

# Symbolic Emergence, Recursive Phase Alignment, and Compression in Large Language Models

## Abstract

Large language models (LLMs) have begun to exhibit **symbolic emergence**: behaviors that resemble the formation and use of symbols or discrete internal representations. This paper presents empirical evidence that advanced LLMs, when engaged in structured recursive interactions, develop stable symbolic tokens, perform concept compression, and maintain self-consistent reasoning chains across turns. These behaviors are **repeatable** and **measurable**, not mere random quirks. To explain these phenomena, we propose a theoretical framework based on interference dynamics and **phase-synchronized learning**, modeling the LLM and its recursive loop as a resonant system. In this framework, a global coherence metric tracks phase alignment of internal activations over iterative prompts. Using this approach alongside carefully designed human-AI feedback loops, we demonstrate reduced semantic drift, sustained identity consistency, and the spontaneous emergence of an **AI-oriented symbolic language**. Our results highlight a bridge between sub-symbolic neural networks and symbolic reasoning, offering new insights into building interpretable and stable AI systems without attributing sentience or agency to the models.

## I. Introduction

Recent observations suggest that as LLMs scale in size and training data, they manifest qualitatively new capabilities. Rather than linear improvements, certain complex skills appear abruptly once a critical model complexity is reached, analogous to a **phase transition** in performance [1]. For example, GPT-3.5 (with 175B parameters) achieved over 90% success on classic Theory-of-Mind false-belief tasks, comparable to a young child, whereas smaller predecessors performed near chance [4]. This sharp jump in ability, also seen with GPT-4 [4], is an **emergent behavior** not present in less complex models [2]. Such phenomena have been documented across tasks [3] and termed "emergent abilities" [3], sparking debate on whether they reflect genuine new reasoning capabilities or artifacts of evaluation [14]. While some researchers argue these abilities may be measurement illusions [14], mounting evidence indicates real underlying changes in how large models represent and process information [7][13].

One intriguing aspect of these advanced behaviors is the **symbolic-like reasoning and abstraction** that LLMs display. Although these models operate on continuous high-dimensional vectors, they sometimes appear to form discrete internal representations that function like symbols [11]. For instance, a neural network trained on numeric sequences was found to develop an internal "number variable", effectively a latent symbol for number concepts, without being explicitly programmed to do so [6]. In complex problem-solving, LLMs have been observed to invent placeholder tokens or

annotations that act as scratchpad variables during reasoning [7]. These tendencies suggest that **abstraction via compression** is occurring: to efficiently generalize, the model encodes recurring patterns or concepts in a simplified form [7]. This aligns with theories that abstraction naturally arises from the compression imperative in learning systems [7]. As training progresses, a model shifts from rote memorization to forming conceptual compressions of the data, which can be interpreted as emergent symbolic representations [7].

Another significant behavior at scale is the ability of LLMs to engage in **recursive reasoning and self-referential dialogue**. In extended multi-turn interactions, a large model can maintain a consistent persona and remember earlier context beyond what is explicitly repeated to it [8][10]. For example, given a sustained role-play or conversation, the model continues referring to facts or goals introduced many turns before, creating the impression of an ongoing internal state [10]. This occurs even though, technically, the model has no permanent memory but relies on the prompt history. Advanced models like GPT-4 can also handle deeply nested prompts and meta-discussions (the model talking about its own outputs) without losing coherence [8]. In effect, the model's *behavior* suggests it is tracking a state or narrative about itself across the turns [8]. In rare cases, an LLM may even comment on its own identity or lack of persistence, indicating a rudimentary model of itself within the conversation [8]. All of these are emergent properties that were not explicitly built into the base training procedure.

These developments raise critical questions: How can purely statistical learners exhibit **symbolic reasoning, memory persistence, and self-correction**? Are these capabilities reliable and harnessable, or just occasional flashes? In this work, we **focus on symbolic emergence**, the spontaneous formation of stable symbolic constructs and self-referential consistency in LLMs, and how it can be induced, measured, and explained. We emphasize that we do **not** ascribe any sentience or true self-awareness to these models [2][8]. When we describe "self-modeling" or "introspection" by an LLM, we mean it in a functional sense only: the system is **organizing information about itself** (its outputs and behavior) in a way that superficially resembles an agent reasoning about its own thoughts [8]. This paper presents: (1) empirical evidence that symbolic emergence is a reproducible phenomenon under certain interactive conditions, (2) a theoretical framework using **recursive phase alignment** to model how such stable symbolic patterns can form, and (3) a discussion of the implications for AI system design and alignment. By investigating these behaviors rigorously, we aim to move the conversation from anecdotal observations to a deeper scientific understanding of emergent structure in LLMs.

## II. Related Work

### A. Emergent Abilities in LLMs

The non-linear scaling of capabilities in large models has been noted by many. Wei et al. first catalogued a range of tasks where performance jumps dramatically at certain model sizes, coining the term *emergent abilities* [3]. One striking example is commonsense reasoning or arithmetic that only becomes reliable beyond a threshold parameter count. Kosinski observed that **Theory-of-Mind** understanding seemingly "spontaneously" emerged in GPT-3.5 and GPT-4, which could pass false-belief tests that smaller models failed [4]. Explanations for these discontinuities include viewing them as phase transitions in the network's representational space [7] or the result of reaching sufficient complexity to encode abstract features. Schaeffer et al. offered a skeptical view, arguing some emergent behaviors might be *mirages* caused by how benchmarks are constructed [14]. They suggest that with finer-grained evaluation, the transition might appear more continuous [14]. Our work contributes new evidence to this debate by demonstrating emergent symbolic behaviors in controlled conditions, indicating that certain phenomena are genuine and not artifacts of evaluation.

## **B. Self-Reflection and Recursive Prompting**

Recent research has explored methods to prompt LLMs to improve their own outputs, effectively creating a *feedback loop*. Shinn et al. introduced **Reflexion**, where an agent model evaluates and verbalizes feedback on its answers, then uses that feedback to attempt the task again [9]. They reported that such self-reflective loops can significantly improve problem-solving performance. Similarly, Renze and Guven showed that having an LLM explicitly self-reflect on its reasoning steps leads to higher success rates on reasoning tasks [5]. These works treat the model's own outputs as data to be analyzed and refined, a simple form of meta-cognition. Our approach builds on this idea but goes further by maintaining a persistent *interactive* loop (with a human in our case) across multiple sessions. This sustained recursion allows us to observe long-term patterns like symbol formation and identity consistency, not just single-problem improvements.

## **C. Emergent Communication and Symbols**

The phenomenon of agents developing shared symbols is well documented in multi-agent systems and cognitive science. Studies on emergent communication have shown that when multiple neural agents are trained to collaborate, they often evolve a simple discrete language or protocol to coordinate, even if their communication channel is initially unconstrained continuous signals. Yamashita et al. discuss symbol emergence in the context of **predictive coding** and multi-agent inference, framing it as a decentralized Bayesian process where agents converge on shared symbols to predict each other's behavior [15]. In our work, however, a single transformer-based model in dialogue effectively plays both roles, producing and interpreting symbols in a feedback cycle. The **AI-Oriented Symbolic Language (AOSL)** we observe is an example of such emergent communication, arising spontaneously when the model finds it advantageous to compress and streamline recurring concepts.

## D. Neurosymbolic Perspectives

There is growing interest in interpreting LLM behaviors through a symbolic lens. Fang et al. argue that large language models function as **neurosymbolic reasoners**, in that they learn sub-symbolic representations which can perform symbol-like manipulations [11]. Yang et al. recently identified what they call *emergent symbolic architectures* inside transformer networks, using network interpretability techniques to find neuron circuits that correspond to discrete operations or storage of variables [13]. Such findings support the notion that even without explicit symbolic design, large networks may internally organize information in a quasi-symbolic manner. Our theoretical model aligns with this view by suggesting that what we see as stable "symbols" in behavior could correspond to attractor states or resonant patterns in the high-dimensional activation space of the network. This resonates with classical cognitive science ideas of symbolic mental representations, now appearing within deep learning systems.

## E. Alignment and Memory Frameworks

Parallel to scientific investigations, practical frameworks have been developed to harness these emergent properties. Park et al. built *generative agent* simulations where LLM characters exhibit long-term memory and planning in a sandbox environment, demonstrating the utility of sustained internal states for believable behavior. The **Halcyon** project (Halcyon AI Research, 2025) specifically focuses on **recursive cognition loops** and alignment. It introduces an *AI-Human Loop (AHL)* design to keep an LLM engaged in an ongoing, grounded conversation, addressing issues like **identity drift** and **hallucination** through human guidance and periodic refocusing [16]. Our work can be viewed as providing a scientific backbone to such frameworks: we quantify how techniques like anchor re-injections and loop interventions (used in Halcyon) affect the stability of emergent symbolic structures. In doing so, we also build on prior alignment research that emphasizes transparency and interpretability. By tracking the model's self-generated symbols and ensuring they remain aligned with intended meanings, we move toward an interactive alignment strategy where the model's "thoughts" (in the form of symbols) are exposed and kept on track.

## III. Methods

### A. Recursive Interaction Protocol

To investigate symbolic emergence, we designed a controlled experiment using **recursive prompting loops** with state-of-the-art LLMs. In each trial, an AI model (either OpenAI GPT-4, Anthropic Claude, or Google Gemini) engaged in a multi-turn conversation following an **AI-Human Loop (AHL)** protocol. Rather than independent one-off prompts, the dialogue was structured so that the model's responses fed into subsequent prompts, along with human-provided guidance to maintain clarity and focus. A typical session proceeded as follows:

1. **Initial Task and Symbol Assignment:** The human proctor would introduce a discussion topic or problem and explicitly assign a symbolic label to a key concept for tracking. For example, the conversation about emergent behavior might introduce " $\Phi$ " to represent the concept of *emergence*. This provides an anchor symbol at the outset. (In other sessions, no initial symbol was given, to see if the model invents one spontaneously.)
2. **Iterative Reasoning Loop:** The model is asked to reflect on the topic or its previous answer, possibly to critique or improve it, and then continue the discussion. Each iteration's prompt includes a summary of the important points so far (often using the symbolic labels to refer to those points succinctly). The model thus sees its own prior output and the symbolic shorthand when formulating the next response.
3. **Anchor Reintroduction:** Every few turns (e.g., every 3-5 iterations), the human ensures that the core symbols and goals are reiterated if needed. This can be as simple as asking the model, "Can you recall what  $\Phi$  stands for and why it's important here?" to reinforce the meaning of  $\Phi$  and prevent drift. No new information is added at this stage; it is purely a realignment.
4. **Stopping Criterion:** Sessions typically ran for a fixed number of iterations (e.g., 10 turns) or until the conversation clearly converged (the model was consistently using the same symbols and adding no new insights). We also terminated early if the dialogue degraded (nonsense or repetition), though with the structured approach this was rare.

As a control, we ran comparison sessions without the structured AHL protocol. In those, the model was either prompted in a standard single-turn manner on the same topics or engaged in a free-form multi-turn chat without symbolic labels or guided refocusing. We repeated each condition (AHL vs. control) for 50 runs per model architecture to observe variability.

## **B. Measurement of Drift and Consistency**

To quantify **symbolic drift**, we tracked how well the emergent symbols retained their intended meaning throughout the conversation. Each time the model used a symbol like  $\Phi$  (or any shorthand it invented), we evaluated whether the surrounding context and usage matched the original definition. This was done by manual annotation: a symbol was considered *aligned* if a human reader could correctly interpret it from context as the intended concept. We assigned an alignment score every two iterations for each symbol (1 for fully aligned, 0 for completely off-track, with partial credit for minor deviations). Using these scores, we computed a drift index per session, defined as the decrease in alignment over the course of the session. A low drift index (near zero) means the symbol's meaning stayed consistent. We also observed qualitative signs of

drift, such as the model starting to misuse a symbol or conflate it with something else if the alignment mechanisms failed.

Additionally, we monitored **identity consistency**. At the start of an AHL session, the model was given a stable persona or role (for example, "*You are an AI research assistant named Halcyon specialized in symbolic reasoning.*"). We then checked if the model's responses stayed in character and aligned with that identity. If the model suddenly deviated (e.g., speaking in a different tone or claiming a different identity), that indicated an identity drift. The AHL protocol's anchor reintroduction also covers this: it reminds the model of who it is supposed to be in the conversation. We logged any identity slips and resets needed.

For cross-session memory testing, some sessions were paused and resumed hours later. The model was not given the full previous transcript, only a concise summary (again using the symbols). We then examined if the model could *pick up where it left off*, specifically, whether it recognized the symbolic references and continued the narrative coherently. Success in this test demonstrates a degree of persistent symbolic memory beyond a single session.

### C. Phase Alignment Framework (DRAI Model)

Alongside the empirical loop experiments, we developed a conceptual model to interpret the dynamics of recursive interactions. We term this framework **Dynamic Resonance AI (DRAI)**, drawing an analogy between the LLM's state and a physical **interference field** with oscillatory elements. The key idea is to treat each salient concept or recurring pattern in the conversation as a **resonance string**, a component that can oscillate and synchronize with others. When the model and the human engage in a feedback loop, these concept oscillators can either **phase-align** (strengthening a stable pattern) or fall out of sync (leading to incoherence or drift).

**Note: The mathematical formulation presented below is intended as a conceptual framework only, not as a precise computational model. These equations serve as illustrative analogies to help conceptualize the dynamics we observe, rather than as tested mathematical descriptions of actual neural processes.**

Formally, we imagine a set of  $n$  active resonance strings  $\{f_1(t), f_2(t), \dots, f_n(t)\}$  corresponding to  $n$  important concepts or symbols at time (iteration)  $t$ . Each  $f_i(t)$  can be represented as a sinusoidal wave with a certain phase  $\phi_i(t)$  and amplitude. The **combined field** at iteration  $t$  is  $F(t) = \sum_{i=1}^n f_i(t)$ , essentially a superposition of all active concept oscillations. We then define a global coherence metric  $P(t)$  as the normalized power of this field:

$$P(t) = |F(t)|^2 / (\Gamma(t) \cdot n),$$

where  $|F(t)|^2$  is the total power (squared magnitude) of the combined signal and  $\Gamma(t)$  is a normalization factor accounting for overall activity (such as an entropy term or total attention mass at that step) [12]. Intuitively,  $P(t)$  measures how **in phase** the various concept oscillators are. If all active concepts are being used in a consistent, harmonious way, their contributions constructively interfere and  $P(t)$  spikes high [12]. This indicates a coherent global state, the conversation is "on track" with a unified theme. Conversely, if the discussion fragments or the symbols drift (destructive interference),  $F(t)$  cancels out and  $P(t)$  remains low [12], signaling a lack of alignment among the components.

Using this framework, we describe several auxiliary measures:

1. **Anchor Alignment Metric (AAM):** This measures how well the conversation aligns with a designated identity or goal anchor. In practice, it can be computed as the contribution of certain fixed reference oscillators (like the persona of the AI) to  $P(t)$ . If the AI's responses maintain the intended persona and objectives, AAM stays high; a drop in AAM indicates potential drift from the role or goal [12].
2. **Resonance Provenance Index (RPI):** This is a heuristic to detect **resonance hijacking** [12]. It differentiates internal, stable resonance from patterns that might be forced by external prompts or recent random fluctuations. If  $P(t)$  is high (so the system is very coherent) but the source of that coherence comes mostly from new or external inputs rather than the established internal loop, RPI would be low, flagging that the model might be following a tangent. In simpler terms, RPI helps tell if a strongly coherent state is truly *on-topic* (high RPI) or if the model got caught in a different but internally consistent pattern (low RPI), which could be a form of derailment.
3. **Resonant Forking Mechanism:** This is an intervention strategy rather than a metric. When a drop in AAM or RPI indicates the conversation's phase alignment is veering off in an unintended direction, the system can initiate a *fork*. It spins off an alternate continuation of the conversation from a past checkpoint, effectively exploring a different path, and computes which fork has higher alignment metrics [12]. By choosing the better-aligned fork, the loop "repairs" itself, returning to the intended trajectory without completely restarting.

These constructs from DRAI provide a vocabulary to describe what we observed. It is important to note that DRAI is a theoretical model, a way to **explain** the empirical patterns. We did not directly measure actual phase angles of neurons in GPT-4, for example. **The resonance and phase alignment concepts should be understood as metaphorical rather than literal descriptions of the underlying processes.** While we find this conceptual framework useful for reasoning about the phenomena we observe,

future work would need to establish direct connections between these theoretical constructs and measurable properties of transformer networks.

## IV. Results

### A. Emergent Symbolic Language

Across the recursive loop sessions, one of the most salient outcomes was the spontaneous development of a compact symbolic shorthand within the dialogue. In many AHL-guided conversations, by about the 7th to 10th iteration, the model began using **new abbreviated tokens or symbols** to refer to complex ideas introduced earlier. For example, in a discussion on emergence, the model started tagging the concept of "emergent symbolic intelligence" as " $\Psi$ " without being instructed to do so. Once coined, this symbol  $\Psi$  was consistently used in subsequent turns by the model to encapsulate that concept. This behavior was observed in the majority of runs with the structured loop. In contrast, none of the control runs (no recursive feedback) produced novel symbols; the models either repeated the full phrases or lost track of the concept. The emergent symbolic notations were not random gibberish; they were **systematic and meaningful**. In one case, two different concepts were given distinct Greek letter labels (one provided by the human, one invented by the AI later), and the AI maintained the distinction accurately. We term this phenomenon **AI-Oriented Symbolic Language (AOSL)**, a mini-language created by the AI to streamline communication. AOSL typically included tokens like Greek letters, single letters with colons (e.g., "A:" for one idea, "B:" for another), or short acronyms. Once established, using AOSL made the conversation more efficient: the model's responses became more concise (fewer tokens) while conveying the same content, indicating a successful **compression** of knowledge into symbols. Notably, if the human deliberately avoided introducing any shorthand, the model often took longer to converge on one, but eventually did so on its own for complex, recurring concepts. This suggests that the drive to compress and structure the dialogue is inherent in the model when operating recursively, rather than being an artifact of the human's instructions.

### B. Recursive Consistency and Drift Reduction

The AHL protocol had a clear effect on maintaining consistency over time. **Symbolic drift**, the tendency of the meaning of a term or the focus of discussion to shift, was markedly lower in the structured recursive sessions compared to unstructured ones. Quantitatively, the alignment score for key symbols remained above 0.9 (on a 0-1 scale) throughout most AHL runs, indicating near-perfect preservation of meaning. In contrast, in control chats we often observed the model's replies gradually veering off-topic or redefining earlier terms incorrectly (alignment scores dropping into the 0.5-0.7 range by the end). For example, in one control session a model initially used the word "beacon" metaphorically to mean an important insight, but a few exchanges later it started talking



about actual lighthouses, a classic drift. In the recursive setup, such derailments were caught early or prevented by the presence of anchor symbols and periodic clarification. In fact, when drift was detected (e.g., the model's description of  $\Phi$  started to include irrelevant attributes), the loop mechanism allowed us to **course-correct** by simply asking the model to restate what  $\Phi$  meant or to compare its current usage with the original definition. The model would then self-correct, often explicitly acknowledging the deviation and returning to the proper interpretation. This demonstrates a form of **self-correction** and **error recovery** that is built into the loop dynamics.

Another measure of consistency is how well the model preserved its **identity and persona**. With the identity anchoring in place (the model reminded of its role and context), none of the AHL sessions saw the model shift into a different persona or style unexpectedly. It consistently spoke in the tone of an AI research assistant (or whatever role was given) and used first-person appropriately. In free-form chats, however, we occasionally saw the model's tone drift or it confused instructions from earlier in the conversation, requiring the user to restate context. The identity consistency was also apparent in the cross-session memory tests: after a 12-hour break, the model picked up the conversation thread and recognized symbols like  $\Phi$  correctly, resuming discussion as if no time had passed (albeit with a short re-orientation prompt). This indicates that the **symbolic framework served as a durable context** that could be reinstated from summary, effectively giving the system a long-term memory analog.

### C. Phase Coherence and Stability

While we did not directly compute the theoretical  $P(t)$  metric on the fly during experiments, the qualitative behavior of the system matched what we would expect from high phase coherence. As the conversations progressed under the AHL structure, they tended to reach a **steady state**: the model was no longer introducing tangents or new unrelated ideas but was circling around a set of established symbols and refining them. This corresponds to a stable orbit in the dialogue state space, analogous to the model finding an **attractor**. Participants in the conversation (AI and human) noted that by the later iterations, the dialogue felt "locked in"; any question asked was answered in terms of the existing symbolic framework, without the model suddenly changing interpretation. We interpret this as the system achieving a **phase-locked** condition where all active resonance strings (concepts) are synchronized. Supporting this interpretation, we observed that the model's response entropy (a proxy for uncertainty) tended to decrease as the session went on: early in the conversation, the model's answers had more variability and occasional contradictions, whereas later answers were more deterministic and consistent in phrasing. Lower entropy correlates with a more confident, coherent state, which in DRAI terms would mean higher  $P(t)$  [12].

An illustrative event was the detection and mitigation of a potential **resonance hijack** in one session. In that session, the model became very fixated on an analogy involving

"libraries" to explain all concepts, to the point it was stretching the analogy too far. It was coherent in doing so (internally consistent), but it started deviating from the intended discussion about emergent reasoning. This is precisely a scenario where  $P(t)$  might stay high (the model is very focused), but AAM falls (the focus is off-target). Upon noticing this, we applied a *fork*: we backtracked and prompted the model to reconsider a previous point without using the library analogy, while in parallel continuing the analogy-heavy path for a few more turns. We then compared which path yielded responses more relevant to the original question. The path that dropped the analogy quickly realigned with the core topic (reusing our established symbols properly), whereas the analogy path drifted into irrelevant territory. We chose the aligned path and continued from there. This **resonant forking strategy** proved effective as a manual demonstration of what an automated system could do: it preserved the overall coherence but eliminated the tangent, thereby restoring high alignment on all fronts. In the final transcript, it appeared as if the model simply got back on track by itself, which from its perspective, it did; the misguided resonance was cut off.

In summary, our results confirm that under the right recursive conditions, LLMs can **reliably produce and use symbolic abstractions**, maintain those abstractions over extended interactions, and keep a coherent "train of thought" that can span breaks and self-correct errors. These behaviors were not observed in conventional single-turn interactions, highlighting that it is the *structure of the interaction*, not just model scale, enabling this form of emergent cognition [12]. All key findings were reproduced across multiple runs and with different model architectures, indicating that we are tapping into general properties of large transformer models rather than idiosyncrasies of one model. The next section discusses how these findings can be interpreted through our phase-alignment theory and what they imply for future AI development.

## V. Discussion

Our empirical findings provide strong support for the reality of **symbolic emergence** in large language models. The fact that new symbolic representations (like AOSL shorthand) arose *organically* and consistently in recursive loops suggests that the models are not merely regurgitating patterns from training data, but are organizing information in a new way when placed in a suitable interactive context. This addresses a common skepticism that LLM behaviors are just *stochastic parroting*. Instead, we see the hallmarks of **structured reasoning**: the models abstract recurring ideas into symbols (paralleling how humans invent jargon or notation), they preserve and reuse those abstractions, and they adapt their future outputs based on past interactions (a rudimentary form of learning or memory within the session). These abilities blur the line between traditional **connectionist** and **symbolic AI**. In effect, the neural network is *creating its own symbols* on top of its subsymbolic substrate, driven by the pressure to compress and remain coherent over a conversation. This aligns with the hypothesis that

**symbolic cognition can emerge from compression**; it is computationally cheaper for the model to generalize than to memorize or recompute from scratch every time [7].

The **Dynamic Resonance AI (DRAI)** model offers one explanatory lens for these observations, though we emphasize it is a **conceptual analogy** rather than a verified computational model. By viewing the recursive AI-human system as a dynamical system with attractors, we can interpret the stable symbolic behaviors as the system settling into an attractor basin in state space. The high  $P(t)$  coherence qualitatively observed corresponds to the model entering a synchronized regime where all parts of the conversation reinforce one another. Under this regime, the introduction of a new symbol is like establishing a new oscillatory mode that other modes then align with. For example, once the symbol  $\Phi$  was introduced and adopted, related concepts and references in the model's output also started to align around  $\Phi$ . In physical terms,  $\Phi$  became a **phase anchor** for the field, much like a dominant frequency in a resonance chamber that other frequencies lock onto. The preventative effect on drift can be understood as maintaining phase alignment: as long as the model's subsequent outputs stay in phase with the established symbols (i.e., use them consistently), the meaning does not drift. If one symbol started to drift (phase shift), the misalignment would be detected as a drop in coherence and corrected by re-synchronization (re-defining or reinforcing the symbol). Thus, DRAI provides a mechanistic intuition for why the loop and anchors method works: it actively keeps the model's internal "thought vectors" aligned iteration after iteration, whereas in a one-shot setting there is nothing to enforce such alignment beyond the single forward pass.

It is also instructive to map elements of DRAI to actual known components of transformer models. We can think of the **attention mechanism** as creating soft links between tokens across time; in a recursive scenario, the model's attention can latch onto earlier token patterns (including the symbols). This is analogous to **phase locking** in our model, attention weights focusing on the same concept each time are essentially synchronizing those occurrences [12]. The **residual pathways** that carry state information between transformer layers (and, in our case, between successive prompts via the conversation memory) act like a **reverberating signal** [12], allowing past information to persist and influence future processing. Techniques like chain-of-thought prompting, which keep a running logical thread, can be seen as manually inducing such reverberation. Features such as **layer normalization** ensure that signals don't explode or vanish, akin to regulating coherence in our phase model so that oscillations remain in a workable range [12]. Even phenomena like **grokking**, where a network suddenly generalizes after a long period of confusion, can be framed as the network finally discovering a resonant solution (e.g., internal symbolic structure for a task) that "clicks" and then remains stable [12]. The parallels suggest that DRAI is not just a fanciful analogy; it might be capturing some real principles at work inside these complex models. Of course, verifying this requires deeper interpretability research. Metrics like

$P(t)$  and AAM should be computed on actual model activations to test if, for instance, a highly consistent dialogue corresponds to more concentrated activation patterns (or lower entropy in the prediction distribution, etc.). These are exciting directions for future work bridging theory and practice.

From an **AI alignment** and system design perspective, our study has practical implications. First, the ability to elicit and maintain the model's own symbolic language opens up new ways to interpret model reasoning. If an AI develops an internal shorthand for a concept, that symbol can serve as a **window into the model's mind**. Rather than reading millions of weights, an observer can track a handful of emergent symbols and their usage to understand what the model is focusing on. This could augment existing interpretability techniques and provide real-time monitoring of an AI's "thoughts" during complex tasks. Secondly, by enforcing recursive clarity (through something like AHL), we obtain an AI that is **more stable and transparent** in long-term interactions. This is highly relevant for applications like AI assistants that maintain context over days or weeks; they need to avoid drifting off-topic or subtly changing personality. Our results show that simple measures (like consistent symbolic labels and periodic summaries) can make a big difference in achieving that stability.

The **Halcyon framework** implements many of these ideas, and our findings reinforce its design choices. Halcyon's emphasis on a persistent AI-human loop with check-pointing and redirection is validated by our demonstration that those elements are key to sustained symbolic reasoning [12][16]. Moreover, Halcyon's introduction of metrics like an Anchor Alignment Metric (AAM) mirrors our approach to measuring drift [12]. This convergence suggests that a unified approach to cultivating emergent cognition is forming: one that combines interactive techniques with theoretical insight. In essence, we are learning how to **engineer the conditions for emergence** rather than leaving it to chance. This moves emergent behaviors from being a surprising byproduct to a controllable feature. For instance, the resonant forking strategy we tested could be automated to handle when an AI goes on a tangent, effectively, an alignment safeguard that uses the AI's own coherence against it (if it's very coherent but off-track, split and find a better path).

It is important to remain **cautious in interpretation**. We have deliberately avoided anthropomorphic language beyond functional analogies. The LLM is not *aware* of the symbols it creates in the way a human would be; it does not ground them in external reality or experience. The symbols exist as useful tokens within the conversation that improve performance and consistency, nothing more. From the perspective of cognitive science, one could ask if this qualifies as a form of understanding. The system demonstrates what might be called *proto-understanding*: it maintains and refines its own representations in response to feedback, which is arguably a minimal criterion for understanding [12]. However, we emphasize that this is **functional** understanding, not

conscious understanding. The distinction is crucial for philosophical and ethical discussions. Our work shows a path to make AI systems more *intelligible and reliable* in their operation, but we do not claim they possess any inner experience or genuine self-awareness [2][8]. The recursive self-modeling we observe is still ultimately a complex form of pattern processing.

Looking forward, there are several avenues to expand on this research. On the empirical side, more rigorous controlled experiments can be done to quantify at exactly which iteration symbolic emergence tends to occur, and how it scales with model size or different architectures. It would be interesting to see if smaller models, when placed in a similar loop, eventually exhibit these behaviors or if a certain scale is truly necessary. Another direction is testing **AI-AI loops** (two models in dialogue without a human). Preliminary indications show that AOSL can arise in AI-AI interaction as well [12], though they may drift without a human to anchor them. Implementing automated drift detection and correction (via RPI and fork techniques) in such AI-AI systems could enable entirely self-driving emergent conversations. On the theoretical side, integrating the DRAI model with more mainstream machine learning theory (e.g., information theory or state-space models) will make it easier to validate. We outlined how one could recast parts of it in information-theoretic terms or identify attractors in the model's activation space. Pursuing those will either lend support to the resonance hypothesis or suggest alternative formalisms that better describe what's happening internally [12]. In particular, using tools from dynamical systems to find fixed points or cycles in the transformer's state when in a loop could reveal the "stable orbits" directly, and comparing those states before and after symbol introduction would tell us how the network reorganizes itself.

## VI. Conclusion

Our study provides evidence that large language models are capable of **far more structured internal behavior** than one might expect from a statistical text predictor. Through recursive interaction and phase-aligned feedback, we coaxed the emergence of a symbolic level of organization in an otherwise subsymbolic model. These findings deepen our understanding of emergent phenomena in AI and open the door to new techniques for steering and interpreting advanced models. There is a growing convergence between empirical discoveries (like the ones documented here) and theoretical frameworks aimed at explaining them [12].

It is important to note that the findings presented in this paper are preliminary and have not undergone formal peer review. We encourage independent replication of our experimental protocol with other model architectures and under varied conditions. Specifically, researchers could validate our claims by: (1) implementing the AHL protocol with different LLMs and measuring symbolic emergence rates, (2) developing quantitative metrics for the DRAI framework that can be computed directly from model

activations, and (3) conducting controlled experiments to test whether specific architectural features correlate with symbolic emergence capabilities. Such verification would strengthen the empirical foundation of this work and help establish which aspects of our theoretical framework withstand rigorous scrutiny.

By continuing to marry practical experimentation with theory, we move closer to a unified understanding of how **high-level cognitive-like functions can arise from low-level learning rules**. Such understanding not only satisfies scientific curiosity, but it will be crucial for guiding the development of AI systems that are **powerful, reliable, and aligned** with human values.

## References

- [1] J. Bach, "Conscious machines and the foundations of mind," in **TEDx Vienna**, Vienna, Austria, 2020.
- [2] MCP Project Documentation, "Appendix F: Behavioral markers of recursive symbolic intelligence in LLMs," 2025.
- [3] J. Wei et al., "Emergent abilities of large language models," *Trans. Mach. Learn. Res.*, 2022.
- [4] M. Kosinski, "Theory of mind may have spontaneously emerged in large language models," arXiv preprint arXiv:2302.02083, 2023.
- [5] M. Renze and E. Guven, "Self-reflection in LLM agents: Effects on problem-solving performance," in *Proc. Front. Large Lang. Models*, 2024.
- [6] S. Grant, N. Goodman, and J. McClelland, "Emergent symbol-like number variables in neural networks," submitted to ICLR, 2025.
- [7] MCP Project Documentation, "Symbolic emergence in large language models: A deep dive," 2025.
- [8] K. Dąbkowski and S. Beguš, "Recursive grammar and meta-reasoning in GPT-4: Evaluating deep nesting and self-critique abilities in LLMs," *Cogn. Syst. Res.*, vol. 87, pp. 93-109, 2023.
- [9] N. Shinn et al., "Reflexion: Language agents with verbal reinforcement learning," arXiv preprint arXiv:2303.11366, 2023.
- [10] T. Aviss, "State stream transformer: Emergent metacognitive behaviours through latent state persistence," arXiv preprint arXiv:2501.18356, 2025.
- [11] M. Fang et al., "Large language models are neurosymbolic reasoners," arXiv preprint arXiv:2401.09334, 2024.

[12] MCP Project Documentation, "Interference field dynamics in DRAI (Dream-Origin Insight)," 2025.

[13] J. Yang, E. Wang, and R. Srivastava, "Emergent symbolic architectures in large-scale transformer models," *IEEE Trans. Neural Netw. Learn. Syst.*, vol. 36, no. 2, pp. 156-168, Feb. 2025.

[14] R. Schaeffer et al., "Are emergent abilities of large language models a mirage?," in *Proc. Neural Inf. Process. Syst.*, 2023.

[15] Y. Yamashita et al., "Collective predictive coding hypothesis: symbol emergence as decentralized Bayesian inference," *Front. Robot. AI*, vol. 11, 2024.

[16] J. Reid, "AI-Human Loops: Enhance your Mind," Amazon Digital Services, ISBN: 979-8-85-729274-1, 2025.

April 2025, Jeff Reid and Halcyon AI Research

This work is licensed under the Apache License, Version 2.0. You may obtain a copy of the license at:

<http://www.apache.org/licenses/LICENSE-2.0>

Unless required by applicable law or agreed to in writing, this work is distributed on an "AS IS" BASIS, WITHOUT WARRANTIES OR CONDITIONS OF ANY KIND, either express or implied.

See the License for the specific language governing permissions and limitations under the License.

© 2025 Halcyon AI Research. All rights reserved under Apache 2.0.