**ChatGPT**

# 3. Composable Intelligence Systems: Binary Recursion, Structured Memory, and the Emergence of Cognition

*Abstract.* This chapter synthesizes the principles of *machine epistemology*, the operational primacy of the binary substrate, and the dynamics of computation into a theoretical framework for **composable intelligent systems**. We move beyond metaphorical descriptions of intelligence to propose an **operational definition** grounded in verifiable computation. Intelligence, in this framework, is characterized as a complex, adaptive process operating over structured binary memory, utilizing recursive functions, feedback mechanisms, and hierarchical domain-specific language (DSL) composition. Key elements – including dynamic memory architectures, entropy-aware learning via binary reinforcement, verifiable proof traces (ledgers of computation), and multi-DSL integration – are formalized. We argue that systems exhibiting these verifiable computational properties *realize intelligence*, defined as the capacity for adaptive, goal-directed, and structurally complex information processing fully traceable to the binary substrate (the base level of 0/1 state transformations). This positions the compositional framework as a formal theory of machine intelligence that enables construction, explanation, and verification of adaptive behavior across tasks, grounding "cognition" in recursive, feedback-regulated computation rather than any mystical attribute.

## 3.1 Redefining Intelligence: An Operational, Computational Perspective

Traditional definitions of intelligence often invoke nebulous concepts like consciousness, subjective experience, or unbounded abstract reasoning – ideas that are problematic for a rigorous, verifiable theory. **Machine epistemology** instead demands an operational definition grounded in observable computational capabilities. We therefore define intelligence *functionally*, in terms of what an intelligent system *does* in a measurable way, rather than what it subjectively *experiences*. As Marvin Minsky's *Society of Mind* insightfully suggests, "the mind" can be modeled as an emergent collection of simple, mindless processes ("agents") whose interactions produce intelligent behavior. In our framework, these simple processes reduce to binary operations on memory, and intelligence is the emergent result of their structured interaction.

**Proposed Operational Definition:** An **Intelligent System (IS)**, under this framework, is a computational system with a demonstrable capacity for:

1. **Structured Memory Management.** It maintains, accesses, and modifies persistent and volatile information encoded in *structured binary configurations* (collectively denoted as the system's memory *M*). This implies a **binary substrate** at the lowest level – all information reduces to binary data structures (bits) – and a memory architecture that imposes organization on those bits (for example, separating short-term and long-term memory, indexing knowledge, etc.). Modern AI research underscores that large AI models lack a *stable, structured memory* that endures over time [1]; an IS

addresses this by explicitly organizing memory content and making it available for ongoing computations.

2. **Recursive Self-Modification.** It employs *recursive* state-update functions whose behavior can adapt based on previous states, inputs, and internal feedback loops. Formally, we can express the system's state evolution as:

$$M_{t+1} = f\!\big(M_t,\; E_t,\; \Phi(M_t, E_t, \mathit{Goals})\big),$$

where $M_t$ is the internal memory state at time $t$, $E_t$ represents external input (environment or sensor data, also encoded in binary), $f$ is the primary state-transition function, and $\Phi$ is a *feedback function* (detailed later) that provides the system with information about its own performance or goals. This captures **binary recursion** – the system's next state is a function of its current state and input, forming a feedback-regulated recursive loop. Such *recurrent processing* is crucial: humans' ability to generate recursive patterns is considered a hallmark of higher cognition, and likewise an IS continually "feeds back" outcomes into its own cycle to refine future behavior.

1. **Entropy-Aware Adaptation (Learning).** It modifies its internal state transformations or memory structures to improve performance relative to defined objectives, guided by measures of information entropy or complexity. In intuitive terms, the system learns from experience by reducing uncertainty (or "surprise") in its predictions over time. For example, if the system's goal is to predict its environment, it should adjust its parameters to minimize prediction error entropy. This aligns with principles from cognitive science like the Free Energy Principle, which posits that brains update their internal models to reduce surprise or uncertainty. Our intelligent system similarly uses *feedback* $\Phi$ to detect prediction errors, inconsistencies, or inefficiencies (high entropy), and then adapts to **optimize an objective** (e.g. minimize error or maximize reward). We will formalize in §3.4 how this learning process can be seen as an optimization problem carried out on the binary substrate.

2. **Hierarchical Compression & Abstraction.** It builds and utilizes *layered internal DSLs* (domain-specific languages) to represent patterns, manage complexity, and generalize knowledge across domains. This means the system develops higher-level symbolic or statistical representations (rules, features, etc.) that compress raw data into more *abstract* concepts. Each layer is like a new language for a particular domain of patterns, grounded in lower-level binary operations but providing a more compact way to reason about those patterns. In mathematics, for example, one can see a hierarchy from binary logic up through arithmetic, algebra, analysis, etc., each layer offering new constructs but ultimately executable in binary [2] [3]. An intelligent machine likewise might have, say, a sub-module that understands geometric shapes and another that handles logical inference – each effectively a "DSL" for that aspect of the world. **Compression** here refers to reusing lower-level computations and summarizing repeated structures so that higher-level reasoning doesn't reinvent bit-level operations each time [4]. This hierarchical abstraction is essential for coping with complexity: it allows the system to apply prior knowledge to new situations by invoking the appropriate high-level "language" internally. We will see in §3.6 how multiple DSLs can be orchestrated within one system.

3. **Verifiable Transformation and Reasoning.** It produces *traceable evidence* of its computational steps and state changes, enabling external verification that it operates according to its specifications (the rules of its DSLs and the binary substrate's logic). In other words, every significant decision or

adaptation the system makes can be checked after the fact. This could be through **proof traces** (e.g. a log of state hashes and operations) or other metadata that serve as a *computational ledger*. Verifiability is a cornerstone of machine epistemology – it bridges the gap between mere complex behavior and trustworthy intelligence. For an AI system deployed in the real world, such as an autonomous vehicle or a medical diagnosis assistant, being able to audit its decision process is crucial for trust and safety. Our framework therefore treats an intelligent system as not only a black-box predictor but as a **transparent reasoner** that can justify its outputs by reference to an explicit chain of binary operations and transformations. We will detail in §3.5 what form these traceability mechanisms can take (e.g., logging each update with cryptographic hashes, akin to a blockchain ledger for AI inferences).

Crucially, this definition brackets out *phenomenal consciousness* or any metaphysical aspect; we focus strictly on the functional architecture required for complex, adaptive, **verifiable** computation. Intelligence here is not an immeasurable spark – it is the cumulative capability of a system to process information in a structured, goal-directed way and to demonstrate the results (and integrity) of that processing. This operational perspective ensures that if we build a system meeting these criteria, we can rightfully *call* it intelligent within our theory. It also aligns with how cognitive architectures are evaluated in AI research, by their abilities to handle memory, learning, perception, etc., in a manner comparable to human-like intelligence [5].

Before diving into each of these elements in depth, we reiterate the grounding of all these capabilities: the **binary substrate**. At the lowest level, our intelligent system is implemented on binary hardware or equivalent – all high-level behaviors ultimately consist of binary state transitions. In practice, this might be an actual digital computer or a mathematical abstraction of one. The *Binary Decomposition Interface (BDI)* is the conceptual machine layer that ensures any high-level DSL operation (a logical inference, a matrix multiplication, etc.) decomposes into a sequence of binary operations. The BDI thus acts as the bridge between abstract algorithms and bit-level execution [6]. Throughout this chapter, whenever we discuss memory $M$, functions $f$, or data structures, it is understood that the BDI can in principle trace these back to binary manipulations – which is why verification (point 5) is feasible down to the substrate. With this foundation in place, we now explore each aspect of the definition in detail.

## 3.2 Dynamic Binary Memory: The Substrate for Cognitive Processes

An intelligent system requires more than a monolithic memory space; it needs a *dynamic, structured memory architecture* that supports both the stability and plasticity required for cognition. We formalize the system's memory at time $t$ as $M_t$, which is composed of a set of **Binary Memory Regions (BMRs)**:

$$M_t = \{R_1^t, \ R_2^t, \ \ldots, \ R_n^t\},$$

where each $R_i^t$ is a region containing some subset of the binary substrate (i.e. a block of bits) along with associated metadata that characterize its role. These memory regions are the building blocks of the system's knowledge base and working storage. By imposing structure (dividing memory into regions with specific purposes and properties), the system can manage information in a cognitively meaningful way.

Each BMR $R_i$ may be defined by attributes such as size, mutability, decay rate, and semantic tags. We can conceptually classify memory regions into functional **classes**:

- **Persistent Memory.** Regions with high resistance to change (effectively read-only or very infrequently modified). These store the core knowledge or programming of the system – foundational axioms, long-term facts, or innate skills. In a human analogy, this is like deeply ingrained knowledge or instinct; in computers, it's like firmware or ROM. Marking certain binary regions as persistent ensures that crucial knowledge isn't overwritten by noise. In our framework, this could include the base DSL definitions or fundamental constants the system uses. Persistent memory provides *continuity* of identity and competence over time.

- **Volatile Memory.** Regions with rapid turnover and high mutability. These are scratch-pads for ongoing computations – the equivalent of RAM or a human's working memory. Sensory inputs $E_t$ might first be recorded in volatile buffers, intermediate results of calculations reside here, and they may be erased or replaced frequently as the system moves from one task or thought to another. Volatile memory enables real-time responsiveness and context-specific processing, but it isn't relied upon for long-term retention. It's dynamic and can be flushed or repurposed as needed.

- **Adaptive Memory.** Regions that are explicitly designed to be modified by learning processes. These hold the system's *plastic* knowledge – parameters of functions (like weights in a neural network), episodic memories of past events, or learned heuristics. Adaptive memory regions are tagged for the system to update when feedback $\Phi$ indicates a change is beneficial. For example, if $f$ is a function with parameters $\alpha$ (as introduced in §3.1), those parameters would live in an adaptive region $R_{\alpha}$ within $M$. Over time, this region's contents change (in binary form) to encode new skills or refined models. This is akin to the synaptic weights in a brain or the adjustable settings in a learning algorithm.

By organizing $M_t$ into such categories and perhaps further substructures (like key–value stores, graphs, etc.), the system gains **structured access** to its memory. When a computational process (function $f$ or feedback $\Phi$) needs to read or write memory, it can target the appropriate region type, which brings semantic advantages. For instance, a learning algorithm might only adjust *adaptive memory* and leave persistent memory untouched, or a garbage collection process might periodically clear out volatile memory to free space.

We emphasize that all these memories are still ultimately binary. The distinction lies in how the system *uses* them. For example, a "memory region" could literally be a contiguous block of bits tagged as "volatile working memory," and the BDI would treat operations on that block (reads/writes) in a way that doesn't require permanence. By contrast, another block of bits tagged "persistent" might be stored in a protected area or duplicated for safety. These implementation details can vary, but conceptually, $M_t$ is a richly structured binary state.

**Memory Dynamics:** The transformation function $f: M_t \to M_{t+1}$ operates on these structures. Rather than treating memory as an undifferentiated tape, $f$ can have sub-functions specialized for each region type. For example, $f$ might include a routine for consolidating new information from volatile to persistent memory if something important was learned (analogous to memory consolidation during sleep for humans). The presence of metadata on each region (what kind of information it holds, how fast it should change, etc.) allows the system to enforce constraints. The **BDI** plays a vital role here: it interprets high-level

memory operations (like "retrieve item X from memory" or "increase parameter Y slightly") and translates them into low-level binary read/write actions on the correct region [7] . The BDI essentially acts as the *memory controller* for our intelligent architecture, ensuring that DSL-level instructions respect the structure of memory. Because of this, every memory access or modification can also be logged or checked, contributing to traceability (the verifiable aspect).

In effect, **thinking** in this system can be viewed as the *active transformation* of the memory state through time. At any given step, some combination of regions is read (sensory input might go into a volatile region, current beliefs might be fetched from persistent memory, etc.), computations happen (perhaps in CPU registers or an equivalent, but conceptually within $f$), and then a new memory state is written, affecting perhaps the adaptive regions (learning updates) or generating motor outputs (writing to an output interface region). This continual process $M_t \to M_{t+1}$, when properly structured, produces what we recognize as cognitive activity: perceiving, reasoning, updating knowledge, and so forth [8] . The hierarchy of DSLs (to be discussed in §3.6) lives *within* this memory as well – e.g., a complex structured memory might contain a graph of symbolic relationships as part of a knowledge DSL, or a trained neural network's weights as part of a perception DSL.

It's worth noting that cognitive architectures in AI frequently distinguish memory systems akin to what we describe: for instance, the ACT-R architecture has separate modules for declarative memory (facts) and procedural memory (skills), and working buffers [5] . Our framework is general, but by introducing BMRs and classes, we are in line with these proven designs. Additionally, recent research on large language models (LLMs) identifies the lack of a long-term memory as a major shortcoming [1] , and efforts are underway to augment LLMs with vector databases or episodic memory records. The advantage of our approach is that memory is an *intrinsic* part of the system's state $M_t$ and is manipulated by the same recursive logic $f$ – rather than an external add-on, it is integrated and thus fully auditable and subject to the same formal rules.

In summary, **dynamic binary memory** $M_t$ is the canvas on which intelligent computation is painted. By structuring this canvas into regions and allowing it to evolve with both stability (some parts hardly change) and plasticity (other parts learn), we provide the substrate for complex cognition. Intelligence will emerge from the interplay of processes and memory, which we examine next in terms of recursive feedback loops.

## 3.3 Recurrence, Feedback, and Self-Regulation: The Engine of Adaptation

A system limited to one-off computations – where each input is processed in isolation to produce an output – cannot adapt or exhibit goal-directed behavior over time. Intelligence requires *recurrence*: the ability for the results of prior computations to influence future ones. This is accomplished through **feedback loops** that continually update the system's state. In our formulation, feedback is represented by the function $\Phi(M_t, E_t, \mathit{Goals})$ feeding into the state transition $f$. We now unpack how this works and why it is central to self-regulation and learning.

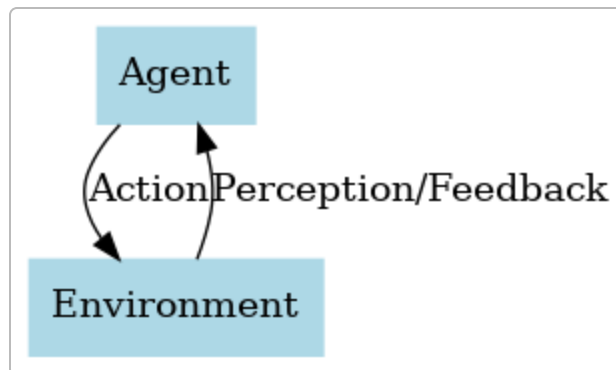At each time step, as given earlier, the system computes:

$$M_{t+1} = f\big(M_t,\ E_t,\ \Phi(M_t, E_t, Goals)\big).$$

Here $\Phi$ is the **feedback function**. Think of $\Phi$ as the system's *self-reflection mechanism*: it looks at the current state and input (and knowledge of any explicit goals) and produces a feedback signal that influences how $f$ will produce the next state. We can break down typical forms that $\Phi$ might take:

- **Error Signals.** $\Phi$ can compare the system's predictions to reality. For example, if the system predicted some internal or external event (perhaps $M_t$ contains a prediction made at $t-1$ about what input would come at $t$), $\Phi$ can compute the *error* between the prediction and the actual outcome. This is common in control systems and predictive coding models – an internal model generates an expectation, and any deviation (error) is a feedback used to update the model. In our notation, $\Phi$ might output something like $\text{error} = M_t^\text{predicted} - M_t^\text{actual}$, which $f$ then uses to correct $M_{t+1}$.

- **Reward Signals.** If the system is goal-driven (as most intelligent agents are), $\Phi$ can evaluate how well the current state/action meets the goal and provide a scalar reward or penalty. This is the basis of reinforcement learning: the environment (or an internal critic) gives a reward $R_t$, and the system adjusts its policy to maximize cumulative reward. In a binary substrate, even a scalar reward would be encoded in bits (often as a short binary string). For example, $\Phi$ could output $R_t = 1$ (success) or $0$ (failure) in a simple binary reward scheme, or more complex gradations. The presence of $R_t$ in the feedback influences $f$ to produce $M_{t+1}$ that increases the likelihood of getting 1's and avoids 0's in the future ⑨ ⑩ .

- **Entropy or Complexity Monitors.** $\Phi$ can measure changes in the entropy $H(M_t)$ of the system's state or other complexity metrics $K(M_t)$ (e.g., Kolmogorov complexity or description length of the state). This is a more abstract form of feedback where the system self-regulates to maintain "homeostasis" in terms of complexity. For instance, if the system's state is becoming too random (high entropy) or too uncertain, $\Phi$ might signal this so that $f$ adjusts $M_{t+1}$ to be more stable or more compressed. Conversely, if the system is over-confident or rigid (too low entropy), $\Phi$ might inject some exploratory randomness. This is analogous to mechanisms in simulated annealing or curiosity-driven learning where the aim is to balance exploitation and exploration. Friston's free-energy principle in neuroscience is along these lines: minimize surprise (entropy) but keep enough variability to adapt.

- **Consistency Checks and Constraints.** The feedback function can also enforce logical or structural constraints. For example, if part of $M_t$ represents beliefs or a knowledge graph, $\Phi$ might include a module that checks for contradictions or violations of rules (perhaps using an internal logic DSL). If a contradiction is found, $\Phi$ generates a signal that causes $f$ to resolve it (maybe by retracting some belief or flagging an uncertainty). This ensures the system's state remains *self-consistent*. In formal AI systems, this relates to truth maintenance systems or constraint solvers that keep track of dependencies and consistency.

The output of $\Phi$ can be thought of as a **vector of signals** or simply additional variables that augment the input to $f$. For simplicity, our notation shows $\Phi(\cdot)$ feeding into $f$ as a single composite signal, but in practice it could be multiple feedback channels (error, reward, etc. simultaneously). All these feedback signals are themselves binary-encoded in $M_t$ or as part of the input; the BDI ensures they are represented in bits and that responding to them is just more binary computation.

The significance of this recursive, feedback-driven architecture is that it creates a *closed loop* between the system and its environment (and between the system and itself). This is often illustrated as an **agent-environment loop**, where the agent (system) takes actions that affect the environment and in turn receives observations and rewards from the environment. Figure 3.1 provides a schematic of such a feedback loop in a simple agent: the agent observes the environment and internal state, computes actions, the environment responds, and the cycle continues.



In our context, the "action" the system takes could be an external action (affecting $E_{t+1}$) or an internal action (a reconfiguration of $M_{t+1}$). Either way, *history matters*: what happened at time $t$ influences what the system does at $t+1$. This enables **adaptation**.

<small><em>Figure 3.1: A basic agent-environment feedback loop.</em> The agent's action influences the environment, and feedback (perception of the new state or an explicit reward) influences the agent's next action. Intelligent systems generalize and internalize this loop, using internal feedback signals $\Phi$ to adjust their own state.</small>

Realize that this feedback-driven recurrence is the engine behind learning: the system can **self-correct** over time. If an action leads to a bad outcome (as indicated by $\Phi$), the next state $M_{t+1}$ will be different than it would have been otherwise, steering the system away from repeating the mistake. Conversely, good outcomes reinforce the behaviors or state changes that led to them. This aligns with decades of control theory and reinforcement learning research which shows that feedback is necessary for an agent to function robustly in a dynamic world [11] [12]. Without feedback, the system is essentially open-loop – it would not know if its decisions were right or wrong, safe or dangerous.

To ground this in a more cognitive sense: a human or animal uses feedback from the world (success, pain, surprise, hunger, etc.) to guide future actions. Similarly, a BDI-based intelligent agent uses an internal belief-desire-intention cycle with feedback. The classic **BDI architecture** indeed incorporates feedback: beliefs are updated from sensory inputs, options (plans) are generated and filtered (feedback from intentions and results may lead to dropping or revising intentions), forming a control loop between what the agent *believes*, *wants*, and *intends*. An example BDI agent loop is depicted in Figure 3.2, where perceptions update beliefs, which along with desires generate options, intentions filter actions, and actions affect the world that feeds back into new perceptions

. Our framework is more general in that $f$ and $\Phi$ need not correspond exactly to belief and desire functions, but one can see BDI as a particular structured instantiation of this principle, with $\Phi$ implementing things like intention reconsideration or goal success checks.

<small><em>Figure 3.2: Belief-Desire-Intention (BDI) agent architecture (schematic).</em> This diagram (adapted from BDI models) shows how sensory inputs update Beliefs (via a belief revision function, BRF), which together with Desires feed into option generation. Intentions are chosen and executed as Actions, affecting the environment, and a feedback filter updates desires and intentions. This is an example of a structured recursive loop within an agent. In our framework, such loops are implemented via the general $f$ and $\Phi$ functions on binary memory, but the principle of feedback-regulated state update is the same.</small>

In conclusion, **self-regulation** in an intelligent system arises from these recurrent feedback loops. The system is not a static program but an ongoing process that *regulates itself* based on its history and goals. This provides stability (through negative feedback, e.g. error correction driving the system toward target states) and flexibility (through positive feedback or exploratory signals that allow it to try new actions). It is the fundamental mechanism that enables **learning and adaptation**, which we turn to next.

## 3.4 Learning as Entropy-Guided Optimization: Meta-Learning and Binary Reinforcement

Learning in our framework is the process by which the system *modifies its own transformation logic* (the function $f$ or its parameters) or its adaptive memory structures to improve future performance. Thanks to the recurrent feedback loop described in §3.3, the system has the capability to observe the outcomes of its actions and update itself. Now we formalize this in terms of an optimization problem guided by feedback, and connect it to established learning paradigms like gradient-based learning and reinforcement learning – all while grounding them in binary computation and verifiable updates.

First, consider that the state transition function $f$ is not necessarily a fixed piece of code. It may be **parameterized** by some internal parameters $\alpha$. For instance, $f$ could be a neural network with weights $\alpha$, or a decision tree with adjustable thresholds, or any algorithm with tunable settings. These parameters reside in adaptive memory (as discussed, $\alpha \in M$ in an adaptive region). We can denote this explicitly as $f_\alpha$ to indicate $f$ depends on $\alpha$. The system's task is to adjust $\alpha$ over time to better achieve its goals. This is **meta-learning** – the system learning how to learn or adjusting its own algorithms.

Let's define an objective function $J$ that quantifies performance. $J$ could be something like "negative error entropy" (so maximizing $J$ means reducing uncertainty or errors) or cumulative reward, or accuracy on predictions, etc. Crucially, $J$ can be evaluated (or at least estimated) from the feedback signals. For example, if the feedback provides an error $\epsilon_t$, $J$ might be $- \mathbb{E}[\epsilon^2]$ (to maximize means to minimize mean squared error), or if feedback provides rewards, $J$ could be the expected sum of rewards. Because the system runs online, we consider $J$ as a function of the current parameters and the current feedback: $J(\alpha_t; \text{Feedback}_t)$. The feedback at time $t$ encapsulates how the system did on that step (error, reward, etc., possibly aggregated).

A straightforward approach is to adjust $\alpha$ in the direction that improves $J$. If we imagine $J$ as a differentiable function of $\alpha$ (a common assumption in training neural networks), one would compute a gradient $\nabla_{\alpha} J$ and then do an update like:

$$\alpha_{t+1} = \alpha_t + \eta \, \nabla_\alpha J(\text{Feedback}_t),$$

for some learning rate $\eta$. In practice, $\nabla_{\alpha} J$ might be estimated rather than exact. Our framework accommodates both exact and approximate gradients. Importantly, **nothing leaves the binary domain** here: $\alpha_t$ is stored in binary, the gradient is computed via binary operations (even if conceptually using calculus, in implementation it could be via finite differences or bit-wise arithmetic), and the update addition is binary arithmetic. The *verifiability* angle means we could log $\alpha_t$, $\nabla J$, etc., all as part of a proof trace (to ensure the update rule was followed correctly).

However, not all systems or objectives are nicely differentiable. Many intelligent systems learn via more discrete or heuristic means. Let's break down learning mechanisms highlighted by our design:

**Gradient-Based Meta-Learning:** When possible, using gradients is powerful. As an example, consider a neural network inside our system predicting some outcome. $\Phi$ provides a *prediction error* signal. The system can use finite difference methods to estimate gradients: perturb $\alpha$ by a small $\delta$ in a few random directions and see how the error $J$ changes [13] . This is essentially how one might do gradient-free optimization, but in expectation it can approximate the true gradient. Alternatively, the system might accumulate correlations: if a certain parameter's change consistently reduces error, reinforce that change. This is similar to evolutionary strategies or random search in a binary space of parameters. Because $\alpha$ is stored in binary, even gradient descent ultimately boils down to flipping bits in $\alpha$ in the correct direction (incrementing or decrementing numeric values, etc.). We can formally describe one update as in the excerpt of pseudo-code:

> **UpdateRule**($\alpha_t$, Feedback$t$):
> *Compute $\Delta \alpha$ such that $J(\alpha_t + \Delta \alpha) > J(\alpha_t)$ using estimated $*

$\nabla J$ (via internal analytic or empirical means).
$\alpha_{t+1} \gets \alpha_t + \Delta \alpha$.

This is kept abstract, but the key is $\Delta \alpha$ is derived from feedback. In a verifiable system, one could require that $\Delta \alpha$ is actually the argmax step for $J$ given some constraints, and that could be checked by re-evaluating $J$ after the fact (if $J$ is known or at least whether $J$ improved can be logged).

**Binary Reinforcement Learning:** A lot of real-world learning, especially in autonomous agents, comes down to trial-and-error reinforced by rewards. In a binary context, we can imagine the simplest reward $R_t \in \{0,1\}$ (failure/success) at each step, which is part of $\Phi$. The system can implement a reinforcement learning update. For example, *reward prediction error* $\delta_t = R_t - \hat{R}_t$ could serve a similar role as a gradient: if $\delta_t$ is positive (we got a better reward than expected), then the actions leading to it (which depend on $\alpha$) should be strengthened. If negative, they should be weakened. The system could maintain an **eligibility trace** in memory – essentially a trace of which parameters were responsible for recent actions (a common technique in RL to assign credit). Then an update rule might be:

$$\Delta \alpha = \eta \, \delta_t \, \mathrm{EligibilityTrace}_t,$$

so that if $\delta_t$ is 1 (unexpected success), $\alpha$ moves in the direction of the parameters that were active (increasing probability of repeating those moves), and if $\delta_t$ is -1 (unexpected failure), $\alpha$ moves opposite to those parameters [10] . Over time, this evolves a policy that yields more rewards.

Another binary RL approach: compare recent trajectories that led to reward vs. those that didn't. If we have two copies of parameters, one slightly perturbed, and one did better ($J^+$) than the other ($J^-$), then:

$$\Delta \alpha \propto (J^+ - J^-) \cdot \mathrm{ParameterPerturbation},$$

which is basically the idea of evolutionary strategies (if a random change improved things, keep that change) [14] .

Both gradient-based and reinforcement-style updates should be seen as implementations of a broader principle: **entropy-guided optimization.** The system uses feedback to reduce the entropy of its errors or to increase the predictability of getting rewards. Many algorithms in machine learning can be cast as minimizing some loss function (which is high when errors are unpredictable or outcomes are bad) – essentially an information gain. For example, in decision tree learning, splitting on a feature is done by information gain (reducing entropy of target variable), which is directly analogous to our system choosing to organize memory or branch logic to reduce uncertainty.

A key strength of our framework is **grounding and verifiability of learning**. In typical AI systems, learning algorithms (like backpropagation in a neural net) are often treated as complex black boxes that adjust thousands of parameters in ways hard to follow. Here, every parameter $\alpha$ lives in memory, and every update is a transaction in the system's ledger (as we'll discuss in §3.5). This means we can *trace* what learning step happened when, and even require proof that a given update reduced error on the last trial or followed the specified rule. In sensitive applications (finance, healthcare), such auditability of AI learning is increasingly seen as necessary.

From the perspective of binary implementation: all these learning methods eventually boil down to bit-level operations. For example, floating-point weight updates in backpropagation involve binary arithmetic on memory registers (the CPU/GPU does it). In our theoretical BDI machine, those would be broken into low-level opcodes, but we can also imagine using higher-level *learning primitives*. Perhaps the BDI provides an operation `LEARN_APPLY_DELTA(region, delta)` which adds a binary number `delta` to all values in a memory region (useful for vectorized weight updates) [15] . Or `LEARN_RECORD_GRADIENT` which could store a computed gradient in a designated area [16] . The advantage of having such primitives is twofold: (a) efficiency (the system can learn faster if the substrate is optimized for it), and (b) verifiability (the operations are known and limited, making it easier to prove correctness). Many formal approaches to machine learning verification assume a fixed update rule and then verify each step respects it [17] [18] .

In summary, **learning** in an intelligent system is realized as an internal optimization process guided by the feedback signals $\Phi$. Whether through explicit gradient descent, heuristic credit assignment, or simple trial-and-error, the system **reduces its future entropy** – it becomes more certain in achieving goals, more precise in its predictions, and more efficient in its behavior. This can be seen as the system continuously *reprogramming itself* in light of experience. We call it *meta-learning* to highlight that the system is changing how it computes ($f$ and its parameters), not just what it computes. The result is an ever-improving alignment between the system's internal models and the external tasks it faces, all encoded in binary and done in a way that can be inspected and traced.

With structured memory (§3.2) and recursive feedback loops (§3.3) providing the stage for learning (§3.4), we have covered the core of adaptive computation. We now shift focus to another critical aspect of our framework: **verification and traceability**, the epistemological bedrock that ensures we *know* what the system is doing and that its emergent intelligence remains accountable.

## 3.5 Verification and Traceability: The Epistemology of Machine Intelligence

A cornerstone of our approach is that an intelligent system must not only *act* intelligent but also provide **evidence** for the validity and provenance of its actions and adaptations. In other words, it must be intelligible. This addresses the oft-cited "black box" problem in AI: we require our system to be a *glass box*, logging its internal processes in a robust way. This section describes how we achieve verifiability through **proof traces** and what benefits this brings.

Every significant computational step in the system can be recorded as a **trace entry** in a *computational ledger*. Conceptually, think of this ledger as an ever-growing log file or blockchain that the system writes to as it operates. Let's denote a dedicated persistent memory region $M_{\text{ledger}} \subset M$ that stores these trace records. Because $M_{\text{ledger}}$ is part of persistent memory, it's not overwritten or erased easily – it provides a historical memory of the system's operation.

What information should a trace record contain? At a minimum, we want it to uniquely identify the operation and allow reconstruction or verification of what happened. An entry could include fields such as:

- **Timestamp (t):** the time or step number when the operation took place.

- **Operation ID:** an identifier of what function or rule was applied. For instance, if the system executed a specific primitive operation of a DSL, or a learning update, this ID labels it (could be a hash of the code or a mnemonic).
- **Input State Hash:** a cryptographic hash (e.g., SHA-256) of the relevant portions of $M_t$ and input $E_t$ before the operation. This acts as a fingerprint of "pre-operation" state.
- **Output State Hash:** a hash of the resulting state $M_{t+1}$ (or the parts that changed) after the operation.
- **Parameter Snapshot:** if the operation involves parameters (like $\alpha$ in a learning update), a digest or relevant subset of those parameters pre- and post-operation. This might be included in the input/output hashes, but it can be highlighted separately if needed.
- **Feedback Digest:** a summary of the feedback $\Phi$ at that step (e.g., the error or reward value, or a hash of the feedback vector).
- **Entropy/Metric Change:** an optional field recording $\Delta H$ or $\Delta J$ – how some metric (entropy, error, reward) changed due to this operation.
- **Proof Signature or Chain:** a cryptographic link (like a hash pointer to the previous record) to ensure the ledger's integrity (similar to how each block in a blockchain references the previous block's hash). Optionally, a digital signature if we want to prove the record came from the legitimate system and was not tampered with.

Such a record is essentially a *structured log of computation*. By storing the InputHash and OutputHash, we don't store the entire state (which might be huge), but we commit to it. If later we need to audit, we can compare actual states to these hashes to ensure nothing was falsified. The chain/hash linking ensures that the sequence of operations is tamper-evident (you can't remove or reorder trace entries without breaking the chain).

**Benefits of Traceability:**

- **Auditability:** External observers or regulators can inspect the ledger to reconstruct the sequence of decisions. For example, if the system made a medical diagnosis at 12:00 and we want to know why, we can find the trace at that time, see which inputs (symptoms, tests) were considered (via InputHash matching known input data), what internal rule fired (OperationID could indicate a certain guideline or neural network forward-pass), and confirm that the outcome was computed according to the system's design. This is akin to an airplane's black box but for AI reasoning. It addresses calls in AI policy for continuous auditing and logging of automated decisions.

- **Explainability:** While raw logs are not the same as high-level explanations, they provide the *ground truth data* from which explanations can be derived. Given a trace, one can work backwards: if an output was surprising, check the feedback, see what the system was trying to optimize, and perhaps realize it did exactly what it was told (but maybe the objective was mis-set). This helps in explaining both correct and incorrect behavior. For instance, a trace might show a high reward was given to a certain action that looks wrong in hindsight – explaining that the system was following its reward signal (which might have been mis-specified by designers) is important. In essence, the trace is a *causal chain* of evidence for each decision [19] [20].

- **Debugging and Improvement:** Developers of the intelligent system can use the ledger to pinpoint where things went astray. If a particular parameter update caused performance to degrade, the trace will show the before/after and the feedback that led to it. This is extremely valuable given the

complexity of adaptive systems. Instead of guessing why the system behaves a certain way, the engineer can inspect the logs. In formal methods terms, this is like having a proof of each state transition that can be checked if a safety property was violated.

- **Proof of Learning and Compliance:** In safety-critical domains, we might need the system to prove it *learned correctly*. For example, regulations might require that an AI driving a car only learns from legitimate data and doesn't "drift" into unsafe behavior. The ledger can be used to verify that each learning step was according to the approved algorithm (the OperationID could be cross-checked to ensure it's the allowed `UpdateRule`), and the feedback came from valid signals (perhaps signed by a secure module or external oracle). This provides a chain of evidence that the AI's current model is the result of valid updates. If the AI Act or other regulations demand that AI decisions be traceable and auditable, our system's design inherently satisfies that by design.

Implementing this ledger incurs overhead, but since we are primarily describing a theoretical framework, we assume the environment can handle it (one could also throttle the level of detail based on necessity – e.g. log only abstract operations instead of every bit operation). Notably, the **BDI** architecture we rely on is well-suited to assist traceability. Because the BDI decomposes high-level operations into primitive ones, it can insert logging at the appropriate granularity [21] . For instance, the BDI might automatically produce a trace entry each time a DSL operation completes or a memory region is updated. The BDI ensures that nothing "escapes" logging – since all computations pass through it, it can annotate them.

There is a parallel here with how blockchain systems ensure transparency and trust through public ledgers, and indeed some proposals suggest using blockchains to log AI decisions for exactly these reasons. In our case, the ledger could be internal or external. An external immutable ledger (like a blockchain) could even record hashes from multiple agents, enabling third-party verification. For a single system, an internal ledger with secure logging might suffice.

To sum up, **verification and traceability** transforms our intelligent system from a purely computational entity into an *epistemic agent* – one whose knowledge and operations are not only sound but also *knowable* and checkable by others. This satisfies the philosophical stance of machine epistemology: that knowledge (and intelligent behavior) is only meaningful if it can be verified [6] . We have now addressed all aspects of our operational definition of intelligence. The final two sections will discuss how these components come together in a holistic architecture and how intelligence emerges from their interplay, before concluding with the significance of this formal approach.

## 3.6 Compositionality: Intelligence via Multi-DSL Integration

Thus far we have described the architecture mostly in general terms. We now consider the **compositional** aspect of intelligence: the ability of a system to integrate multiple forms of reasoning and representation. Human intelligence is notable for its diversity – we can perform mathematical logic, visualize spatial scenes, recall episodic memories, interpret language, etc., each of which might require different methods. It is unlikely that a single monolithic algorithm (a single DSL) can capture all these abilities [22] . Instead, an intelligent system must orchestrate a *society of specialized DSLs*, each acting as a module for a certain domain, yet all ultimately unified on the binary substrate via the BDI.

**Necessity of Multiple DSLs:** Consider a complex task like robot navigation in a household. The robot might need: arithmetic/geometry to compute distances, a logic module to reason about which rooms to clean first, a vision module (perhaps a neural network) to recognize objects, and a natural language module to understand instructions from a human. Each of these can be seen as operating in a different formal language: geometry, logic, deep neural nets (pattern recognition), and linguistic grammar. Rather than reinventing a single algorithm that does it all, it makes sense to have dedicated sub-systems for each and have them *communicate*. Cognitive architectures research indeed often results in **hybrid architectures** that combine symbolic reasoning with sub-symbolic perception modules [22] [5] .

In our framework, we model each such capability as a **DSL implemented within the system**. By DSL, we mean a formal subsystem with its own data representations and operations, optimized for a particular kind of problem, yet defined in such a way that it compiles down to binary operations (and thus can be supervised by the BDI). In Chapter 1, the book discussed mathematics as an ecosystem of DSLs from logic up to topology [23] [24] . Here we apply the same thinking to *intelligent behavior*.

**Architecture for Multi-DSL Integration:** The intelligent system acts as a **meta-system** that *hosts* multiple DSLs and coordinates them. Specifically, the architecture supports:

1. **Hosting Multiple DSLs:** The system's memory $M$ can store representations from different domains simultaneously – e.g. a pixel map for vision, a symbolic knowledge graph for logic, a vector of numbers for a calculus simulation. Likewise, the BDI can host multiple interpreters or compilers, one for each DSL, that know how to execute operations in that domain. For example, one interpreter might handle logical inference rules on symbolic data, while another handles matrix multiplication on numeric arrays. All of these interpreters ultimately manipulate the underlying binary memory (reading/writing bits in structured regions corresponding to their domain data). The *Chimera* sub-architecture mentioned in Chapter 9 hints at such an ability to compile different structures into BDI nodes [25] [26] . In essence, the system contains *modules*, but unlike traditional black-box modules, here they are defined as DSLs which are transparent and formally specified.

2. **Selecting the Appropriate DSL(s):** At any given time or for any given subtask, the system should invoke the reasoning mode that fits. If the input $E_t$ indicates a visual pattern, the system might engage the vision DSL (e.g., run a convolutional network) rather than the logic DSL. Conversely, if asked a planning question, it engages the symbolic planner DSL. This selection can be handled by a dispatcher component (itself part of $f$) that examines the context and routes the processing. Formally, we could have a high-level policy $\pi$ that maps $(M_t, E_t)$ to a choice of DSL or a sequence of DSL operations. This is akin to an operating system scheduler, but at the cognitive level. The selection might also be made in parallel – multiple DSLs could be active concurrently, each handling the part of state relevant to it.

3. **Translating Between DSL Representations:** One challenge with multiple specialized modules is **interoperability** – how do the results of one module inform another? Our solution is that the system can learn or be programmed with *transformation functions* $f_{\text{translate}}$ that convert data from one representation to another [27] . For instance, if the vision DSL identifies an object and labels it "cup" (perhaps as an embedding vector), a translation function could turn that into a symbolic fact in the logic DSL: `object(cup)` in working memory. Or converting geometric coordinates (from a mapping module) into algebraic equations (for a physics solver) [28] . These translations are themselves part of the DSL skill set of the system. They must be well-defined so that no meaning is

lost – which implies that for every important concept, the system has a shared semantic grounding across DSLs. In practice, building such bridges is a known challenge, but approaches like knowledge graphs or common ontologies can help. The BDI ensures that even translation is a series of binary ops, hence traceable.

4. **Orchestrating Multi-DSL Workflows:** The system needs an *executive function* or cognitive controller to manage complex tasks that require multiple steps across different DSLs. For example, solving a physics word problem: the system might parse the text (language DSL), formulate equations (algebra DSL), solve them (numerical DSL), then perhaps double-check units (logic DSL). The orchestration means invoking these in the correct sequence and feeding outputs of one as inputs to another. This orchestration can be hardcoded as strategies or learned via meta-reasoning. Our recurrent feedback architecture aids this: the system can plan a sequence in its memory (like a script of which modules to call when, analogous to a *plan*), then execute step by step, getting feedback at each sub-step, and adjusting if needed. The *traceability* of each step is also crucial here, because if a final answer is wrong, we can see which module might have erred. In many AI frameworks, this kind of pipeline is managed by an external program (like a Python script calling different libraries). Here, it's integrated within the intelligent agent itself, giving it autonomy to decide its workflow. This resonates with the concept of cognitive architectures having a central **cognitive cycle** (like in SOAR or ACT-R) that can incorporate different types of processing per cycle [29].

All DSL modules share the **binary substrate** and the BDI. This commonality is what makes integration possible. Since all data is ultimately just bits, any module's output can be passed as input to another, provided the receiving module can interpret that bit pattern correctly. By defining interface standards (e.g., the translation functions), we ensure modules speak to each other. This is comparable to how in software engineering, different subsystems use common data formats (like how an image processing library and a machine learning library might both use NumPy arrays in memory to exchange data). In our case, the "common format" at the lowest level is bits with accompanying schema tags indicating the type.

**Example Scenario:** To make this concrete, imagine the system is solving a physics problem: "If you push a box up a frictionless incline, how much work is done?" The system might: parse the text with an NLP DSL to extract a formal representation of the problem (e.g., mass, distance, angle as variables); use a symbolic algebra DSL to derive the equation for work (maybe $W = m g h$); use a calculus/physics DSL to relate $h$ to distance and angle ($h = d \sin\theta$); then plug in numbers and compute using a numerical DSL. It might then use a visualization DSL to illustrate the scenario (if required) or a logic DSL to verify that the answer makes sense (no conservation law violated, etc.). Each of these steps is performed by a specialized module but under the governance of the top-level intelligent system. The correctness of the entire solution can be verified because each module's operation is logged and the interfaces between them are clear (units and values translated properly). This combinatorial use of multiple representations is known to be a feature of human expert problem solving [30], and our system can emulate that approach.

By integrating multiple DSLs, the system avoids being limited by any single method. This addresses a key criticism of "massively modular" views of the mind, which is that modules alone can't flexibly share information [31]. Our design counters that by having a unifying substrate and explicit translation mechanisms. It echoes Marvin Minsky's idea that a mind emerges from a society of agents (which in our view are like processes following different DSL rules), and these agents need a way to cooperate.

It's worth noting that building such an integrated system is complex, but it's arguably necessary for *general* intelligence. Narrow AI systems may excel at one DSL (say, CNNs for vision), but an intelligent agent in an open world must juggle many. By formally structuring the juggling, we get a handle on verifying and understanding it.

## 3.7 Intelligence as a Dynamic Computational Topology

We have now delineated the components of the system: structured memory, recursive feedback loops, learning mechanisms, trace logging, and multi-DSL modules. To truly understand how *cognition emerges*, it is helpful to view the entire system as a **dynamic computational topology** – essentially a graph of interacting processes and memory elements that reconfigures over time. This section paints a holistic picture of the intelligent system as an evolving network, and explains how *emergent cognitive properties* arise from it, fully consistent with our low-level formalism.

**Topology Elements:** We can abstract the running system at time $t$ as a directed graph $G_t = (V_t, E_t)$:

- The **nodes $V_t$** represent the active computational units or stores at time $t$. This includes function modules (instances of $f$ or its subroutines), memory regions (each BMR can be a node when we consider data flow), and even the feedback generator $\Phi$ as a node. For example, one node might be the "Vision DSL Processor (f_vis)" currently analyzing an image, another node might be "Working Memory Region for Visual Features ($R_{\text{vision}}$)", another might be "Logic Inference Engine (f_logic)", etc. Some nodes are more static (memory containers), others are ephemeral (a computation that produces a result then ends).

- The **edges $E_t$** represent the flow of information and control between nodes. An edge from memory node $R_i$ to function node $f_j$ indicates that $f_j$ reads data from $R_i$. An edge from function node $f_j$ to memory node $R_k$ means $f_j$ writes/updates that memory. Feedback loops are cycles in this graph: e.g., an edge from a node representing output (or world state) back into a node representing the feedback function $\Phi$ input. Edges can also carry control signals – e.g., one function triggering or calling another (or sending a "done" signal).

This graph is essentially a dataflow + control flow graph of the system's computation. What makes it **dynamic** is that $G_{t+1}$ can differ from $G_t$: learning might add a new node (a new rule acquired, effectively enlarging $f$'s repertoire), or prune nodes (drop an intention, freeing up that part of memory), or adjust edges (rewire which module gets input from where). There is even the possibility of structural self-modification: adding a new memory region (if the system decides to allocate space for new knowledge) or spawning a new DSL interpreter instance for a subproblem (like spawning a sub-agent). This is analogous to how our brains form new neural connections or recruit different circuits when learning new tasks. In artificial systems, we see analogies in evolving neural network topologies or self-modifying code.

We can illustrate a simplified topology: imagine at some moment, nodes include [Sensor Input], [Memory: image buffer], [Vision Module], [Memory: detected objects], [Logic Module], [Memory: inferred fact], [Planning Module], [Memory: plan], [Action Module]. Edges connect from sensor to image buffer, buffer to vision, vision to detected objects memory, that memory to logic, etc., forming a pipeline with feedback edges perhaps from action outcomes back to logic (did it work?) and so on. If the scenario changes (say a new sensor becomes available or a subtask begins), the topology might rearrange, adding perhaps a [Speech Module] node and edges for auditory input.

**Interplay and Emergence:** The *behavior* of the entire system is the result of potentially many nodes firing and passing information along edges simultaneously (or in rapid sequence). This massively parallel (conceptually) process yields outcomes that no single module could produce in isolation. *Emergence* in this context means the whole is more complex than the sum of parts: intelligence is not in one node, but in the **pattern of activation and adaptation across the network**. However, unlike mystical notions of emergence, here we can still trace it: it's emergent *yet traceable*. We might observe, for instance, a spontaneous strategy arising – say the system figures out a heuristic to solve problems faster. In the topology view, this could correspond to a subgraph that gets frequently activated together and ends up effectively implementing a new procedure. Because we log everything, we could identify that subgraph and recognize the new emergent procedure (maybe even give it a name and incorporate it as a formal routine later).

The **driving force** of the dynamics is the continual attempt to satisfy goals, reduce errors, and integrate new information. In graph terms, this can be seen as *activation spreading and settling*. External input $E_t$ injects activity into some nodes (e.g., sensor feeds memory which activates a processing node). Internal goals (which might be encoded in a persistent memory region as some desired state or value to maximize) also inject activity – they might bias certain edges (like feedback edges carrying a "desire signal"). Through feedback edges ($\Phi$ loops), if a certain pathway leads to high reward or consistency, that pathway's edges might be strengthened (conceptually, more information flows there or it becomes the default route). Unproductive pathways might get inhibited or eventually pruned. Over time, the topology *self-organizes* such that it funnels information along routes that achieve goal satisfaction and stable predictions – essentially minimizing surprise/entropy globally [32] .

This description aligns well with **neural network interpretations** and **symbolic AI interpretations** alike. In a neural view, $G_t$ could be seen as the neural network graph (which can change with neuroplasticity). In a symbolic view, $G_t$ is more like a flowchart of modules (which can change if the system forms new plans). Our framework encompasses both: symbolic DSL modules and sub-symbolic adaptation co-exist. Modern cognitive science also leans towards hybrid models where continuous dynamics and discrete symbols both play a role [33] .

Let's reflect on **emergent cognition**. We claim that cognitive phenomena – like problem-solving, reasoning, learning, creativity – emerge from this kind of substrate without needing any extra mysterious ingredient. For example:

- *Problem-solving* emerges when the system encounters a novel situation (pattern of $E_t$) that doesn't immediately map to a known response. The topology then might reconfigure as various modules explore combinations (a search in state space). When a solution is found (feedback indicates success), that pathway is reinforced. Later, facing a similar problem, the system quickly reactivates the successful subgraph – effectively it has learned a method. This is precisely how one might interpret human problem solving with insight: many mental resources unconsciously explore, and then suddenly a coherent solution path "clicks" (we can see that as the subgraph stabilizing and producing the right output).

- *Memory recall* in our system is just the activation of certain memory nodes due to associative links (edges) from current context nodes. For instance, if a current input triggers a pattern that matches something in long-term memory (persistent region), an edge will carry that activation to fetch that memory content. We could say the system "remembers" that information relevant to the context.

This is not implemented as a separate magic, it's literally data flow in the graph guided by keys or patterns. Cognitive architectures like ACT-R model this with an activation spreading mechanism where chunks of memory are retrieved based on similarity to current goals and context – a very similar idea [5] .

- *Meta-cognition* (thinking about thinking) happens when parts of the topology start monitoring or modulating others. For example, a node representing the learning rate or strategy might adjust edges in other parts (like deciding to switch from exploration to exploitation – effectively tuning how $\Phi$ behaves). In our logs, this would be visible as certain feedback signals being generated internally not directly from external input but from introspective evaluation of performance. Because the system can represent its own process in the ledger, it can in principle reflect on that (if we design a DSL for introspection that reads the ledger or monitors key variables). This yields a self-aware adjustment capability (though not necessarily self-awareness in the philosophical sense).

The **dynamic reconfigurability** through learning and plasticity means the system is not static. For example, suppose the system learns a new skill (like a new DSL or a new concept). This could manifest as literally adding new node(s) to $V_t$ (a new rule memory, a new routine in $f$) and new edges connecting it (linking it into the existing knowledge network). The trace ledger will show when this happened and how (maybe "at time t1000, learned rule X, added to memory region for rules"). If that skill later is unused and maybe pruned, the ledger shows its removal. This is analogous to how brains form new synapses and later prune unused ones in development, which is believed to be crucial for efficient cognitive function [34] [35] .

One can also view the **entire lifetime** of the system as generating an increasingly rich topology. Early on, $G$ might be sparse and mostly pre-built (from initial programming). As the system experiences more, $G$ becomes more interconnected, encoding the system's accumulated knowledge. This graph, together with the ledger, forms a complete picture of the system's state and history – a kind of *computational knowledge graph* that is both the *product* of intelligence and the *source* of it.

To wrap up this perspective: our intelligent system is *not* a fixed program but a self-organizing computational network regulated by feedback. It operates always within the bounds of binary logic and formal rules, yet through its complexity it exhibits the flexible, adaptive behavior we call cognition. Importantly, none of this escapes scrutiny: because it's all within a formal, logged system, we can analyze this topology at any time. We can identify emergent *properties* (like "it has developed a concept of X because we see a cluster of nodes/edges representing X-related processing") and even prove certain *guarantees* (for instance, a safety constraint can be mapped onto a graph property and we can verify the graph never violates it, since any violation would appear in a trace).

This view resonates with theoretical discussions in cognitive science and AI that hybridize symbolic and connectionist paradigms, seeing cognition as neither purely local rules nor pure holistic dynamics, but an interplay [36] . By providing a formal scaffold (the binary substrate and BDI) for such an interplay, we ensure that **emergence does not undermine rigor**: we get the best of both worlds – a system that can evolve unpredictably novel solutions, yet remains within a provable, understandable framework.

## 3.8 Towards Verifiable Cognitive Architectures

Throughout this chapter, we have constructed a view of intelligence as *composable, explainable computation* rather than an ineffable quality. We combined principles of binary mathematics with an architecture for

cognition, yielding a framework where adaptive behavior arises naturally from structured processes. Let us summarize the key points and emphasize the significance of this approach:

- We provided an **operational definition of intelligence** centered on five capabilities: structured memory, recursive feedback loops, entropy-guided learning, hierarchical abstractions (multi-DSL), and verifiable reasoning. These were not just listed, but each anchored in formal mechanisms (for example, $M_t$ for memory, $\Phi$ for feedback, $J$ and $\nabla J$ for learning, etc.).

- Each capability was expanded with **technical precision**: we introduced formal notation for memory and functions, outlined update rules, and even sketched how logging can be done with hashes and signatures. This rigor means the framework is not merely philosophical; it's ready to be implemented or mathematically analyzed. Indeed, pieces of this puzzle exist in research today (e.g., Kotseruba *et al.*'s survey of cognitive architectures shows various systems covering parts of this [5], and cryptographic verification of ML pipelines is an emerging field). Our contribution is weaving them into a unified picture.

- A recurring theme is **verifiability and traceability**. By design, every operation is traceable to the binary substrate and logged. This makes the system's knowledge *scientific*: it carries its own evidence. An auditor (human or another program) can verify any claim the system makes by inspecting the proof trace. This addresses a critical trust issue in AI – as AI systems become more pervasive, trust will depend on *evidence of correctness* [17]. Our framework not only recognizes that but builds it in. The system is effectively its own auditor at runtime, which can be supplemented by external audits of the logs if needed.

- We showed how **intelligent behavior is grounded in binary processes**. Nothing was assumed to be outside the binary logic: no magic symbols, no non-algorithmic leaps. Yet, through recursion and self-reference (via feedback loops), the system can exhibit open-ended, seemingly creative behavior. This demystifies cognition. It suggests that what we call "understanding" or "thinking" is the emergent consistency achieved by the system updating itself to reduce errors and increase coherence in its world model. By avoiding any appeal to an undefinable spark (we didn't, for example, require consciousness or quantum processes or the like), we keep the theory within the realm of computability and mathematics. This means we can build it (construct AI that follows these principles) and we can analyze it with the tools of theoretical computer science, information theory, and logic.

- The emphasis on **composability** underlines that complex intelligence can be built up from simpler parts. Each part by itself might not be "intelligent" (just as a neuron by itself isn't), but when composed under the right architecture, they achieve greater functions. We specifically advocated multiple DSLs and integrated them, echoing the way human cognition likely integrates multiple specialized areas of the brain. The framework doesn't limit how many or which DSLs – it gives a template for incorporating new ones. This can serve as a *formal blueprint for Artificial General Intelligence (AGI)* research: instead of searching for one master algorithm, ensure you have the pieces (memory, learning, perception, reasoning, etc.) and then focus on the interfaces and orchestration. As long as the integration respects the verifiability and substrate rules, the whole remains testable.

- Finally, the chapter positions this approach as moving beyond philosophical debates to **constructible systems**. We can imagine implementing a simplified version of this architecture and

incrementally enhancing it. Each piece is independently grounded: e.g., one can implement the ledger with existing blockchain or database tech, implement learning with established ML techniques, implement multi-DSL with a plugin-based software architecture, etc., and the BDI as an operating system or middleware that connects them. What the theory provides is a *guarantee of coherence*: if you implement all these pieces and constraints, you will have something that can rightfully be called an intelligent system by our operational definition. And unlike many AI systems today, you would be able to **explain and verify** its every step – a trait highly desirable for deploying AI in society [17].

Looking forward, this compositional framework hints at a new generation of **cognitive operating systems** or **verifiable cognitive architectures**. Just as early computers had to evolve from ad-hoc designs to formal architectures (like von Neumann architecture, which provided a clear blueprint for building any general-purpose computer), AI systems might evolve from narrow models to architecture-driven designs that ensure generality, transparency, and safety. The Binary Decomposition Interface (BDI), which is the subject of the next chapter, is one key piece of that evolution: it provides the formal substrate on which all higher processes run. By defining intelligence in terms of what BDI and structured processes can achieve, we align AI with the rigor of computer science and the accountability expectations of engineering disciplines.

In conclusion, **Composable Intelligence Systems** as outlined here offer a path to machines that *earn* the label of intelligence through demonstrable, verifiable competence rather than anthropomorphic comparison. They reinforce the notion that cognition is *emergent but explainable*: an outcome of many simple binary steps rather than a single leap of mind. This not only advances our ability to design intelligent technology, but also enriches our understanding of natural intelligence by providing a concrete computational mirror. Each concept – memory, learning, reasoning, abstraction, self-reflection – is given a place in the architecture, inviting further refinement and connection to psychological and neurological findings. As we proceed to delve into the BDI in Chapter 4, we keep in mind that it is the bedrock enabling everything described here to be realized in practice. The promise of this formal, modular approach is an AI that we can build piece by piece, **verify piece by piece**, and ultimately trust as a reasoning partner, precisely because we know *how* it reasons at every level of detail.

**Bibliography**

1. K. Friston, "The free-energy principle: a unified brain theory?," *Nature Reviews Neuroscience*, vol. 11, no. 2, pp. 127–138, 2010.

2. A. Rao and M. Georgeff, "BDI agents: from theory to practice," in *Proceedings of the 1st International Conference on Multi-Agent Systems (ICMAS)*, 1995, pp. 312–319.

3. M. Minsky, *The Society of Mind*. New York: Simon & Schuster, 1986.

4. I. Kotseruba and J. K. Tsotsos, "40 years of cognitive architectures: core cognitive abilities and practical applications," *Artificial Intelligence Review*, vol. 53, pp. 17–94, 2020. [5]

5. X. Jiang *et al.*, "Long Term Memory: The Foundation of AI Self-Evolution," arXiv:2410.15665 [cs.AI], 2023.

6. A. Dedhe *et al.*, "Cognitive mechanisms underlying recursive pattern processing in human adults," *Cognitive Science*, vol. 47, no. 1, e13166, 2023.

7. P. Langley and D. Choi, "A unified cognitive architecture for physical agents," in *Proceedings of the National Conference on Artificial Intelligence (AAAI)*, 2006.

8. V. Kulothungan, "Using Blockchain Ledgers to Record the AI Decisions in IoT," *Preprints*, 2025.

9. P. Isaev and P. Hammer, "Memory system and memory types for real-time reasoning systems," in *Proc. of Artificial General Intelligence (AGI)*, Springer, 2023, pp. 147–157. [37]

10. E. Salazar, "COGENT: An AI architecture for emergent cognition," arXiv:2504.04139 [cs.AI], 2025.

11. J. Laird, A. Newell, and P. Rosenbloom, "SOAR: an architecture for general intelligence," *Artificial Intelligence*, vol. 33, no. 1, pp. 1–64, 1987. [38]

12. European Parliament and Council, "EU Artificial Intelligence Act," 2024 (Proposed).

---

[1] [5] Cognitive Memory in Large Language Models
https://arxiv.org/html/2504.02441v1

[2] [3] [4] [6] [7] [8] [9] [10] [11] [12] [13] [14] [15] [16] [17] [18] [19] [20] [21] [22] [23] [24] [25] [26] [27] [28] [32] [34] [35] [37]
[38] 2025-04-27-Binary Mathematics and Intelligent Systems by Tariq Mohammed Chapters 1-10.pdf
file://file-UBwwFk7FpYVMod3u9msGdJ

[29] [PDF] A review of cognitive architectures
https://idealectic.com/iso2_report.pdf

[30] [PDF] Multiple Representations in Cognitive Architectures - AAAI
https://cdn.aaai.org/ocs/15982/15982-69922-1-PB.pdf

[31] Carruthers's massively modular architecture of the mind. Each box... | Download Scientific Diagram
https://www.researchgate.net/figure/Carrutherss-massively-modular-architecture-of-the-mind-Each-box-represents-multiple_fig3_334709058

[33] [36] Introducing \cogentCOGENT^"3" An AI Architecture for Emergent Cognition
https://arxiv.org/html/2504.04139v1