction of matching cells (0.0 to 1.0)."""

def get_shape(grid: np.ndarray) -> tuple:
    """Return (height, width) of the grid."""

def get_colors(grid: np.ndarray) -> set:
    """Return the set of unique colors present in the grid."""
```

## Appendix C: Training Hyperparameters

### C.1 Phase 1: Base Pretraining

Hyperparameter	V3-mini	V3-full
Data	FineWeb-Edu ~25B tokens	FineWeb-Edu ~100B tokens
Optimizer (matrices)	Muon	Muon
Optimizer (embeddings)	AdamW (beta1=0.9, beta2=0.95)	AdamW (beta1=0.9, beta2=0.95)
Matrix learning rate	0.015	0.01
Embedding learning rate	0.25	0.15
LR schedule	WSD (10% warmup, 80% stable, 10% decay)	WSD
Weight decay	0.1	0.1
Precision	BF16	BF16
Batch size (per device)	8	64 (8 per device x 8 devices)
Sequence length	2048	4096
Gradient clipping	1.0	1.0

Hyperparameter	V3-mini	V3-full
Dropout	0.0	0.0
z-loss weight	1e-4	1e-4
Softcap	$\tanh(\text{logits}/15)*15$	$\tanh(\text{logits}/30)*30$

## C.2 Phase 2: Vision + State Alignment

Hyperparameter	V3-mini	V3-full
Data	LLaVA-Pretrain 558K + v2-665K	Same + expanded multimodal
Optimizer	AdamW	AdamW
Learning rate	2e-4	1e-4
Epochs	1	1
Trainable	Projector, FiLM, state token, attn temp, feedback, router	Same
Frozen	Backbone, SigLIP, expert weights	Same
Batch size	16	128
State dropout	15%	15%

## C.3 Phase 3: SFT

Hyperparameter	V3-mini	V3-full
Data mix weights	SmolTalk 1.0, LLaVA 1.0, convos 10.0, identity 5.0, state 3.0, ARC 3.0, func 2.0	Same
Optimizer	AdamW	AdamW
Learning rate	1e-4	5e-5
Trainable	All	All
Batch size	8	64
Adversarial probe weight	0.01	0.01
State dropout	15%	15%

## C.4 Phase 4c: GRPO

Hyperparameter	V3-mini	V3-full
Completions per prompt	4	8
KL penalty	0.05	0.03

Hyperparameter	V3-mini	V3-full
Reward: personality	0.4	0.3
Reward: reasoning	0.3	0.3
Reward: ARC correctness	0.3	0.4
Temperature	0.7	0.7
Batch size	4	32

## C.5 Phase 5a: GDiT Pre-training

Hyperparameter	V3-mini	V3-full
Synthetic examples	100K	500K
Masking ratio	Uniform(0, 1)	Uniform(0, 1)
Optimizer	AdamW (beta1=0.9, beta2=0.99)	AdamW
Learning rate	3e-4	1e-4
LR schedule	Cosine decay	Cosine decay
Warmup steps	500	2000
Batch size	64	256
Weight decay	0.01	0.01
Gradient clipping	1.0	1.0
State conditioning	Enabled (random state vectors)	Enabled
MoE load balance weight	0.01	0.01

## C.6 Phase 5d: Joint Fine-tuning

Hyperparameter	V3-mini	V3-full
Learning rate	5e-5	2e-5
Trainable	All components	All components
Data	Full mix + ARC + synthetic grids	Same + expanded
Batch size	4	32
Gradient clipping	1.0	1.0
Backbone LR multiplier	0.1x	0.1x
GDiT LR multiplier	1.0x	0.5x



Hyperparameter	V3-mini	V3-full
Router LR multiplier	0.5x	0.3x

### C.7 Phase 6: Self-Play

Hyperparameter	V3-mini	V3-full
Attempts per task	5	5
Batch size (tasks)	100	500
Fine-tune frequency	Every 500 solutions	Every 2000 solutions
Fine-tune LR	2e-5	1e-5
Temperature	0.7	0.7
Max self-play rounds	50	100
Mode router exploration	0.3 (annealed to 0.05)	0.2 (annealed to 0.05)

## Appendix D: V3-mini Memory Budget

### D.1 Training Memory (BF16)

Component	Formula	Memory
<b>Model parameters</b>		
Backbone (1.8B params)	1.8B x 2 bytes	3.6 GB
MoE experts (450M params)	450M x 2 bytes	0.9 GB
GDiT + encoder (80M params)	80M x 2 bytes	0.16 GB
SigLIP (400M, frozen FP16)	400M x 2 bytes	0.8 GB
State engine (0.9M params)	0.9M x 2 bytes	0.002 GB
Router (0.5M params)	0.5M x 2 bytes	0.001 GB
<b>Subtotal: model</b>		<b>5.5 GB</b>
<b>Optimizer states</b>		
AdamW (2 moments per trainable param)	~2.1B x 2 x 4 bytes	16.8 GB
(SigLIP excluded --- frozen)		
<b>Subtotal: optimizer</b>		<b>16.8 GB</b>
<b>Activations</b>		

Component	Formula	Memory
Backbone (batch=8, seq=2048, 24 layers)	~6 GB (with gradient checkpointing)	6.0 GB
GDiT (batch=8, grid=30x30, 8 layers)	~0.5 GB	0.5 GB
SigLIP (batch=8, inference only)	~0.3 GB	0.3 GB
<b>Subtotal: activations</b>		<b>6.8 GB</b>
<b>Gradients</b>		
Same size as trainable params	~2.1B x 2 bytes	4.2 GB
<b>Subtotal: gradients</b>		<b>4.2 GB</b>
<b>Miscellaneous</b>		
Gradient checkpointing buffers	~1 GB	1.0 GB
CUDA workspace / fragmentation	~2 GB	2.0 GB
<b>Subtotal: misc</b>		<b>3.0 GB</b>
<b>Total training memory</b>		<b>~36 GB</b>
<b>Available (DGX Spark)</b>		<b>128 GB</b>
<b>Headroom</b>		<b>~92 GB</b>

## D.2 Inference Memory

Configuration	Model	KV Cache	GDiT	SigLIP	Total
BF16 (full precision)	5.1 GB	1.5 GB	0.16 GB	0.8 GB	<b>7.6 GB</b>
INT8 (quantized)	2.5 GB	1.5 GB	0.08 GB	0.8 GB	<b>4.9 GB</b>
INT4 (aggressive quant)	1.3 GB	1.5 GB	0.04 GB	0.8 GB	<b>3.6 GB</b>
INT4 + no vision	1.3 GB	1.5 GB	0.04 GB	0	<b>2.8 GB</b>
INT4 + no vision + short ctx	1.3 GB	0.4 GB	0.04 GB	0	<b>1.7 GB</b>

V3-mini at INT4 with reduced context length fits in 1.7 GB --- within the memory budget of a Raspberry Pi 5 with 8 GB RAM or any modern smartphone.

## D.3 Training Throughput Estimates (DGX Spark)

Phase	Samples/sec	Time per Epoch	Notes
Phase 1 (pretraining)	~800 tok/s	~25-30 days	BF16, full backbone
Phase 2 (alignment)	~200 samp/s	~1 hour/epoch	Only small modules train

Phase	Samples/sec	Time per Epoch	Notes
Phase 3 (SFT)	~50 samp/s	~4 hours/epoch	Full model, batch=8
Phase 5a (GDiT pretrain)	~300 grids/s	~6 min/epoch	GDiT only, small
Phase 5d (joint)	~30 samp/s	~4 hours	All components, batch=4
Phase 6 (self-play gen)	~5 tasks/s	Ongoing	Inference + GDiT + search

*This document presents the complete Engram V3 "Cognitive Mesh" architecture as of February 23, 2026. It combines proven techniques (MaskGIT discrete diffusion, DiT transformers, MoE, GRPO, world models) with novel integrations (state-conditioned denoising, unified generation/simulation, learned reasoning mode routing, metacognitive unmasking verification) to address the fundamental autoregressive bottleneck in spatial reasoning. The architecture is designed for two scales --- V3-mini (2.5B) as a research prototype on consumer hardware, and V3-full (190B) for competitive ARC-AGI performance --- with a clear training pipeline, cost estimates, and performance projections. The code is forthcoming. The architecture is designed. What it needs is training time and, for V3-full, cloud compute.*