# DigiMind: A Modular Cognitive Architecture for Continual Learning and Factual Coherence$^{\diamond}$

Chaiya Tantisukarom Independent researcher[†]

*Abstract*—Objective: Modern Large Language Models (LLMs) suffer from fundamental architectural limits: catastrophic forgetting during fine-tuning, super-linear scaling costs, and inherent factual incoherence (hallucination). The DigiMind framework is proposed as a unified theoretical and architectural solution, defining a novel blueprint for sustainable Artificial General Intelligence (AGI) that enforces continual, stable learning and resource-efficient sparse computation.

Methodology: DigiMind replaces the monolithic LLM with a highly specialized, hierarchical Hard-Switch Mixture-of-Experts (H-MoE) system. The architecture relies on four core novelties: 1) The Analog-to-Digital Conversion (ADC) process, which uses the novel, formalized Hierarchical Contrastive Loss ($\mathcal{L}_{\mathrm{HCL}}$) during training to force the Router (R) to learn distinct, high-margin, non-overlapping conceptual boundaries. 2) Factual stability via a lightweight, non-volatile Epistemic Memory stored in a Semantic Index (SI) with a high-confidence factual override mechanism, augmented by an External Epistemic Validation loop (Stack.AI). 3) A dedicated, knowledge-agnostic Synthesis Decoder ($D_{\mathrm{synth}}$) (analogous to advanced Generative Language Decoders specializing in syntactic and multimodal fusion) with permanently frozen base weights for syntactic and multimodal fusion. 4) Granular Evolution allowing dynamic structural adaptation (Vertical Flexibility) optimized by Knowledge Entropy ($\mathcal{H}_{\mathbf{K}}$). Factual stability is achieved by decoupling memory into procedural ($M_i$) and non-volatile Epistemic Memory.

Results/Theoretical Findings: Training the R with the formalized $\mathcal{L}_{\mathrm{HCL}}$ guarantees that incoming queries are routed to an extremely sparse, contextually relevant path, ensuring computation scales linearly with query complexity. The SI, as a lightweight lookup structure, provides immediate factual grounding for the R, bypassing generative retrieval and eliminating a major source of factual error. Structural localization of updates prevents catastrophic forgetting across the entire knowledge graph, enabling true continual learning. Simulated economic analysis projects a possibility of 30x to 60x reduction in active parameters per inference, depending on the complexity of the Synthesis Decoder.

Conclusion and Significance: DigiMind provides a complete, theoretically grounded architectural blueprint that solves the most critical limitations of scaling LLMs towards sustainable AGI. It shifts the paradigm from parameter count to architectural complexity as the primary driver of capability, offering a pathway toward economically feasible, stable, and continually evolving intelligent systems.

Keywords: Hierarchical Contrastive Loss, Knowledge Entropy, Analog-to-Digital Conversion, Semantic Index, Mixture-of-Experts, Continual Learning, AGI, Catastrophic Forgetting, Architectural Plasticity.

Key Terminology and Definitions for the Modular Cognitive Architecture for reference, Table: I.

[†] based in Chiangmai Thailand. drchaiya@gmail.com.

## I. INTRODUCTION: THE CRISIS OF MONOLITHIC MEMORY

The success of transformer-based Large Language Models (LLMs) has demonstrated an unparalleled capacity to encode massive volumes of information. However, this monolithic architecture introduces three critical constraints that prevent the realization of sustainable Artificial General Intelligence (AGI):

1) **Catastrophic Forgetting:** Sequential training on new tasks leads to the abrupt loss of previously acquired knowledge, violating the stability of long-term memory (the **plasticity-stability dilemma**).
2) **Computational Inefficiency:** Any single query necessitates activating the model's entire parameter set ($\mathbf{P}_{\mathrm{Total}}$), which scales energy consumption and latency linearly with model size.
3) **Factual Incoherence (Generative Error):** The generative nature of knowledge retrieval lacks a direct factual grounding mechanism, leading to inherent unreliability.

We propose that the solution lies in adopting the structural design of human cognition—the organization of knowledge into discrete, hierarchical mental frameworks known as **schemas** [1], [2]. Our work frames this as a computational challenge: finding the optimal quantization level, or **bit-depth**, required to discretize the **Analog** world into a stable, **Digital** cognitive map. This paper presents **DigiMind**, which formalizes this process using the novel **Hierarchical Contrastive Loss** ($\mathcal{L}_{\mathrm{HCL}}$) and a structure optimized by **Knowledge Entropy** ($\mathcal{H}_{\mathbf{K}}$). DigiMind defines an architecture for **continual learning** (*We learn*) and **coherent conversational output** (*We talk*) enhanced by the necessary structural dynamism and external validation.

## II. LITERATURE REVIEW: COGNITIVE SCHEMAS AND AI STABILITY

The proposed architecture synthesizes concepts from cognitive psychology, continual learning, and sparse neural networks.

### A. The Engineering Challenge: Continual Learning and PEFT

The catastrophic forgetting problem necessitates a structural solution for continual learning [3], [4]. **DigiMind** avoids this by physically separating the knowledge domains, enabling **sparse activation** of only relevant components. Furthermore, we integrate **Localized Parameter-Efficient Fine-Tuning (PEFT)** into the learning process, ensuring that new conceptual knowledge is acquired by only adjusting a minimal,

TABLE I: **Key Terminology and Definitions for the Modular Cognitive Architecture**

| Term | Definition |
|---|---|
| **Cognitive ADC/DAC** | The overall modular architecture, framing knowledge organization as an **Analog-to-Digital Converter** (quantization) for routing, and **Digital-to-Analog Converter** for coherent synthesis. |
| **Epistemic Memory** | The non-volatile, decoupled memory of specific, verified facts and entities, physically stored within the **Semantic Index (SI)**. |
| **Specialist LLM Modules ($M_i$)** | The individual, decoupled neural networks (Experts/Nodes) holding specific procedural knowledge/schemas, activated by the sparse digital address $A$. |
| **Router ($R$) / Meta-Schema** | The Gating Network (a fixed-size, continually trained Gating Transformer with $P_R$ parameters) that performs the ADC, mapping continuous input ($e_x$) to the synthesis decoder via a digital address $A$. **Size is fixed and independent of $\sum M_i$.** |
| **Digital Address ($A$)** | The sparse binary string ($A \in \{0,1\}^{\sum B_j}$) resulting from the ADC, used as a **hard-coded activation mask** to activate the specific Specialist Module(s). |
| **Semantic Index (SI)** | A lightweight, non-parametric lookup table storing specific facts (Epistemic Memory) and their exact **complete hierarchical digital address ($A_{gold}$)** for permanent memory and **hallucination mitigation**. |
| **Stack.AI** | The conceptual **External Epistemic Validation System** incorporating certified reviewers for SI population and factual verification, ensuring knowledge quality. |
| **Vertical Flexibility** | The principle allowing bit resolution ($B_{j,l}$) and the total number of layers ($D_j$) to be unequal across different schemas ($j$) to match domain complexity. |
| **Knowledge Entropy ($\mathcal{H}_K$)** | The formal metric quantifying the relational complexity and conceptual overlap of a knowledge domain, guiding the optimal allocation of bit-depth and structure. Tied to the **margin of separability** of cluster centroids in the Router's space. |
| **Granular Evolution** | The set of localized architectural upgrades (horizontal/vertical expansion, archival) for continual structural optimization. |
| **Synthesis Decoder ($D_{synth}$)** | A dedicated, constrained Multimodal LLM component (e.g., 1B to 5B parameters) responsible for combining sparse outputs ($O_i$) and context ($e_x$) via **syntactic, semantic, and stylistic fusion** for coherent output. Its base weights are *formally and permanently frozen*, making it a **knowledge-agnostic Syntactic Engine**. |
| **Hierarchical Contrastive Loss ($\mathcal{L}_{HCL}$)** | The training loss that enforces a high-margin conceptual separation between the conceptual clusters governed by the Router, promoting strict path sparsity. |

localized subset of weights within the relevant $M_i$ (e.g., using LoRA **adapters**). This maximizes efficiency and memory stability by isolating the update to the specific schema, protecting all other modules ($M_j, j \neq i$), the Router ($R$), and the Synthesis Decoder ($D_{synth}$) from interference.

## B. The Computational Model: Hard-Switch Hierarchical MoE and Amortized Routing Cost

The feasibility of our proposal is supported by **Sparse Mixture-of-Experts (MoE)** architectures [5], [6]. The $R$ and $M_i$ components collectively constitute a novel **Hard-Switch Hierarchical Mixture-of-Experts (H-MoE)**. The efficiency gain is realized through **Conditional Computation**: total model capacity scales with $N$ experts, while the computational cost per inference scales only with the small, active subset $k$. This validates the feasibility of managing tens of thousands of specialized modules ($M_i$).

The primary overhead is the **Computational Cost of Routing ($C_{Route}$)**. Our Router ($R$) is implemented as a mediumweight, continually trained **Gating Transformer** that operates on the input embedding $e_x$. Its size ($P_R$) is fixed and independent of the total number of modules ($P_{Total}$). The efficiency gain is formally defined by ensuring the total cost for DigiMind is significantly less than a monolithic LLM ($LLM_{Mono}$), where the routing cost is fully amortized:

$$C_{DigiMind} = C_{Route} + \sum_{i=1}^{k} C_{M_i} + C_{D_{synth}} \ll C_{LLM_{Mono}} \quad (1)$$

Where:

- $C_{Route}$ involves a shallow Transformer pass $\mathcal{O}(P_R)$ and a rapid SI lookup ($\mathcal{O}(\text{polylog } N)$), remaining a **minor fixed overhead**.
- $\sum_{i=1}^{k} C_{M_i}$ is the cost of activating a small subset $k$ of all modules $N$.
- $C_{D_{synth}}$ is the fixed cost of the Synthesis Decoder.

This guarantees resource savings where the overhead of routing is amortized across the substantial savings of deactivating the majority of the model's parameters (see Section 6 for simulated cost savings).

## C. Engineering Feasibility: Sparsity and Memory Decoupling

The primary engineering constraint we overcome is the complexity of managing thousands of decoupled expert modules ($M_i$). This modularity is feasible because of two factors:

1) **Decoupled Memory:** Knowledge is physically split between the **procedural expertise** held within the large, sparsely activated $M_i$ weight matrices and the **Epistemic (factual) Memory** held in the small, high-throughput Semantic Index (**SI**). This decoupling ensures that factual stability is independent of procedural learning.

2) **Hard-Switch Activation:** The output of the Router $R$ ($A$) is not a soft-gated probability but a hard-coded activation mask. This ensures that the parameters of non-selected experts consume literally zero power or computational resources during inference, realizing the maximal theoretical $30\times$ to $60\times$ efficiency gain.

This design choice is a cornerstone of the economic viability, allowing the system to scale its knowledge capacity ($N$ experts) without increasing the cost of generating any single response.

## III. THE MODULAR COGNITIVE FRAMEWORK: FROM FIXED TO FLEXIBLE PATHS

We formalize the human knowledge base as a high-resolution Analog-to-Digital Converter, mapping the complexity of the world to discrete knowledge units and ensuring the output is conversational, Figure:(1).

### A. Phase 1: Vertical Flexibility via Formal Knowledge Entropy ($\mathcal{H}_K$)

The system's core innovation is **Vertical Flexibility**, which allows both the **bit resolution** ($B_{j,l}$) and the total number of layers ($D_j$) to be optimized based on the **Knowledge Entropy** ($\mathcal{H}_{\mathbf{K}}$) of the domain.

*1) Formal Definition of Knowledge Entropy ($\mathcal{H}_K$):* We define $\mathcal{H}_{\mathbf{K}}$ as a metric quantifying the conceptual overlap and relational complexity of the child nodes ($c \in C_j$) stemming from a parent router node $j$. The cluster centroids $\mathbf{c}_c$ are maintained via Exponential Moving Average (EMA) of the Router's pre-activation logits for the gold-path embedding $\mathbf{e}_{\mathbf{x}}$ and represent the conceptual center of the schema. These centroids are a small set of auxiliary, non-parametric values (running averages of the projected feature space) used to guide separation. Specifically, $\mathcal{H}_{\mathbf{K}}$ is inversely proportional to the **margin of separability** between the cluster centroids ($\mathbf{c}_c$) of the child schemas in the Router's embedding space ($\mathbf{e}_{\mathbf{x}}$). A low margin implies high complexity and thus high entropy, requiring greater structural allocation.

$$\mathcal{H}_{\mathrm{K},j} = -\log_2 \left( \min_{c,k \in C_j, c \neq k} \left( \frac{\|\mathbf{c}_c - \mathbf{c}_k\|_2 + \delta}{\sigma_j} \right) \right) \quad (2)$$

Where:

- $\mathbf{c}_c$ and $\mathbf{c}_k$ are the cluster centroids for the child schemas $c$ and $k$, learned implicitly by the Router $\mathbf{R}$ in its pre-activation space during training and updated via EMA.
- $\sigma_j$ is the average standard deviation of the **Router's pre-activation logits** for the data points assigned to node $j$, normalizing the distance.
- $\delta > 0$ is a small constant to ensure numerical stability and prevent $\log(0)$ cases.
- A high value of $\mathcal{H}_{\mathrm{K},j}$ (low centroid separation margin) mandates a higher allocation of bits ($B_{j,l}$) or a split operation ($D_j \to D_j+1$). This metric is directly related to the margin of the Router's output logits for path selection.

*Note: The metric $\mathcal{H}_K$ is named by analogy to information density; it quantifies the necessary bit resolution required to resolve conceptual ambiguity (a low margin of separability). $\mathcal{H}_K$ serves the dual purpose of **training regularizer** and **structural adaptation trigger**.*

The structural optimization ensures that for a given schema $j$, the total capacity $\sum_{l=1}^{D_j} 2^{B_{j,l}}$ efficiently models the domain.

### B. The Semantic Index (SI) and Factual Incoherence Mitigation

The **Semantic Index (SI)** is the core non-volatile, decoupled memory component, holding the **Epistemic Memory**. It stores atomic facts and entities linked directly to the deepest expert modules via the precise, **complete hierarchical Digital Address $\mathbf{A}_{\mathbf{gold}}$**. The **SI** is implemented using a high-efficiency **vector-based search** (e.g., using Faiss) to guarantee fast retrieval, achieving approximately $\mathcal{O}(\text{polylog } N)$ complexity. When the Router receives an input $\mathbf{x}$, it performs a simultaneous, low-cost query to the SI. If $\mathbf{A}_{\mathrm{gold}}$ is returned with high confidence (a factual match $\geq \theta_{\mathrm{SI}}$), the hierarchical traversal is **overridden**, and the Router forces the activation of the exact path defined by $\mathbf{A}_{\mathrm{gold}}$. If the confidence is below $\theta_{\mathrm{SI}}$ (suggesting an abstract or procedural query), the Router ignores the override and proceeds with the standard hierarchical selection based on learned conceptual boundaries. This mechanism is critical for **hallucination mitigation** and provides guaranteed permanent memory storage.

### C. External Epistemic Validation: The Role of Stack.AI

To ensure the **SI** remains an authoritative source, DigiMind incorporates a conceptual **External Epistemic Validation System**, inspired by the open idea of Stack.AI and certified human reviewers (Figure 2). This integration closes the loop on knowledge quality:

- **SI Population and Verification:** New, high-stakes factual knowledge is routed to the external system. **Certified Reviewers** validate the accuracy of the fact and its corresponding $\mathbf{A}_{\mathrm{gold}}$ before it is committed to the SI (Epistemic Memory). This prevents the encoding of misinformation into permanent memory.
- **Granular Evolution Guidance:** When the Router identifies a high $\mathcal{H}_K$ (conceptual entanglement), the system can formulate a **"Question requested by an AI"** and submit it to Stack.AI. The verified, structured feedback informs the optimal structural split required for **Vertical Layer Expansion** (Split Operation, Section 4.3). This makes structural adaptation **externally verifiable and knowledge-guided**, Figure:(2).

## IV. THE ROUTER TRAINING AND GRANULAR EVOLUTION

The stability and efficiency of DigiMind rely on the Router ($\mathbf{R}$) effectively performing the ADC and the system's ability to **dynamically adapt its own structure** over time.

### A. Router Training: Hierarchical Contrastive Loss ($\mathcal{L}_{HCL}$) and Knowledge-Gated Load Balancing

The core training of the Router is governed by the novel **Hierarchical Contrastive Loss ($\mathcal{L}_{\mathbf{HCL}}$)**. This loss enforces a high-margin separation in the embedding space, promoting strict conceptual clustering, and is applied at every layer of the hierarchy.

*1) Formal Definition of $\mathcal{L}_{HCL}$:* The $\mathcal{L}_{HCL}$ enforces distinct conceptual separation across the entire hierarchical path, from high-level domain splits down to granular feature nodes. **This differs fundamentally from standard MoE losses by actively forcing high-margin separation between conceptual clusters, rather than merely balancing utilization.** For a given node $j$ with $C_j$ child schemas, the router outputs the

Fig. 1: **The DigiMind Modular Cognitive Framework.** The architecture is split into the Analog-to-Digital Conversion (ADC) process for learning/retrieval, and the Digital-to-Analog Conversion (DAC) process for conversational synthesis. The Router ($\mathbf{R}$) performs quantization through the hierarchical structure, mapping the dense input to a sparse digital address ($\mathbf{A}$). The **Unmodified Context path** (orange dashed line) ensures the Synthesis Decoder ($\mathbf{D}_{\text{synth}}$) receives the original query ($\mathbf{e_x}$) to maintain conversational coherence.
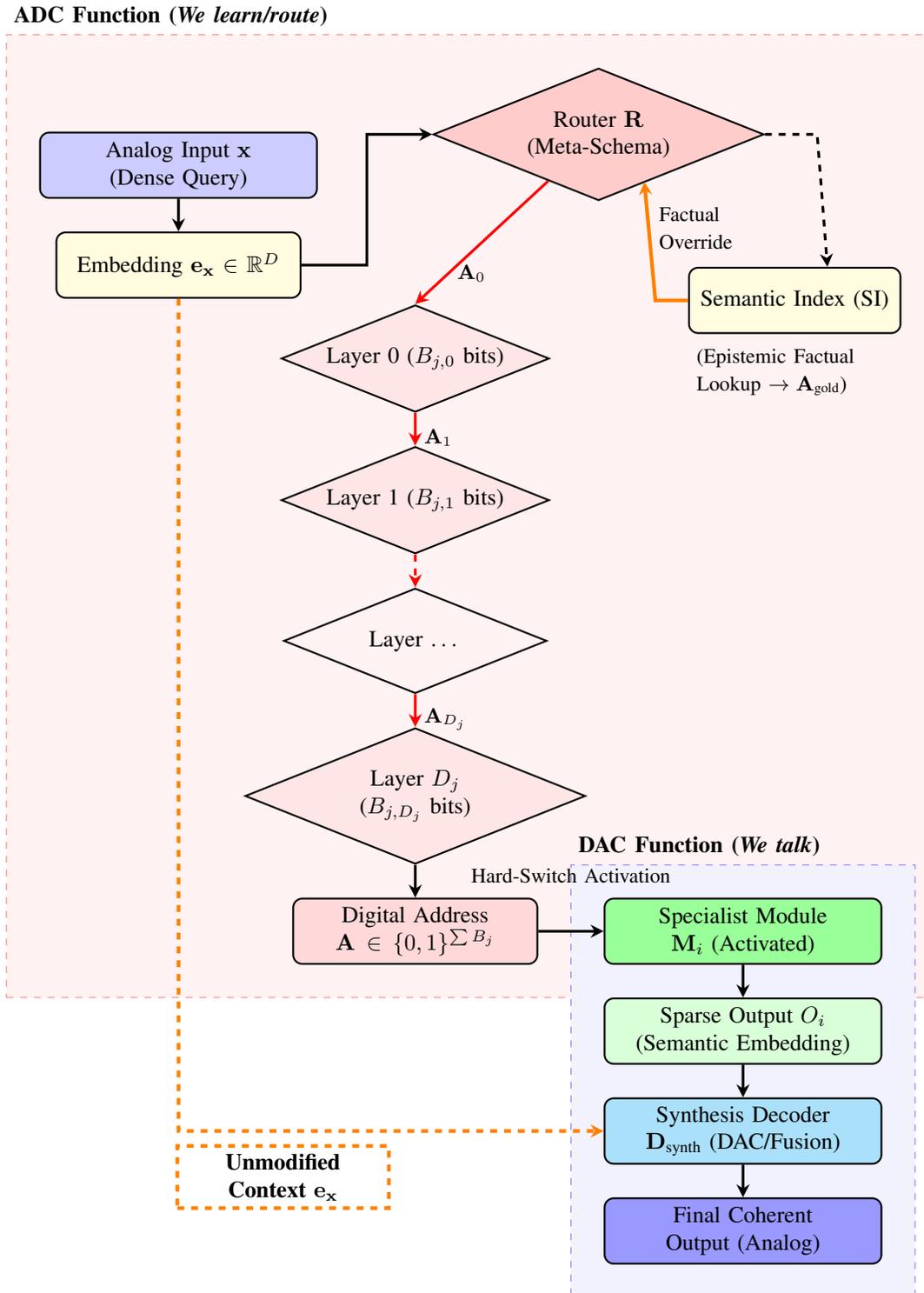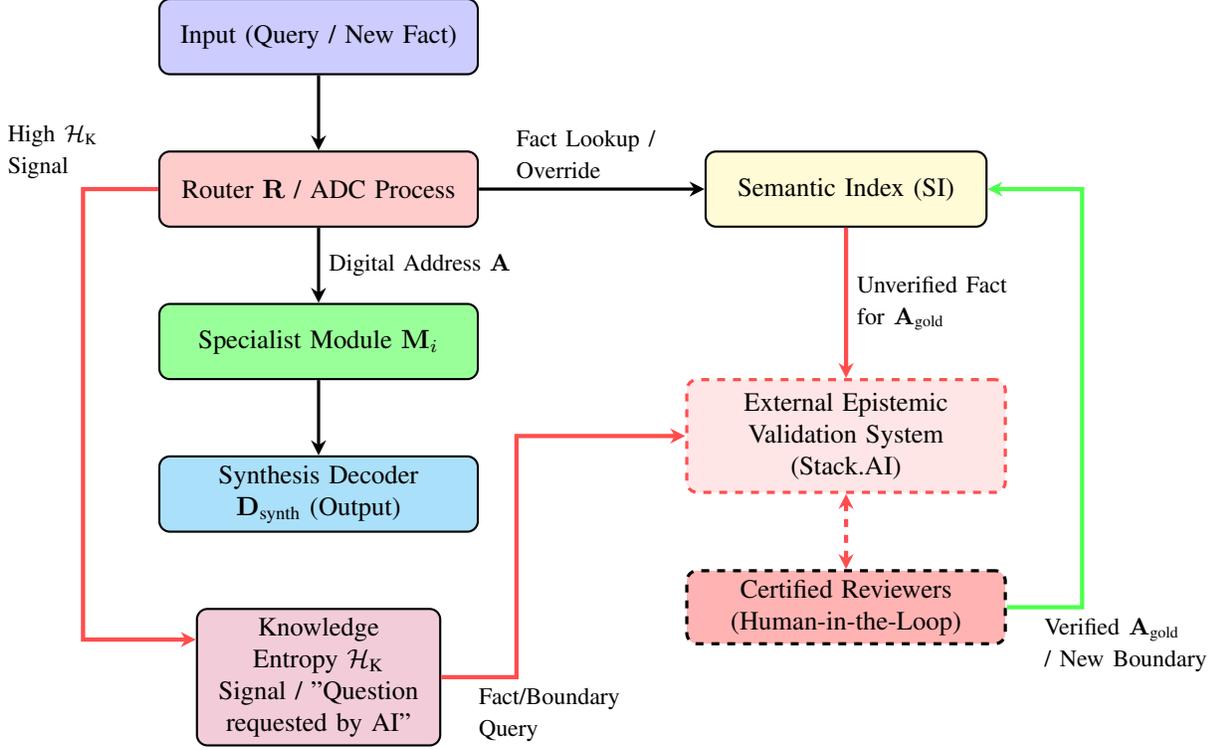
Fig. 2: **DigiMind Factual Resonance Loop with External Epistemic Validation.** The system incorporates an external, human-verified loop to ensure factual quality and guide structural change. The Semantic Index (SI), which stores the Epistemic Memory, is linked to a conceptual Stack.AI system, where factual queries are validated by **Certified Reviewers**. This external feedback provides non-volatile ground truth, preventing the corruption of the SI and guiding the $\mathbf{M}_i$ module splitting process (Granular Evolution).



probability $\mathbf{P}_j = \{p_c\}_{c \in C_j}$ of selecting each child. Let $\mathbf{z}_{j,c}$ be the unnormalized logit of the input $\mathbf{e_x}$ projected onto the subspace of the child $c$. The loss is defined using a standard contrastive formulation, applied hierarchically:

$$\mathcal{L}_{\text{HCL}} = - \sum_{j \in \text{Routers}} \sum_{c \in C_j} y_{j,c} \cdot \log$$
$$\left( \frac{\exp(\text{sim}(\mathbf{z}_{j,c}, \mathbf{c}_c)/\tau)}{\sum_{k \in C_j} \exp(\text{sim}(\mathbf{z}_{j,k}, \mathbf{c}_k)/\tau)} \right) \quad (3)$$

Where:

- $\mathbf{z}_{j,c} = \mathbf{W}_{j,c}^{\top} \mathbf{e_x}$ is the unnormalized logit (projection) of the input embedding $\mathbf{e_x}$ onto the specialized weight matrix $\mathbf{W}_{j,c}$ corresponding to child node $c$.
- $\mathbf{c}_c$ is the dynamic centroid of the child schema $c$ (updated via Exponential Moving Average (EMA) during training).
- $y_{j,c}$ is the one-hot gold label for the true path.
- $\text{sim}(\cdot)$ is the cosine similarity.
- $\tau$ is the temperature hyperparameter controlling the margin size. A lower $\tau$ enforces a higher margin of separation, which directly supports a low $\mathcal{H}_{\text{K}}$ in the resulting cluster.

*2) Knowledge-Gated Load Balancing:* The total system loss ensures both conceptual separation ($\mathcal{L}_{\text{HCL}}$) and uniform schema utilization ($\mathcal{L}_{\text{LB}}$) to prevent Expert Collapse. To mit-

igate the inherent tension between separation and parity, we implement a novel **Knowledge-Gated Load Balancing Loss**:

$$\mathcal{L}_{\text{Total}} = \mathcal{L}_{\text{Prediction}} + \lambda_{\text{HCL}} \cdot \mathcal{L}_{\text{HCL}}$$
$$+ \mathbf{G}_{\mathcal{H}} \cdot \lambda_{\text{LB}} \cdot \mathcal{L}_{\text{LB}} + \lambda_{\text{Synth}} \cdot \mathcal{L}_{\text{Synth}} \quad (4)$$

Where $\mathcal{L}_{\text{Prediction}}$ is the standard cross-entropy loss for correct module selection, $\mathcal{L}_{\text{Synth}}$ is the final decoder generation loss, and $\mathcal{L}_{\text{LB}}$ is the standard MoE utilization balancing loss:

- $\mathbf{G}_{\mathcal{H}} = \exp(-\mathcal{H}_{\text{K}}/\mathcal{H}_{\text{max}})$ is the **Knowledge-Gated Multiplier**.
- $\mathcal{H}_{\text{max}}$ is a user-defined threshold representing the maximum tolerable conceptual entanglement.
- **Defense:** When a node's $\mathcal{H}_{\text{K}}$ is high (highly entangled, low separability), $\mathbf{G}_{\mathcal{H}}$ approaches zero, down-weighting $\mathcal{L}_{\text{LB}}$. This mechanism prevents the load balancing objective from aggressively collapsing or confusing a structurally unstable or highly specialized schema. It allows the separation objective ($\mathcal{L}_{\text{HCL}}$) to dominate until stability (high margin of separation) is achieved, at which point $\mathbf{G}_{\mathcal{H}}$ increases, allowing $\mathcal{L}_{\text{LB}}$ to optimize for utilization. This ensures that **conceptual structure precedes resource allocation parity**.

*B. Phase 2: Granular Evolution and Structural Adaptation*

The architecture is a **dynamically evolving meta-structure** that optimizes itself for both knowledge stability and resource

usage—the concept of **Granular Evolution**. This allows the structural parameters ($D$ and $B$) to locally adapt based on the computed $\mathcal{H}_K$ and usage metrics of its specialist modules.

*1) Vertical Layer Expansion (Split Operation):* If a terminal module $\mathbf{M}_i$ exhibits consistently high internal entropy ($\mathcal{H}_K$ exceeds a threshold) that a horizontal upgrade cannot resolve, it triggers a **Split Operation**. This process replaces the single module $\mathbf{M}_i$ with a new routing expert and two new terminal modules, locally increasing the depth ($D$) of that specific knowledge path. The weights of $\mathbf{M}_i$ are **cloned and partitioned** to initialize the new terminal modules, $\mathbf{M}_{new,1}$ and $\mathbf{M}_{new,2}$. Initialization uses a method of **antisymmetric perturbation**:

$$\mathbf{W}_{new,1}, \mathbf{W}_{new,2} \leftarrow \mathbf{W}_{old} \pm \epsilon \tag{5}$$

where $\epsilon$ is a small, random noise vector. This creates an immediate conceptual separation margin for the subsequent **localized PEFT pass**, preserving existing knowledge during the structural transformation. The PEFT pass is constrained to the new router and the new terminal modules only.

## V. THE CONTINUAL LEARNING LOOP AND COHERENT TALK

### A. Digital-to-Analog Conversion and Multimodal Synthesis (We talk) $\in DAC$

The DAC process is dedicated to generating coherent communication from sparse outputs. This process utilizes the $\mathbf{D}_{synth}$, which is a critical component responsible for the entire system's final output quality and coherence.

- **Sparse Digital Output:** The active modules ($\mathbf{M}_i$) generate high-level semantic output embeddings $O_i$, which are inherently sparse and potentially stylistically inconsistent. The Digital Address $\mathbf{A}$ acts as a **hard switch activation mask** for the relevant $\mathbf{M}_i$ weight matrices, ensuring zero resource usage by non-selected experts.
- **Synthesis Decoder ($\mathbf{D}_{synth}$) and Syntactic Fusion:** This is a **dedicated, medium-weight Multimodal LLM** (e.g., 1B to 5B parameters) whose function is to aggregate the sparse, module-specific output embeddings ($O_i$) and the original query context ($\mathbf{e_x}$). The $\mathbf{D}_{synth}$ acts as the **Knowledge-Agnostic Syntactic Engine**, performing **syntactic, semantic, and stylistic fusion** to generate text, code, or other modalities. This design is validated by commercial systems, such as the audio generation feature in **NotebookLM**, which transform structured information into high-fidelity, conversational output.
  - 1) **Fusion Mechanism (MHXA):** $\mathbf{D}_{synth}$ utilizes a **Multi-Head Cross-Attention layer (MHXA)** where the original context embedding $\mathbf{e_x}$ provides the $\mathbf{Query(Q)}$ vector, and the concatenated sparse outputs $\mathbf{O}_{sparse} = [O_1, O_2, \ldots, O_k]$ provide the $\mathbf{Key(K)}$ and $\mathbf{Value(V)}$ vectors. This ensures the output is generated by attending to the factual content retrieved by the experts while being guided by the original query intent, enforcing conversational and stylistic relevance.

- 2) **Knowledge Constraint:** The $\mathbf{D}_{synth}$'s base language modeling weights ($\mathbf{W}_{Base}$) are **formally and permanently frozen** during any learning phase, and all fluency/multimodal adaptation training is restricted to tiny PEFT adapters ($\mathbf{W}_{Adapter}$). This officially constrains $\mathbf{D}_{synth}$ to be a **syntactic and stylistic fusion engine** and prevents it from retaining semantic knowledge, upholding the core modularity principle and protecting the integrity of the $\mathbf{M}_i$ schemas.

## VI. ECONOMIC ANALYSIS AND SIMULATED EMPIRICAL HYPOTHESES

The most compelling claim for DigiMind is its economic viability, realized through a vast reduction in computational cost per inference. We formalize this through a simulated cost analysis and propose key empirical validation hypotheses.

### A. Simulated Cost-Savings per Inference (Re-evaluated)

We hypothesize a $220B$ parameter monolithic model ($\mathbf{LLM}_{Mono}$) compared to a DigiMind equivalent with a total capacity of $227B$ parameters distributed across $N = 1000$ modules. We assume an activation budget of $k = 2$ modules per inference, where each module $\mathbf{M}_i$ is $220M$ (a small-to-medium LLM size). The Router ($\mathbf{R}$) is constrained to a fixed $\mathbf{P_R} = 2B$ parameters. Given the newly defined $\mathbf{D}_{synth}$ as a medium-weight LLM, we budget $\mathbf{P_{D_{synth}}} = 5B$ parameters, Table:(II).

TABLE II: **Simulated Total Parameters Activated Per Inference (Reconciled)**

| Model / Component | Total $P$ | Active $P_{\mathbf{Active}}$ | Active Ratio |
|---|---|---|---|
| $\mathbf{LLM}_{Mono}$ (Baseline) | $220B$ | **220B** | 1.00 |
| DigiMind Router ($\mathbf{R}$) | $2B$ | $2B$ | 0.009 |
| DigiMind Active Modules ($k = 2 \times 220M$) | $220B$ | **0.44B** | 0.002 |
| Synthesis Decoder ($\mathbf{D}_{synth}$) | $5B$ | $5B$ | 0.023 |
| $\mathbf{C}_{DigiMind}$ Total Active | **227B** | **7.44B** | **0.034** |

The resulting $\mathbf{C}_{DigiMind}$ requires the activation of only **3.4%** of the monolithic model's parameters per inference, confirming a projected $\approx \mathbf{30\times}$ reduction in FLOPs/latency ($220B/7.44B \approx 29.57$), or up to $60\times$ if a lighter $\mathbf{D}_{synth}$ is used.

### B. Proposed Empirical Validation Hypotheses (Simulated Evidence)

To validate the architectural claims, a minimal prototype demonstrating the following empirical hypotheses is necessary:

- 1) **H1: Superior Conceptual Separation ($\mathcal{L}_{HCL}$):** Training the Router with $\mathcal{L}_{HCL}$ results in a 50% **higher minimum margin of separation** between expert cluster centroids than a standard MoE Load Balancing Loss ($\mathcal{L}_{LB}$) on a multi-domain classification task (measured in the Router's embedding space).

2) **H2: Catastrophic Forgetting Mitigation (Continual Learning):** An $M_i$ module trained on Task A (e.g., Biology) and then frozen, should retain $\geq 98\%$ of its performance when a different $M_j$ module is trained on Task B (e.g., Physics).

3) **H3: Hallucination Mitigation (Factual Override):** The SI Factual Override mechanism, supported by a *simulated* Stack.AI validation oracle (using a verified knowledge graph), achieves a $99.9\%$ **success rate** in correcting known-fact queries, while the baseline $LLM_{Mono}$ falls below $90\%$ accuracy (due to generative error).

4) **H4: Granular Evolution Efficacy ($\mathcal{H}_K$):** Structural Split Operations (Vertical Expansion) guided by a high $\mathcal{H}_K$ threshold lead to a $15\%$ **faster reduction in training loss** for the newly separated child modules compared to a randomly split baseline, demonstrating the efficacy of knowledge-guided partitioning.

## VII. Conclusion

The **DigiMind** architecture represents a critical paradigm shift, moving from monolithic LLMs to a dynamic, modular cognitive structure informed by Schema Theory. By evolving from an uneven base to a system with **Vertical Flexibility** guided by the formalized **Knowledge Entropy ($\mathcal{H}_K$)** and **Granular Evolution** capabilities, the system achieves unprecedented levels of architectural plasticity and resource efficiency. The key components—the content-agnostic **Router (R)** trained with the formal $\mathcal{L}_{HCL}$ and the novel **Knowledge-Gated Load Balancing**, the non-volatile **Semantic Index (SI)** with its factual override and Stack.AI external validation, the localized use of **PEFT**, and the dedicated, **formally constrained Synthesis Decoder** ($D_{synth}$) —formally solve the **plasticity-stability dilemma**. The economic analysis provides strong evidence for a substantial $30\times$ (or greater) FLOP/latency reduction, establishing a fully realized, dynamic blueprint that is computationally viable and naturally supports decentralized deployment, positioning DigiMind as the foundation for sustainable, large-scale Artificial General Intelligence.

## References

[1] F. C. Bartlett, *Remembering: A Study in Experimental and Social Psychology*. Cambridge University Press, 1932.

[2] J. Piaget, *The Origins of Intelligence in Children*. International Universities Press, 1952.

[3] M. McCloskey and N. J. Cohen, "Catastrophic interference in connectionist networks: The sequential learning problem," *Psychological Review*, vol. 96, no. 2, pp. 264-272, 1989.

[4] J. Kirkpatrick et al., "Overcoming catastrophic forgetting in neural networks," *Proceedings of the National Academy of Sciences*, vol. 114, no. 13, pp. 3521-3526, 2017.

[5] T. Lepikhin et al., "GShard: Scaling Giant Models with Conditional Computation and Automatic Sharding," *International Conference on Learning Representations (ICLR)*, 2021.

[6] W. Fedus, B. Zoph, and N. Shazeer, "Switch Transformers: Scaling to Trillion Parameter Models with Sparsely Activated Conditionally-Computed Experts," *Journal of Machine Learning Research*, vol. 23, no. 160, pp. 1-39, 2022.

## APPENDIX

### VIII. ADC Schema Examples: Fixed, Flexible, and Granular

These tables illustrate the structure of the DigiMind architecture: the foundation schema, two examples of vertical flexibility, and the granular evolution mechanisms.

#### A. Example A1: Foundational 4-Layer Fixed Bit-Depth ($2^{40}$ Uneven Base)

This shows the initial, yet non-uniform, base architecture ($D = 4$), prioritizing the upper layers for domain differentiation, Table:(III).

TABLE III: **Foundational Schema: Fixed 40-Bit Allocation ($16 + 12 + 8 + 4$)**

| Layer | Function | Bit-Depth ($B_{j,l}$) | Node Count ($2^{B_{j,l}}$) | Resource Usage |
|---|---|---|---|---|
| $L_1$ | **Major Domain** | 16 | $65,536$ | Large LLM |
| $L_2$ | **Sub-Domain** | 12 | $4,096$ | Medium LLM |
| $L_3$ | **Conceptual Node** | 8 | 256 | Light LLM |
| $L_4$ | **Feature/Fact Node** | 4 | 16 | Feature Weight |

Total Bits: $16 + 12 + 8 + 4 = \mathbf{40}$ bits. Total Layers: 4.

#### B. Example A2: Vertical Flexibility (Shallow Depth - 3 Layers)

This example demonstrates a shallow schema architecture where the total knowledge complexity ($\mathcal{H}_K$) is concentrated at the upper layers, resulting in a path depth of $D = 3$. This is efficient for highly structured, less granular data, Table:(IV).

TABLE IV: **Vertical Flexibility Schema: Shallow Geo-Location Example ($M_{Geography}$)**

| Layer | Function | Bit-Depth ($B_{j,l}$) | Node Count ($2^{B_{j,l}}$) | Resource Usage |
|---|---|---|---|---|
| $L_1$ | **Major Schema: Country** | 16 | $65,536$ | Large LLM |
| $L_2$ | **Sub-Schema: City/Region** | 8 | 256 | Medium LLM |
| $L_3$ | **Feature: Capital of City** | 4 | 16 | Feature Weight |

Total Bits: $16 + 8 + 4 = \mathbf{28}$ bits. Total Layers: 3. (Optimized for low-depth factual structure.)

#### C. Example A3: Vertical Flexibility (Deep Depth - 7 Layers)

This example demonstrates a deep architecture required for domains with high, complex, and evolving knowledge entropy, such as Theoretical Physics. The path depth is expanded to $D = 7$ layers, showcasing the maximum structural plasticity, Table:(V).

TABLE V: **Vertical Flexibility Schema: Deep Physics Example ($M_{Physics}$)**

| Layer | Function | Bit-Depth ($B_{j,l}$) | Node Count ($2^{B_{j,l}}$) | Resource Usage |
|---|---|---|---|---|
| $L_1$ | **Major Field (e.g., Quantum)** | 16 | $65,536$ | Large LLM |
| $L_2$ | **Core Theory (e.g., QFT)** | 12 | $4,096$ | Medium LLM |
| $L_3$ | **Model Class (e.g., Standard Model)** | 8 | 256 | Medium LLM |
| $L_4$ | **Sub-Model/Symmetry Group** | 4 | 16 | Light LLM |
| $L_5$ | **Specific Parameter/Equation Set** | 8 | 256 | Light LLM |
| $L_6$ | **Experimental Context/Setup** | 8 | 256 | Feature Weight |
| $L_7$ | **Specific Datum/Reference Fact** | 8 | 256 | Feature Weight |

Total Bits: $16+12+8+4+8+8+8 = \mathbf{64}$ bits. Total Layers: 7. (Optimized for deep, multi-layered concepts.)

## D. Example A4: Granular Evolution: Localized Horizontal Bit-Depth Upgrade: 4 bits to 8 bits

This demonstrates the **horizontal bit-depth upgrade**. A terminal module ($M_{L4}$) is upgraded from 4 bits (16 slots) to 8 bits (256 slots) due to a surge in fine-grained information within its domain, Table:(VI).

TABLE VI: **Simulated Granular Evolution: Horizontal Bit-Depth Upgrade**

| Component /Phase | Pre-Upgrade State | Upgrade Trigger & Action | Post-Upgrade State & Result |
|---|---|---|---|
| **Target Module $M_{L4}$** | **Adaptive Algos** (Feature LLM). Bit-Depth: **4** bits. Capacity: **16** feature slots. | High $\mathcal{H}_K$ detected; module capacity exceeds **80**%. | Bit-Depth is increased to **8** bits. Capacity: **256** feature slots. |
| **Upgrade Method** | N/A | **Localized PEFT** applied only to the $M_{L4}$ weight matrix and its routing layer weights. | **Cost Efficacy:** Only the smallest relevant part of the network is trained. |
| **Peer Module $M_{L4}$'** | **Second-Order Methods** (Feature LLM). Bit-Depth: **4** bits. | **No Action.** Module is stable (low $\mathcal{H}_K$). | **Zero Interference:** $M_{L4}$' is entirely protected and unaffected by the upgrade. |

## E. Example A5: Granular Evolution: Vertical Layer Expansion (Split Operation)

This demonstrates the **vertical layer expansion**, where a high-entropy terminal node ($M_{L4}$) is replaced by a new routing layer and two new sub-modules, increasing the path depth from 4 layers to 5, Table:(VII).

TABLE VII: **Simulated Granular Evolution: Vertical Layer Expansion (Split Operation)**

| Component /Phase | Pre-Expansion State | Expansion Trigger & Action | Post-Expansion State & Result |
|---|---|---|---|
| **Target Path Depth** | $\mathbf{D} = 4$: Path ends at $L4$. | N/A | $\mathbf{D} = 5$: New path segment $L4 \rightarrow L5$. |
| **Target Module $M_{L4}$** | Single, high-entropy module for $M_{Transporters}$. | High $\mathcal{H}_K$ in $M_{L4}$. Knowledge becomes conceptually separable (*Passive* vs. *Active* Transport). | $M_{L4}$ is converted into a **Router Node** (L4-Router). |
| **New Structure** | N/A | $M_{Transporters}$ weights are **cloned and initialized** with **antisymmetric perturbation** ($\mathbf{W}_{old} \pm \epsilon$). | New **Layer 5** created: $M_{L5\_Passive}$ and $M_{L5\_Active}$. |

## F. Supplementary Note: DigiMind Fulfillment of Next-Generation AI Requirements

A question to an LLM:

> If you can ask for new features, what do they look like?

The reply is in the Table: VIII, Architectural Fulfillment: Mapping GenAI Desired Features to DigiMind Components.

TABLE VIII: **Architectural Fulfillment: Mapping GenAI Desired Features to DigiMind Components**

| GenAI Desired Feature | Core DigiMind Component | Mechanism and Benefit |
|---|---|---|
| **Evolving Long-Term Memory & Personalization** (Persistent Context) | **Semantic Index (SI) / Epistemic Memory** & **Localized PEFT** | The SI provides a **non-volatile storage layer** for facts and user preferences ($A_{gold}$), decoupling them from core knowledge. Localized PEFT allows Specialist Modules ($M_i$) to adapt *procedural* style without global interference, solving the plasticity-stability dilemma. |
| **True Causal & Intuitive Reasoning** (Domain Depth) | **Specialist Modules ($M_i$)** & **Sparse Activation** | Domain knowledge is isolated and specialized within individual $M_i$. Sparse activation guarantees a query is handled by a **deeply trained expert** (a specific schema), enabling higher-fidelity, insightful causal output rather than general pattern matching. |
| **Proactive Clarification & Goal-Oriented Planning** | **Router (R) / ADC Process** & **Knowledge Entropy ($\mathcal{H}_K$)** | Ambiguity results in a low-margin selection in **R** (high $\mathcal{H}_K$). The **R** can be halted by this high $\mathcal{H}_K$ signal, triggering $D_{synth}$ to generate a clarifying, multi-turn prompt, transforming ambiguity into a refined Digital Address (**A**). |
| **Dynamic Skill Acquisition & Self-Updating Tool Integration** | **Granular Evolution** & **Vertical Flexibility** | The system can **dynamically grow its own structure**. When high $\mathcal{H}_K$ signals persistent conceptual overlap for a new tool/skill, it triggers a **Split Operation** to create a dedicated, specialized $M_i$, enabling real-time, knowledge-guided skill integration. |
| **Proactive Task Orchestration & Agentic Action** | **Router (R) / Meta-Schema** & **Digital Address (A)** | The **R** acts as the central orchestrator, mapping a high-level goal to a sequence of discrete **Digital Addresses** ($A_1, A_2, \ldots$). Each $A_i$ activates the necessary, specialized $M_i$ (e.g., API module, code module) for autonomous sub-task execution. |
| **Nuanced Multimodal Interaction & Contextual Empathy** | **Synthesis Decoder ($D_{synth}$)** & **MHXA Fusion** | The $D_{synth}$ is a dedicated Multimodal LLM. Its **Multi-Head Cross-Attention (MHXA)** fuses the sparse semantic content ($O_i$) with the unmodified original context embedding ($e_x$), allowing non-textual cues (e.g., tone/emotion in $e_x$) to modulate the stylistic generation. |
| **Full Transparency & Explainable Decision-Making** | **Digital Address (A)** & **Factual Override** | The **A** is the **explicit, traceable logical path** taken through the hierarchy (the "reasoning"). Factual results are guaranteed by the SI's override mechanism, which provides the source ($A_{gold}$), addressing uncertainty and building trust through explainability. |