

Structural Homological Resolution with Dual-Stream Parameter Compression: A Unified Framework for Mathematical Reasoning under Extreme Resource Constraints

Abstract

This paper proposes a novel mathematical reasoning framework—Structural Homological Resolution with Dual-Stream Parameter Compression (SHR-DS)—aimed at addressing the fundamental challenge of building high-performance mathematical reasoning models under extreme hardware constraints (8 GB VRAM, limited system memory, no hardware expansion). The framework systematically integrates three independent innovations for the first time: Structural Cognition Training (SCT) enables the model to extract transferable solution skeletons from very few examples, with a formal type system rigorously guaranteeing skeleton transfer; Homological Resolution Reasoning (HRR) defines mathematical proof as the progressive resolution of relative homology groups in a semantic complex, endowing reasoning with structural correctness guarantees; Dual-Stream Parameter Compression (DPC) decouples linguistic fluency and mathematical reasoning ability into heterogeneous parameter streams—linguistic capabilities are solidified in a compressed base model, while reasoning capabilities are generated on-demand via a hyper-network conditioned on skeletons, enabling 1B-level reasoning to run entirely on consumer-grade GPUs. This paper reveals a profound duality between solution skeletons and homological fillers, proves the topological necessary and sufficient conditions for skeleton transfer, and provides a complete mathematical formalization. Based on calibrated data from adjacent existing technologies, we estimate that under 100M–1B parameter scale, SHR-DS can improve sample efficiency by 10–50 \times , increase mathematical reasoning accuracy by 20%–35%, while retaining over 95% of the dialogue capability of the source language model. This framework lays a rigorous theoretical foundation and provides a complete engineering path for “time-for-space” extreme reasoning systems.

Keywords: Mathematical reasoning; algebraic topology; homological resolution; skeleton extraction; hyper-network; parameter-efficient training; few-shot learning; knowledge decoupling

1 Introduction

1.1 Background and Motivation

Large language models have made rapid progress in mathematical reasoning tasks. From Chain-of-Thought (CoT) to inference-time search (Tree-of-Thought, Monte Carlo Tree Search), from continuous relaxation models (e.g., Potential-Guided Reasoning, PGR) to reinforcement learning-driven self-play (DeepSeek-R1), reasoning performance continues to be pushed to new heights. However, these advances are almost universally built upon the “scaling hypothesis”: larger model parameter counts, bigger training data, and higher computational budgets. Deploying a model capable of deep mathematical reasoning typically requires billions or even tens of billions of parameters, along with multiple high-end GPUs (e.g., A100, H100) and hundreds of GB of VRAM.

For resource-constrained researchers, independent developers, or edge applications, this paradigm presents an almost insurmountable barrier. This paper confronts an extreme constraint setting: how can one construct a model that is both capable of fluent natural dialogue and deep mathematical reasoning using only a single consumer-grade GPU with 8 GB VRAM (e.g., NVIDIA RTX 3060 Ti), limited system memory (even less than 32 GB), and without large amounts of additional labeled data? The only abundant resource is time.

This constraint forces us to fundamentally reconsider the concept of a “model” in the deep learning era. In the traditional paradigm, a model is a massive floating-point parameter matrix that, once trained, permanently resides in VRAM, with linguistic capabilities (fluency, common sense, instruction following) and reasoning capabilities (logical deduction, mathematical operations) indiscriminately encoded within the same dense vector space. When VRAM is limited, the hard upper bound on parameter scale directly caps the model’s capability ceiling. Existing parameter-efficient methods (such as LoRA, quantization, pruning, distillation) can compress model size, but often at the cost of sacrificing reasoning depth or language quality, and none address the fundamental possibility that “language and reasoning can be decoupled.”

1.2 Core Insights and Design Philosophy

The core insight of this paper is: *linguistic fluency and mathematical reasoning rigor can be heterogeneously decoupled at the parameter level.* Linguistic capability requires coverage of open-domain diversity, is inherently highly redundant, and can be inherited from large-scale pre-trained models and compressed via mature techniques such as quantization, low-rank decomposition, and static pruning, retaining near-complete language quality at minimal VRAM cost. In contrast, mathematical reasoning ability is highly structured, follows specific problem-type logic (skeletons), and the “parameters” it requires can be dynamically generated at inference time by a lightweight hyper-network conditioned on the current problem’s skeleton, then discarded after use.

Furthermore, we re-examine the mathematical essence of “reasoning” itself. We propose that mathematical proof is not a probabilistic “next-token prediction” but a systematic resolution of topological obstructions between premises and conclusions. Through the language of algebraic

topology, the correctness of reasoning can be liberated from the contingency of probabilistic models and transformed into a structurally verifiable guarantee.

Based on these insights, we design a unified framework SHR-DS, whose design philosophy can be summarized as “triple decoupling”:

1. **Parameter decoupling of language and reasoning:** The language base permanently resides in VRAM; reasoning parameters are generated transiently on demand.
2. **Morphological decoupling of knowledge and parameters:** Mathematical problem-solving knowledge is stored in symbolic “solution skeleton” form in a skeleton library, instantiated into neural network parameters by a hyper-network only when needed.
3. **Criterion decoupling of reasoning correctness and statistical prediction:** The legitimacy of reasoning steps is guaranteed by algebraic topological boundary conditions and type system unification; the neural network is only responsible for generating candidate fillers, not for the ultimate correctness of the reasoning.

1.3 Contributions

Specifically, the contributions of this paper include:

- **Theoretical unification:** Proposes and rigorously proves the “Skeleton-Filler Duality Theorem,” establishing a mathematical equivalence between Structural Cognition Training (SCT) and Homological Resolution Reasoning (HRR); introduces a typed solution skeleton quintuple, transforming skeleton transfer from vague semantic matching to a decidable type unification problem.
- **Architectural innovation:** Designs the Dual-Stream Parameter Compression (DPC) architecture, for the first time solidifying linguistic capabilities as a compressed base model via collaborative quantization, low-rank decomposition, and structured pruning, while virtualizing mathematical reasoning capabilities as dynamic expert parameters generated on-demand by a hyper-network conditioned on skeletons, achieving an order-of-magnitude reduction in VRAM footprint.
- **Training paradigm:** Proposes Four-Tier Cognitive Progressive Training, extending Multi-dimensional Self-Elaboration (MSE) into four levels: surface imitation → structural reconstruction → skeleton abstraction → topological cognition, systematically extracting ever more essential mathematical structures from data.
- **Comprehensive estimation and comparison:** Based on publicly available quantitative data from adjacent technologies, provides detailed numerical estimates for SHR-DS across dimensions including data efficiency, reasoning accuracy, language retention, and VRAM footprint, and analyzes essential differences from methods such as MoE, PGR, and MCTS.

1.4 Paper Organization

The remainder of this paper is organized as follows: Section 2 reviews limitations of related work and positions our approach; Section 3 presents the complete theoretical foundation, including the semantic complex, solution skeleton type system, and Skeleton-Filler Duality Theorem; Section 4 elaborates the Dual-Stream Parameter Compression architecture and Four-Tier Cognitive Training paradigm; Section 5 provides theoretical analysis including computational complexity,

convergence, and generalization bounds; Section 6 presents systematic performance estimation and ablation studies; Section 7 discusses applicability boundaries, limitations, and synergy with self-evolution sandboxes; Section 8 concludes and outlines future directions.

2 Related Work and Limitations

2.1 Reasoning Enhancement Methods for Large Language Models

Chain-of-Thought and Autoregressive Generation. Since the introduction of Chain-of-Thought prompting by Wei et al. (2022), generating intermediate reasoning steps to improve final answer accuracy has become a mainstream paradigm. However, these methods still essentially treat reasoning as a sequence generation problem; the model merely learns to imitate surface patterns of reasoning chains in training data without gaining a true understanding of logical structure. This leads to two fatal flaws: first, there is no guarantee of reasoning correctness—a single erroneous intermediate step can cause the entire reasoning chain to collapse; second, data efficiency is extremely low, requiring thousands of examples of the same problem type to internalize a solution pattern.

Inference-Time Search. Tree-of-Thoughts (Yao et al., 2024) and subsequent Monte Carlo Tree Search (MCTS) methods expand the reasoning space at inference time through search, trading extra computation for higher quality solutions. However, search still relies on heuristic scoring functions defined by probabilistic models, lacking structural guarantees of reasoning correctness; simultaneously, the expansion of the search tree in VRAM grows exponentially with width, conflicting with our hardware constraints.

Continuous Relaxation and Energy Models. Potential-Guided Reasoning (PGR) defines reasoning as Langevin dynamics evolution on an energy landscape, offering elegant theoretical formulation. Nevertheless, training the potential function itself still requires massive data, and purely continuous dynamical systems struggle to precisely handle discrete logical hard constraints (e.g., “division by zero”). Moreover, PGR does not solve the storage problem of model parameters—its reasoning capability remains bound to a dense parameter matrix.

Reinforcement Learning and Self-Play. Work exemplified by DeepSeek-R1 enhances reasoning capabilities through RL algorithms such as GRPO. However, such methods assume an already sufficiently powerful base model (billions of parameters), and their training itself requires enormous VRAM and computational resources, making migration to our extreme constraint setting difficult.

2.2 Parameter-Efficient and Model Compression Techniques

Quantization. Post-training quantization methods such as GPTQ can compress FP16 models to 4-bit, significantly reducing VRAM footprint while maintaining reasonable language quality. However, quantization alone cannot improve reasoning capability; the reasoning depth of the compressed model is often comparable to or slightly worse than the original model.

Low-Rank Adaptation and Decomposition. LoRA adds low-rank trainable matrices to a frozen large model for parameter-efficient fine-tuning, but it is essentially an incremental adaptation method and does not reduce the base model’s VRAM footprint. Low-rank decomposition

(e.g., SVD reconstruction) can compress weight matrices, but language quality retention is highly sensitive to the choice of rank, and the compressed model’s performance on structured reasoning tasks may degrade significantly.

Pruning. Structured pruning removes neurons with low contribution, reducing model size. However, pruning is destructive—once pruned, certain fine-grained linguistic capabilities may be permanently lost.

All of the above methods treat linguistic and reasoning capabilities as an inseparable whole for compression, leading to a “see-saw effect”: increased compression inevitably comes with loss of capability in some aspect. The dual-stream decoupling proposed in this paper fundamentally circumvents this dilemma.

2.3 Mixture of Experts (MoE)

MoE increases model capacity without significantly increasing computation by sparsely activating a subset of expert networks. However, all experts in an MoE still physically reside in VRAM; even if only two experts are activated for a given inference, the remaining six still occupy VRAM space. This violates our ultra-low VRAM constraint. Furthermore, MoE’s gating routing is based on statistical features of input tokens rather than the logical structure of the task, thus cannot provide any guarantee of reasoning correctness.

2.4 Neural-Symbolic Methods and Structural Learning

Recent neural-symbolic learning attempts to combine the pattern recognition capabilities of neural networks with the logical deduction capabilities of symbolic systems. However, most work remains at the level of using symbolic knowledge as additional training signals or constraints, without achieving true “symbol-guided parameter generation.” The skeleton type system and hyper-network generator in this paper can be viewed as a deep neural-symbolic fusion: symbolic skeletons tell the model not only “what to do” but also precisely “what parameters are needed to do it.”

2.5 Distinctive Positioning of This Work

Compared to all the above work, the core differences of SHR-DS are:

- For the first time, decouples linguistic and reasoning capabilities at the parameter level heterogeneously, rather than compromising within a single parameter space.
- For the first time, introduces both algebraic topology (homological resolution) and formal type theory into a reasoning framework, structurally guaranteeing reasoning correctness.
- For the first time, proposes the paradigm of on-demand generation of reasoning parameters via a hyper-network, decoupling knowledge capacity from VRAM footprint.

The following sections unfold the complete theory of this framework step by step.

3 Theoretical Foundations

This chapter establishes the mathematical foundations of SHR-DS, sequentially introducing the semantic complex and reasoning obstructions, the type system for solution skeletons, and the core Skeleton-Filler Duality Theorem.

3.1 Semantic Complex and the Topological Model of Reasoning

We first embed mathematical statements and logical relations into an algebraic topological structure.

Definition 3.1 (Semantic Complex) *Let L be the set of all statements in a target mathematical domain, and $A \subset L$ the set of atomic propositions (indivisible logical assertions). A semantic complex K is an abstract simplicial complex satisfying:*

- **Vertex set** $V(K) = A$, each 0-simplex corresponds to an atomic proposition.
- **Edge set** $E(K)$ encodes binary logical relations: an unordered pair $\{u, v\} \in E(K)$ if and only if there exists a direct logical implication, equivalence, or mutual exclusion relation between propositions u and v .
- **Higher-dimensional simplices:** For $k \geq 2$, a k -simplex $\sigma = \{v_0, v_1, \dots, v_k\} \in K$ if and only if the $k + 1$ atomic propositions have an irreducible $(k + 1)$ -ary logical constraint. For example, the modus ponens triple $\{A, A \rightarrow B, B\}$ forms a 2-simplex (triangle), representing the tripartite logical completeness among them.

Construction of the semantic complex is performed automatically by extracting atomic propositions and analyzing logical relations from a large corpus of mathematical texts; the detailed algorithm is given in the Appendix. In the current prototype, the number of vertices is on the order of $10^4 \sim 10^5$, and the total number of simplices is on the order of $10^5 \sim 10^6$, which can be completed in a few minutes on a CPU in a one-time preprocessing step.

Definition 3.2 (Complex Embedding of a Proposition) *Given a mathematical proposition p (which may be a premise, conclusion, or intermediate lemma), it induces a subcomplex X_p in the semantic complex K , spanned by all atomic propositions on which p logically depends and the logical relations among them:*

$$X_p = \{\sigma \in K \mid \text{all vertices } v \in \sigma \text{ are logical dependencies of } p\}.$$

The “logical dependency” relation is determined automatically by matching the syntactic tree of the proposition against the atomic proposition library.

Definition 3.3 (Reasoning Obstruction) *Given a set of premises $P = \{p_1, p_2, \dots, p_n\}$ and a conclusion C , define the reasoning obstruction as the relative homology group:*

$$\mathcal{O}(P, C) = H_*(X_P \cup X_C, X_P),$$

where $X_P = \bigcup_{i=1}^n X_{p_i}$ is the union subcomplex of the premises, and X_C is the subcomplex of the conclusion. Homology coefficients are taken over \mathbb{Z}_2 .

Intuitively, non-trivial relative homology classes correspond to “topological holes” between premises and conclusion—logical gaps that cannot be filled by the known information from the premises. In the simplest example: given premises $\{A, A \rightarrow B\}$ and conclusion B , $X_P \cup X_C$ contains the triangle $\{A, A \rightarrow B, B\}$, which exactly fills a hole in X_P ; therefore the relative homology group is trivial.

Theorem 3.4 (Topological Criterion for Reasoning) *The conclusion C is logically derivable from the premises P only if $\mathcal{O}(P, C) \cong 0$ (all homology groups are trivial). When the semantic complex K is complete for the target mathematical theory (i.e., contains all possible logical constraints), this condition is also sufficient.*

Proof sketch: If a non-trivial relative homology class $[z] \in H_k(X_P \cup X_C, X_P)$ exists, where z is a relative cycle whose boundary lies entirely within X_P but z itself is not in the image of the chain complex of X_P , this means there exists an irreducible k -dimensional obstruction between X_P and X_C that cannot be “filled” by the premises alone. Logically, this corresponds to an ineliminable independence—the conclusion cannot be derived from the premises. Conversely, if all homology classes are trivial, by Whitehead’s theorem for finite complexes, the inclusion $X_P \hookrightarrow X_P \cup X_C$ is a homotopy equivalence, so X_C can be contracted into X_P , which logically implies that C follows from P . Completeness guarantees that the semantic complex captures all logical possibilities of the theory, making topological connectedness equivalent to logical derivability. \square

This theorem reduces the determination of mathematical derivability to a computable algebraic topology problem—homology groups can be computed via Smith normal form of boundary matrices in $O(n^\omega)$ time (where n is the number of simplices, $\omega \approx 2.37$).

3.2 Homological Resolution and Reasoning Fillers

The reasoning process, from a topological perspective, is the progressive addition of “fillers” to eliminate reasoning obstructions.

Definition 3.5 (Reasoning Filler) *Let the current knowledge complex be K_t . A k -reasoning filler Δ_k is a $(k + 1)$ -simplex satisfying:*

- **Boundary condition:** $\partial\Delta_k \subseteq K_t$, i.e., its boundary already exists entirely in the current knowledge.
- **Resolution condition:** Attaching Δ_k to K_t ($K_{t+1} = K_t \cup \{\Delta_k\}$) eliminates at least one non-trivial homology generator.

Theorem 3.6 (Algebraic Effect of Filler Operation) *The above attachment operation reduces the rank of the k -th homology by exactly 1:*

$$\text{rank } H_k(K_{t+1}) = \text{rank } H_k(K_t) - 1.$$

Proof: This follows directly from the long exact sequence of attachment for simplicial complexes; adding a single $(k + 1)$ -simplex affects only the k -th and $(k + 1)$ -th homology, and since $\partial\Delta_k$ is already a boundary in K_t , the operation kills one generator of the k -th homology. \square

Theorem 3.7 (Topological Interpretation of a Reasoning Step) *A standard reasoning step (e.g., “from A and $A \rightarrow B$, infer B ”) is topologically equivalent to: identify a 1-dimensional hole (whose boundary is formed by A , $A \rightarrow B$, and B), then attach the corresponding 2-simplex to fill it.*

Theorem 3.8 (Homological Resolution Sequence of a Reasoning Chain) *A complete and correct reasoning chain is equivalent to a finite obstruction resolution sequence:*

$$\mathcal{O}_0 \xrightarrow{\Delta_1} \mathcal{O}_1 \xrightarrow{\Delta_2} \dots \xrightarrow{\Delta_m} \mathcal{O}_m \cong 0,$$

where $\mathcal{O}_t = H_*(K_t \cup X_C, K_t)$ is the remaining obstruction after step t , and each step adds a reasoning filler Δ_t that eliminates at least one non-trivial generator. The length m is called the topological complexity of the reasoning.

3.3 Type System for Solution Skeletons

The HRR framework above provides a formal model of reasoning, but does not answer the question “how to efficiently find the correct sequence of fillers.” This is the core contribution of Structural Cognition Training (SCT). We extract solution skeletons from very few examples and directly map them to composite fillers.

Definition 3.9 (Mathematical Type System \mathcal{T}) *We introduce a formal type system for mathematical objects and operations, including:*

- **Base types:** \mathbb{R} (real numbers), \mathbb{Z} (integers), \mathbb{B} (booleans), \mathbb{R}^n (n -dimensional real vectors), $\text{Matrix}_{m \times n}$, etc.
- **Composite types:** Product types $T_1 \times T_2$ (e.g., a quadratic equation can be represented as $\text{Quad} = (\mathbb{R}_{\neq 0}, \mathbb{R}, \mathbb{R})$), sum types $T_1 \oplus T_2$ (e.g., the three cases for roots of a quadratic: two real roots, one double root, no real roots), function types $T_1 \rightarrow T_2$ (representing mathematical operations).
- **Type variables:** Placeholders that can be instantiated with concrete types, e.g., α, β .

Definition 3.10 (Typed Solution Skeleton) *A solution skeleton S is a 5-tuple:*

$$S = (\tau_{in}, \tau_{out}, \Sigma, \Gamma, \Delta),$$

where:

- $\tau_{in} \in \mathcal{T}$ is the input type, describing the mathematical structure the problem must satisfy.
- $\tau_{out} \in \mathcal{T}$ is the output type, describing the mathematical structure of the answer.
- $\Sigma = (\sigma_1, \sigma_2, \dots, \sigma_m)$ is a finite sequence of operations, each σ_i satisfying $\sigma_i : \Gamma_{i-1} \rightarrow \Gamma_i$, i.e., a type-preserving operation from type context Γ_{i-1} to Γ_i .
- Γ is the type context, recording all intermediate variables, their types, and dependency relationships, forming a directed acyclic graph.
- Δ is the composite filler corresponding to this skeleton (explicitly constructed by Theorem 5).

Axiom 3.11 (Type Criterion for Skeleton Transfer) *A skeleton S can be transferred from an example e to a new problem q if and only if:*

1. **Type matching:** There exists a type substitution $\theta : \text{TypeVar} \rightarrow \mathcal{T}$ such that $\theta(\tau_{in}) = \text{type}(q)$, and $\text{type}(q)$ is a closed type (contains no free type variables).
2. **Instantiation preserves types:** The operation sequence $\Sigma[q, \theta]$ after substitution, under type context Γ , type-derives a return type of $\theta(\tau_{out})$.
3. **Instantiation correctness:** The actual computed result matches the standard answer.

Conditions 1 and 2 can be mechanically verified by a symbolic type checker in constant time, requiring no neural network involvement. This makes skeleton transfer a decidable, zero-probability-noise process.

3.4 Skeleton-Filler Duality Theorem

We now present the most important theoretical result of this paper, which establishes a strict duality between SCT and HRR.

Theorem 3.12 (Skeleton-Filler Duality Theorem) *Let $S = (\tau_{in}, \tau_{out}, \Sigma, \Gamma, \Delta)$ be a solution skeleton satisfying the upper-closed condition: for each step i , the subcomplex $X_{pre(i)}$ corresponding to the premise type is completely covered by the composite filler $\bigcup_{j=1}^{i-1} \Delta_{\sigma_j}$ from previous steps. Then in the semantic complex K , there exists a uniquely determined composite filler Δ_S (explicitly constructible) satisfying:*

- **Boundary correspondence:** $\partial\Delta_S$ corresponds exactly to the topological gap between the premise subcomplex described by τ_{in} and the conclusion subcomplex described by τ_{out} .
- **Resolution equivalence:** Attaching Δ_S to the current knowledge complex is equivalent to eliminating the obstruction generator g_T uniquely determined by (τ_{in}, τ_{out}) in the relative homology group.
- **Composability:** If a complex problem requires the combination of r skeletons S_1, \dots, S_r , then the composite filler $\Delta_{S_1}, \dots, \Delta_{S_r}$ glued according to matching interface types eliminates all obstructions of the complex problem if and only if the output type of each skeleton unifies with the input type of the next skeleton.

Proof (constructive):

Step 1: Mapping primitive operations to basic fillers. Consider an operation σ_i in the skeleton, which transforms objects in premise type context Γ_{i-1} into new objects in Γ_i . In the semantic complex, the atomic propositions involved in Γ_{i-1} span the subcomplex $X_{pre(i)}$, and those involved in Γ_i span $X_{post(i)}$. Define the basic filler corresponding to σ_i as:

$$\Delta_{\sigma_i} = \text{ConvexHull}(X_{pre(i)} \cup X_{post(i)}) \setminus X_{pre(i)},$$

i.e., the smallest convex subcomplex containing all these atomic propositions, minus the part already covered by the premises. Since $\partial\Delta_{\sigma_i} \subseteq X_{pre(i)} \cup X_{post(i)}$ and $X_{pre(i)}$ is already contained in current knowledge by Γ , Δ_{σ_i} satisfies the boundary condition.

Step 2: Upper-closed condition guarantees gluing legality. The upper-closed condition ensures that at each step of the skeleton sequence, the premises of the current operation have already been constructed in previous steps. Therefore, when we construct the composite filler $\Delta_S =$

$\bigcup_{i=1}^m \Delta_{\sigma_i}$ in the order of Σ , the boundary of each attachment operation already exists in the expanded complex, introducing no new obstructions. By the Mayer–Vietoris sequence and the upper-closed condition, one can prove that $\partial\Delta_S$ consists exactly of the difference between the “first premise subcomplex” and the “final conclusion subcomplex,” which equals the topological gap determined by $(\tau_{\text{in}}, \tau_{\text{out}})$.

Step 3: Resolution equivalence. By Theorem 4, a generator of reasoning obstruction is represented by a representative cycle. By construction, this cycle is exactly bounded by $\partial\Delta_S$. Hence, attaching Δ_S is equivalent to adding that boundary in the chain complex, thereby killing this generator.

Step 4: Composability. When multiple skeletons are combined, if the output type of skeleton j unifies with the input type of skeleton $j+1$, then the conclusion subcomplex of the former exactly covers the premise subcomplex of the latter, and the boundaries at the glue point automatically match. The boundary of the entire glued composite filler remains only the input of the first skeleton and the output of the last skeleton, exactly corresponding to the reasoning obstruction of the overall complex problem. \square

Corollary 3.13 (Topological Necessary and Sufficient Condition for Skeleton Transfer)

Skeleton S can be transferred from example e to new problem q if and only if:

1. *The type matching condition (Condition 1 of Axiom 1) holds.*
2. *Attaching the composite filler Δ_S corresponding to S to the initial knowledge complex of q completely resolves its reasoning obstruction $\mathcal{O}(P_q, C_q)$.*

Both conditions are mechanically verifiable: Condition 1 by a type checker, Condition 2 by computing the homology after one attachment. Skeleton transfer thus becomes a precisely verifiable operation.

3.5 Persistent Homology and Reasoning Quality

Persistent homology provides a natural quantitative evaluation tool for the reasoning process.

Definition 3.14 (Reasoning Barcode) *For a reasoning process $\{\Delta_1, \dots, \Delta_m\}$, the reasoning barcode is the “birth-death” diagram of all homology generators:*

- **Birth time** b_i : *the reasoning step at which generator g_i first appears.*
- **Death time** d_i : *the step at which g_i is killed by a filler.*
- **Persistence** $\rho_i = d_i - b_i$.

Theorem 3.15 (Barcode Criterion for Reasoning Quality) *An “elegant” reasoning satisfies:*

- **Minimal total persistence:** $\sum_i \rho_i$ *is minimized, equivalent to having no redundant steps.*
- **Monotonic death:** $\forall i < j, d_i \leq d_j$, *no logical backtracking.*
- **Complete resolution:** *No generator with infinite persistence remains in the final barcode; the reasoning is complete.*

In SHR-DS, the persistent homology barcode of a skeleton also records its “birth” (first successful extraction) and “death” (replaced by a more general skeleton or deemed erroneous) during training, used to evaluate the long-term effectiveness of the skeleton library.

4 Dual-Stream Parameter Compression Architecture

Section 3 established the topological model of reasoning and the formal type system for skeletons, and proved the duality between skeletons and fillers. However, these theoretical achievements have not yet solved a core engineering bottleneck: how to fit a full 1B-scale language model along with its mathematical reasoning capability into 8 GB of VRAM? This chapter proposes the Dual-Stream Parameter Compression (DPC) architecture, which fundamentally decouples the parameter requirements for linguistic fluency and mathematical reasoning.

4.1 Core Idea: Heterogeneous Storage of Language Stream and Reasoning Stream

Traditional dense models indiscriminately encode all capabilities into a single huge parameter matrix. This matrix must fully reside in GPU VRAM at all times, regardless of whether the current task is casual conversation, translation, poetry writing, or calculus derivation. Our key observation is that linguistic capability and reasoning capability have fundamentally different requirements for parameter storage density, redundancy, and retrieval patterns.

- **Linguistic capability (Language Stream)** needs to cover open-domain vocabulary collocations, syntactic templates, and common-sense associations; it is inherently highly redundant—many neurons and connections contribute only marginally to final generation quality. Therefore, the language stream can be aggressively compressed via quantization, low-rank decomposition, and structured pruning, and once compressed, can permanently reside in VRAM for use by all tasks.
- **Reasoning capability (Reasoning Stream)** is completely different. It is highly dependent on the specific mathematical structure of the current problem: solving a quadratic equation, computing a definite integral, and proving a geometric theorem each require entirely distinct “expert knowledge.” But these expert knowledge modules are not needed at all times—when solving a quadratic, we have no need for calculus expert parameters. Therefore, the reasoning stream naturally fits a “generated on demand, discarded after use” transient existence.

We use a lightweight hyper-network generator G_ϕ (approximately 10M parameters) to replace the storage of full reasoning parameters. Given the symbolic description of a solution skeleton, G_ϕ can generate all the specialized parameters (MLP weights, attention projections, etc.) required for that skeleton within milliseconds. After reasoning completes, these parameters are immediately released, leaving no VRAM footprint. What is truly permanently stored are only the generator itself and a symbolic skeleton library (less than 10 MB).

4.2 Language Stream: Compressed Language Base

The goal of the language stream is to extract, from a pre-trained 1B dense language model (e.g., Qwen2.5-1.5B), a language base B_{lang} that is as small as possible while retaining as much language quality as possible. We adopt a three-stage collaborative compression pipeline: “quantization \rightarrow low-rank decomposition \rightarrow structured pruning.”

4.2.1 Pre-trained Parent Model Selection

We select a 1B-scale model that has been well pre-trained on both general corpora and mathematical corpora as the parent model. This model should possess fluent Chinese generation capability, basic English capability, and support for mathematical symbols and L^AT_EX. In experiments, we recommend Qwen2.5-1.5B or a general variant of DeepSeek-Math-1.3B. Denote its full-precision (FP16) parameter matrix as $W_{\text{orig}} \in \mathbb{R}^{d \times d}$ (taking an FFN layer as an example; actual layer sizes may vary; for simplicity we use a unified notation here).

4.2.2 Step 1: GPTQ 4-bit Quantization

Quantization maps continuous floating-point weights to discrete integer grids, significantly reducing storage bits per parameter. We employ GPTQ (Frantar et al., 2023) for column-wise post-training quantization. For a column w_j of weight matrix W , given calibration inputs X (from general corpora, about 128 samples), GPTQ greedily quantizes columns one by one and propagates the error to remaining columns:

1. Perform quantization operation Q on column w_j :

$$\hat{w}_j = Q(w_j), \quad Q : \mathbb{R} \rightarrow \mathbb{Z}_{2^4}$$

2. Compute quantization error $\epsilon_j = w_j - \hat{w}_j$.
3. Propagate error to remaining unquantized columns:

$$w_{>j} \leftarrow w_{>j} - \frac{\epsilon_j}{[H^{-1}]_{jj}} \cdot (H^{-1})_{:,j}$$

where $H = 2XX^T + \lambda I$ is the approximate Hessian matrix, and λ is a damping coefficient.

The objective is to minimize the output error:

$$\arg \min_{\hat{W}} \|W_{\text{orig}}X - \hat{W}X\|_F^2$$

After GPTQ-4bit quantization, model volume is reduced by approximately 75% (16 bit \rightarrow 4 bit). Language quality loss is typically $< 2\%$ (in WikiText-2 perplexity). Denote the quantized model parameters as \tilde{W}_q .

4.2.3 Step 2: Low-Rank Decomposition

Quantized weight matrices still have high redundancy, especially in the intermediate dimensions of FFN layers. We further apply low-rank decomposition to each weight matrix:

$$\tilde{W}_q \approx AB, \quad A \in \mathbb{R}^{d \times r}, B \in \mathbb{R}^{r \times d}$$

where $r \ll d$ is the chosen rank. The decomposition is performed via truncated SVD:

1. Compute the SVD of \tilde{W}_q : $\tilde{W}_q = U\Sigma V^T$.

2. Retain the first r largest singular values and corresponding singular vectors:

$$A = U_{:,r} \Sigma_{r,r}, \quad B = \Sigma_{r,r} V_{:,r}^T$$

The choice of rank r follows the constraint “language quality degradation below 2%.” In our experiments, $r \approx d/8$ achieves a good balance between compression (approximately $3\times$ additional compression) and quality loss. Matrices A and B are both stored in 4-bit quantization.

4.2.4 Step 3: Structured Pruning

After quantization and low-rank decomposition, many neurons still contribute little to language fluency. We apply structured pruning to remove these neurons, further reducing the total parameter count. Define the importance score of neuron i as:

$$s_i = \|\tilde{W}_{:,i}\|_2 \cdot \mathbb{E}_{x \sim \mathcal{D}_{\text{lang}}} [\|h_i(x)\|_2]$$

where $\tilde{W}_{:,i}$ is the weight column vector connecting to neuron i , and $h_i(x)$ is the activation value of that neuron on input x . This score considers both weight magnitude and actual activation strength, preferentially removing “small-weight and rarely-activated” neurons.

Sort by s_i ascending, and remove the least important p proportion of neurons ($p \approx 0.3$). After pruning, weight matrices of adjacent layers are cropped accordingly, and we fine-tune for 100 steps to recover possible accuracy loss. Final language quality loss is $< 3\%$.

4.2.5 Final Configuration of the Language Base

The stacking effect of the three compression stages ($75\% \times \sim 60\% \times 70\%$) reduces the parameter count to approximately 10% of the original 1B model, i.e., about 100M parameters (FP16) or about 500 MB (4-bit). This language base B_{lang} will permanently reside in GPU VRAM, responsible for text understanding, fluent generation, instruction following, and natural language explanation of mathematical solution steps for all tasks. As we will see in the next subsection, this 500 MB permanent overhead is entirely acceptable.

4.3 Reasoning Stream: Hyper-Network Skeleton Factory

The goal of the reasoning stream is to instantly obtain specialized neural network parameters corresponding to a required solution skeleton when a problem arises, use them, and then release them. We achieve this through a lightweight hyper-network generator G_ϕ and a symbolic skeleton library.

4.3.1 Skeleton Symbol Library

The skeleton library stores the 5-tuple definitions of all known solution skeletons: $S = (\tau_{\text{in}}, \tau_{\text{out}}, \Sigma, \Gamma, \Delta)$. Stored in textual or directed acyclic graph format, each skeleton occupies about 1–5 KB. A library of 1000 skeletons totals less than 10 MB, residing in CPU memory or disk, loaded on demand. Sources of the skeleton library include:

- Automatic extraction from a small number of examples (SCT Stage 2).
- Manual design by human experts (for core mathematical operations such as quadratic formula, integration by parts, etc.).
- Automatic discovery and verification via a self-evolution sandbox (Pillar 4).

4.3.2 Encoder: Skeleton Symbol to Condition Vector

The encoder E_ψ transforms the symbolic description of a skeleton into a fixed-dimensional condition vector $c \in \mathbb{R}^{d_c}$ for use by the generator. E_ψ is a small Transformer encoder (approximately 1M parameters), whose input is a textual or graphical representation of the skeleton:

- **Textual representation:** The skeleton step sequence Σ is expanded into natural language descriptions (e.g., “Step 1: compute discriminant $\Delta = b^2 - 4ac \dots$ ”).
- **Graphical representation:** The variable dependency graph in the type context Γ is encoded as an adjacency matrix of a DAG, plus type embeddings for each node.

The encoder outputs the mean pooling of the sequence as the condition vector:

$$c = E_\psi(S).$$

4.3.3 Hyper-Network Generator

The hyper-network generator G_ϕ (approximately 10M parameters) receives the condition vector c and a random noise seed $z \sim \mathcal{N}(0, I)$, and generates all reasoning parameters W_{expert} required for skeleton S . The scale of W_{expert} is determined by the difficulty of the current problem and the complexity of the skeleton, typically ranging from 12–128M parameters (corresponding to the upper bound of activated parameters).

Generation architecture: G_ϕ consists of multiple conditional batch normalization deconvolution layers that progressively map the low-dimensional condition vector and noise to a high-dimensional parameter space. Suppose we need to generate a weight matrix of shape $d_{\text{in}} \times d_{\text{out}}$. The generation process is as follows:

1. Concatenate condition vector c and noise z , map via a linear layer to an initial feature map $H_0 \in \mathbb{R}^{h_0 \times w_0}$.
2. Pass through L conditional deconvolution layers, each:

$$H_\ell = \text{ReLU}(\text{CondBN}(\text{Deconv}(H_{\ell-1}), c))$$

where CondBN is conditional batch normalization, whose scale and shift parameters are dynamically generated from c by a small MLP:

$$\text{CondBN}(x, c) = \gamma(c) \cdot \frac{x - \mu}{\sigma} + \beta(c)$$

3. The final layer output is reshaped to target shape $d_{\text{in}} \times d_{\text{out}}$ and passed through an appropriate activation function (e.g., tanh) to keep generated weights within a reasonable numerical range.

Role of noise: The noise seed z allows the generation of functionally equivalent but numerically slightly different parameter instances for the same skeleton. This is constrained during

training by a skeleton consistency loss (see Section 5), ensuring that the generated parameters are functionally consistent rather than numerically identical.

Content addressability: After sufficient training, G_ϕ becomes content-addressable: given the symbolic description of the same skeleton, it always generates parameters that can correctly execute the reasoning steps of that skeleton (in a probabilistic sense). This is jointly guaranteed by the skeleton consistency loss and contrastive learning.

4.3.4 Caching Strategy

For frequently used skeletons (e.g., the top 100 most common skeletons), their generated parameters can be cached to disk or CPU memory. When the same skeleton is encountered again, parameters are loaded from cache instead of regenerating. Cache management uses an LRU (Least Recently Used) policy, with total cache size controlled within 1–2 GB on disk (assuming 1000 skeletons \times average 50M parameters \times 4-bit \approx 2.5 GB). Since reasoning stream parameters do not permanently occupy GPU VRAM, this disk overhead is entirely acceptable and naturally diminishes as hardware improves (larger SSDs).

4.4 Runtime Assembly and SHR Integration

The language base B_{lang} and the reasoning stream generator G_ϕ constitute the static components of the DPC architecture. We now describe their dynamic collaborative workflow when handling a concrete mathematical problem.

Step 1: Problem parsing and type inference. The user inputs a mathematical problem q (in natural language). The language base B_{lang} first parses q into an internal representation, and through an additional lightweight type inference head (approximately 0.5M parameters, attached above the output layer of the language base) predicts the mathematical type $\text{type}(q) \in \mathcal{T}$. The type inference head is trained in a supervised manner on (problem text, type annotation) pairs, which can be generated in bulk by synthetic data.

Step 2: Skeleton matching. Search the skeleton library for skeletons whose input type τ_{in} unifies with $\text{type}(q)$. Retrieval uses a standard type unification algorithm (e.g., a variant of Hindley-Milner), executed entirely by a symbolic engine without neural network approximation:

- Initialize type substitution $\theta = \emptyset$.
- For each skeleton S in the library, attempt to unify $\tau_{\text{in}}(S)$ with $\text{type}(q)$:
 - If both are the same base type, succeed.
 - If one is a type variable α , add $\theta(\alpha) = \text{the other}$.
 - If both are composite types $T_1 \times T_2$, recursively unify corresponding components.
 - Otherwise fail.
- If successful, the skeleton matches; otherwise proceed to the next skeleton.

If only partial matching is needed (complex problems may require multiple skeletons), record all successfully matched skeletons along with their type substitutions.

Step 3: Parameter generation or loading. For each matched skeleton S , check whether its generated parameters are already cached. If cached, load W_{expert} from disk or CPU memory

to GPU. Otherwise, generate via $G_\phi(E_\psi(S), z)$. Generation time is approximately 5–20 ms (depending on the required parameter scale; a measured estimate is provided later). Parameters for multiple skeletons can be generated in parallel.

Step 4: Dynamic network assembly. The language base B_{lang} and one or more currently loaded/generated reasoning experts $W_{\text{expert}}^{(1)}, W_{\text{expert}}^{(2)}, \dots$ need to be fused into a unified reasoning network. We adopt a *bypass injection* approach that avoids major modifications to the language base structure.

Specifically, for each Transformer layer of the language base, add a “reasoning bypass” alongside the FFN sublayer. Let the output of the ℓ -th FFN layer of the language base be $h_{\text{lang}}^{(\ell)}$. The reasoning bypass consists of the corresponding layer of the currently activated reasoning expert, with output $h_{\text{expert}}^{(\ell)}$. The two are fused via a learnable gating mechanism:

$$h_{\text{out}}^{(\ell)} = \alpha^{(\ell)} \odot h_{\text{lang}}^{(\ell)} + (1 - \alpha^{(\ell)}) \odot h_{\text{expert}}^{(\ell)}$$

where $\alpha^{(\ell)} = \text{sigmoid}(W_{\text{gate}}^{(\ell)} [h_{\text{lang}}^{(\ell)}; h_{\text{expert}}^{(\ell)}])$ is a gating vector, learned by a small linear transformation. The gating layer is trained during Stage 2 (structural reconstruction) to learn, based on the context of the current token, whether to rely more on the language stream or the reasoning stream.

When a problem requires the combination of multiple skeletons, the expert modules corresponding to each skeleton are connected into the bypass in the order specified by skeleton interface types, with intermediate variable transfers guaranteed compatible by the type context Γ .

Step 5: Homology-guided reasoning. The assembled network executes the HRR obstruction resolution algorithm (Section 3.2). The language base is responsible for translating each reasoning step into natural language explanation (“compute discriminant $\Delta = b^2 - 4ac = 25 - 24 = 1 > 0$, therefore there are two real roots”), while the reasoning expert performs the actual symbolic computation (either implicitly in the continuous representation of the network or by calling an external mathematical engine).

During reasoning, the persistent homology barcode tracks the resolution status of each obstruction in real time. If the router (Pillar 1) detects that the remaining obstruction requires additional skeleton support, it can interrupt the current process and re-enter skeleton matching and parameter generation (Steps 2–4), achieving dynamic deep routing.

Step 6: Release. After reasoning completes, output the final answer and the full reasoning chain. All loaded or generated reasoning stream parameters are immediately released from GPU VRAM, leaving only the language base, hyper-network generator G_ϕ , and encoder E_ψ . GPU VRAM returns to the baseline state containing only the language base and generator (approximately $500 + 20 = 520$ MB), ready for the next problem.

4.5 VRAM and Time Budget

The following gives precise VRAM footprint estimates at inference time (based on 256M total parameters + 128M activated parameters, FP16/4-bit mixed precision):

Table 1: Inference VRAM footprint breakdown

Component	Size	Location
Language base B_{lang} (4-bit, after low-rank + pruning)	~500 MB	GPU VRAM
Hyper-network generator G_ϕ (FP16)	~20 MB	GPU VRAM
Encoder E_ψ (FP16)	~2 MB	GPU VRAM
Currently activated reasoning expert (FP16, max 128M params)	~256 MB	GPU VRAM
KV cache (16384 context, 24 layers, hidden=1024)	~150 MB	GPU VRAM
Skeleton library (1000 skeletons, text/graph format)	<10 MB	CPU memory/disk
Skeleton parameter cache (optional, 1000 \times 50M \times 4-bit)	~2.5 GB	Disk
Total GPU VRAM footprint	< 1 GB	
Additional during training (optimizer states ZeRO-Offload to CPU)	~8 GB	CPU memory

In contrast to an uncompressed 1B dense model (FP16 inference requires approximately 6–8 GB), SHR-DS reduces inference VRAM by about an order of magnitude without sacrificing mathematical reasoning depth.

Regarding time overhead, reasoning stream parameter generation is the only additional latency source. Generating 128M parameters once takes about 5–20 ms (measured estimate on modern consumer GPUs with similarly scaled hyper-networks). For a single mathematical problem (typical reasoning steps 10–50, total time a few seconds), this overhead is negligible. Even in extreme long-chain reasoning (hundreds of steps, requiring multiple skeleton parameter generations), cumulative generation latency is still on the order of hundreds of milliseconds, far less than human thinking and typing time. During training, each batch incurs an additional generator forward/backward pass; total training time increases to 1.5–2.5 times that of standard training, but the improvement in data efficiency (10–50 \times) means that total training time may actually be shorter (see Section 6 performance estimates).

5 Training Paradigm: Four-Tier Cognitive Progression

The DPC architecture in Section 4 solves the problem of “how to fit the model into VRAM,” but does not answer “how the model learns to reason.” This chapter proposes the Four-Tier Cognitive Progressive Training, extending Multidimensional Self-Elaboration (MSE) into four increasingly essential cognitive levels: surface imitation \rightarrow structural reconstruction \rightarrow skeleton abstraction \rightarrow topological cognition, systematically extracting ever more essential mathematical structures from data.

5.1 Training Components and Data Flow

Training of SHR-DS involves three learnable components:

- **Language base** B_{lang} : updated in Stages 1 and 2; from Stage 3 onward, most parameters are frozen, only fine-tuning the gating layers and type inference head.
- **Hyper-network generator** G_ϕ **and encoder** E_ψ : introduced in Stage 3, fully trained in Stage 4.
- **Gating fusion layers** W_{gate} : introduced in Stage 2, fine-tuned throughout all stages.

Training data includes:

- **General corpus** $\mathcal{D}_{\text{general}}$: from Wikipedia, books, high-quality web texts, maintaining language fluency.
- **Mathematical corpus** $\mathcal{D}_{\text{math}}$: competition problems, textbook problems, solution processes, accounting for about 60%, distributed in a difficulty pyramid (basic 50%, intermediate 30%, hard 15%, competition 5%).
- **Skeleton-annotated data** $\mathcal{D}_{\text{skel}}$: a small number (approximately 300–500 problems) of (problem, reasoning chain, solution skeleton) triples, manually or automatically annotated, used only for cold-start of skeleton extraction.

5.2 Stage 1: Surface Imitation

Goal: Enable the language base B_{lang} to acquire basic linguistic fluency and the most superficial problem-solving patterns.

Training method: Standard autoregressive language modeling on $\mathcal{D}_{\text{general}} \cup \mathcal{D}_{\text{math}}$. The mathematical portion uses the complete format “problem \rightarrow step-by-step reasoning \rightarrow answer.” This stage does not involve skeletons, generators, or any reasoning stream components; it trains only the language base itself.

Loss function:

$$\mathcal{L}_1 = -\frac{1}{|\mathcal{D}|} \sum_{(x,y) \in \mathcal{D}} \log P_{B_{\text{lang}}}(y | x)$$

where x is the prefix sequence and y is the target next token.

Training details: Under the 100M configuration, Stage 1 takes about 1–2 days (3060 Ti, DeepSpeed ZeRO-2 + CPU Offload). Upon completion, the language base should be able to generate fluent Chinese and output “plausible-looking” mathematical solution steps, but with low accuracy (simple problems about 40–50%, complex problems $< 20\%$).

5.3 Stage 2: Structural Reconstruction

Goal: Enable the model to “reverse-engineer the process from the answer,” thereby understanding causal dependencies between reasoning steps.

Training method: Freeze most parameters of the language base (except gating fusion layers and the type inference head). Each training sample includes problem q , final answer a , and

the complete reasoning chain $r = (r_1, r_2, \dots, r_T)$. The input is “problem + answer + reverse prompt,” and the model must autoregressively generate the reasoning chain r .

This stage introduces the reasoning stream for the first time. Since G_ϕ is not yet trained, the reasoning expert parameters use fixed small-scale parameters generated by a randomly initialized G_ϕ (e.g., 12M), or directly use a small, skeleton-unconditioned shared expert.

Loss function:

$$\mathcal{L}_2 = -\frac{1}{|\mathcal{D}_{\text{math}}|} \sum_{(q,r,a)} \sum_{t=1}^T \log P_{\text{assembled}}(r_t | q, a, r_{<t})$$

where $P_{\text{assembled}}$ is the joint output of the language base plus current reasoning expert after gated fusion.

Auxiliary task: Simultaneously train the type inference head to predict the correct $\text{type}(q)$ from problem q . Use cross-entropy loss:

$$\mathcal{L}_{\text{type}} = -\frac{1}{|\mathcal{D}_{\text{math}}|} \sum_{(q,\tau)} \log P_{\text{type_head}}(\tau | q)$$

Type annotations τ can be automatically generated by existing large models or manually annotated.

Training details: Approximately 1–2 days. Upon completion, the model can generate reasonably correct reasoning chains (given that the answer is known). This corresponds to learning “verification” rather than “exploration.”

5.4 Stage 3: Skeleton Extraction and Generator Meta-Training

Goal: Strip out transferable solution skeletons from reasoning chains and train the hyper-network generator G_ϕ to generate correct reasoning expert parameters conditioned on skeletons.

Training method: Performed in two alternating sub-stages.

5.4.1 Sub-stage 3A: Skeleton Extraction

Use a small amount of skeleton-annotated data $\mathcal{D}_{\text{skel}}$ (300–500 samples) as cold-start to train the skeleton encoder E_ψ to extract the skeleton S from (problem, reasoning chain) pairs.

Skeleton extraction loss:

$$\mathcal{L}_{\text{skel}} = -\frac{1}{|\mathcal{D}_{\text{skel}}|} \sum_{(q,r,S^*)} \log P_{E_\psi}(S^* | q, r)$$

where S^* is the annotated ground-truth skeleton. The output of E_ψ is an autoregressive generation of the components of the skeleton quintuple (input type, output type, step sequence, type context).

After cold-start, E_ψ can be used to automatically generate skeleton pseudo-labels for larger-scale unannotated mathematical data, expanding the training set.

5.4.2 Sub-stage 3B: Hyper-Network Meta-Training

For each skeleton S and its multiple example instances $\{q_1, q_2, \dots, q_k\}$ ($k \geq 2$), train G_ϕ to generate reasoning parameters that correctly solve these examples.

Procedure:

1. Encode skeleton S with E_ψ to obtain condition vector $c = E_\psi(S)$.
2. Sample noise $z \sim \mathcal{N}(0, I)$, generate parameters $W = G_\phi(c, z)$.
3. Assemble the reasoning network using the generated W , perform forward reasoning on example q_i to obtain predicted answer \hat{a}_i .

Task loss: Standard cross-entropy ensuring the generated parameters correctly solve the problems.

$$\mathcal{L}_{\text{task}} = -\frac{1}{k} \sum_{i=1}^k \log P_W(\hat{a}_i = a_i \mid q_i)$$

Skeleton consistency loss: Enforces that parameters generated for the same skeleton from different examples are functionally consistent. For two different noise seeds z_1, z_2 generating parameters W_1, W_2 for the same skeleton S , their activation patterns on the same set of input tokens should be highly similar. We impose a cosine similarity loss on intermediate representations:

$$\mathcal{L}_{\text{consist}} = \frac{1}{L} \sum_{\ell=1}^L (1 - \cos(\text{Rep}_\ell(W_1), \text{Rep}_\ell(W_2)))$$

where $\text{Rep}_\ell(W)$ is the output representation of the network with parameters W at layer ℓ on a fixed set of random inputs. This loss forces the generator to capture the essence of the skeleton rather than non-essential variations due to random noise.

Total hyper-network loss:

$$\mathcal{L}_{\text{hyper}} = \mathcal{L}_{\text{task}} + \lambda_{\text{consist}} \mathcal{L}_{\text{consist}}$$

λ_{consist} is recommended to be 0.1–0.3.

Training details: Stage 3 takes about 2–4 days. Upon completion, G_ϕ should be able to generate high-quality reasoning parameters for skeletons seen during training, and have some generalization capability to unseen skeletons (as long as their type combinations are within the training distribution). At this point, model accuracy on simple problem types should rise to 70–80%, and on intermediate types 50–60%.

5.5 Stage 4: Topological Cognition – HRR Integration

Goal: Endow the model with topological understanding of skeletons—skeletons are not just solution steps but composite fillers that resolve homological obstructions. Train the model to understand the topological meaning of skeletons and use persistent homology to monitor reasoning quality at inference time.

Training method: Integrate the HRR framework from Section 3 with neural network training.

5.5.1 Filler Generation Training

For each training sample, decompose its reasoning chain into a filler sequence $\Delta_1, \dots, \Delta_m$. The training objective is to enable the model (language base + current reasoning expert) to correctly predict the next filler Δ_t given the current knowledge complex K_{t-1} and remaining obstruction \mathcal{O}_{t-1} .

HRR filler loss:

$$\mathcal{L}_{\text{HRR}} = \frac{1}{m} \sum_{t=1}^m [-\log P(\Delta_t^* | K_{t-1}, \mathcal{O}_{t-1}) + \lambda_{\text{boundary}} \cdot \mathcal{L}_{\text{boundary}}(\Delta_t^*, K_{t-1})]$$

where Δ_t^* is the ground-truth correct filler (derived from the skeleton by Theorem 5), and the boundary loss ensures the predicted filler satisfies the boundary condition:

$$\mathcal{L}_{\text{boundary}}(\Delta, K) = \|\partial\Delta - \text{proj}_K(\partial\Delta)\|^2$$

Here $\partial\Delta$ and $\text{proj}_K(\partial\Delta)$ are represented in some continuous embedding (e.g., average pooling of atomic proposition embeddings in the semantic complex). In practice, this loss can be implemented as a regularizer that forces the filler’s boundary representation to be consistent with the already-existing part of the current knowledge complex.

5.5.2 Skeleton-Filler Alignment Loss

To further strengthen the topological interpretation of skeletons, we introduce an alignment loss that requires the composite filler Δ_S directly mapped from skeleton S (Theorem 5) to be consistent with the filler sequence actually generated and used by the neural network.

$$\mathcal{L}_{\text{align}} = \sum_{S \in \mathcal{L}_{\text{skel}}} \|\text{Embed}(\Delta_S) - \text{Embed}(\Delta_{\text{neural},S})\|^2$$

where $\text{Embed}(\cdot)$ encodes a filler (a set of simplices) into a fixed-dimensional vector (via a graph neural network or simple pooling), and $\Delta_{\text{neural},S}$ is the filler actually constructed by the network when reasoning with skeleton S .

5.5.3 Persistent Homology Regularization

During training, periodically (every 100 steps) compute the reasoning barcode of the current model (on a validation set) and add the total persistence $\sum_i \rho_i$ as a regularization term to the loss, encouraging the model to learn concise reasoning paths:

$$\mathcal{L}_{\text{barcode}} = \frac{1}{|\mathcal{D}_{\text{val}}|} \sum_{q \in \mathcal{D}_{\text{val}}} \sum_i \rho_i(q)$$

Training details: Stage 4 takes about 1–3 days. Upon completion, the model’s performance on complex reasoning problems should show a qualitative leap (expected accuracy improvement of 10–20 percentage points), and reasoning chains will be significantly more concise with fewer logical jumps.

5.6 Joint Training Objective and Curriculum

The total training loss of SHR-DS is a weighted sum of the four stages, but in practice, it is introduced via a curriculum schedule:

- **Early (first 30% of training):** Only \mathcal{L}_1 (Stage 1), weight $\alpha_1 = 1$, others 0.
- **Early-mid (30%–50% of training):** Introduce $\mathcal{L}_2 + \mathcal{L}_{\text{type}}$, $\alpha_1 = 0.3$, $\alpha_2 = 0.7$.
- **Mid-late (50%–80% of training):** Introduce $\mathcal{L}_{\text{skel}} + \mathcal{L}_{\text{hyper}}$, $\alpha_1 = 0.1$, $\alpha_2 = 0.3$, $\alpha_3 = 0.6$.
- **Late (80%–100% of training):** Introduce $\mathcal{L}_{\text{HRR}} + \mathcal{L}_{\text{align}} + \mathcal{L}_{\text{barcode}}$, $\alpha_1 = 0.05$, $\alpha_2 = 0.15$, $\alpha_3 = 0.3$, $\alpha_4 = 0.5$.

Joint loss:

$$\mathcal{L}_{\text{SHR-DS}} = \alpha_1 \mathcal{L}_1 + \alpha_2 (\mathcal{L}_2 + \mathcal{L}_{\text{type}}) + \alpha_3 (\mathcal{L}_{\text{skel}} + \mathcal{L}_{\text{hyper}}) + \alpha_4 (\mathcal{L}_{\text{HRR}} + \mathcal{L}_{\text{align}} + \mathcal{L}_{\text{barcode}})$$

This curriculum simulates the cognitive development process of humans learning mathematics: first learn to speak and imitate (Stage 1), then learn to reverse-engineer the process from the answer (Stage 2), then abstract general solution methods from concrete problems (Stage 3), and finally understand the “why” of the methods—i.e., their topological essence (Stage 4).

6 Theoretical Analysis

This chapter provides rigorous theoretical analysis of the computational feasibility, training convergence, and generalization capability of the SHR-DS framework. These analyses not only prove the favorable mathematical properties of the framework but also provide well-founded hardware configuration and training time estimates for engineering implementation.

6.1 Computational Complexity Analysis

The runtime overhead of SHR-DS mainly comes from three aspects: (1) homology group computation on the semantic complex; (2) forward propagation of the hyper-network generator to produce reasoning parameters; (3) type unification during skeleton matching. We give complexity estimates for each.

Semantic complex construction and homology computation. The number of vertices $|A|$ in the semantic complex K is typically on the order of 10^4 – 10^5 , and the total number of simplices $n = |K|$ is on the order of 10^5 – 10^6 . Construction is a one-time offline operation: extracting atomic propositions and logical relations from the corpus, complexity $O(|A|^3)$, completed in a few minutes on a single CPU.

Computation of the relative homology group $H_*(X_P \cup X_C, X_P)$ requires building a boundary matrix and computing its Smith normal form. Let the number of simplices involved locally be n_{local} (typically much smaller than n , as only the neighborhoods of premises and conclusion are involved). The complexity of Smith normal form of a boundary matrix is $O(n_{\text{local}}^\omega)$, where $\omega \approx 2.37$ is the exponent of matrix multiplication. In practice, using specialized persistent homology libraries (e.g., PHAT or Ripser), homology computation for $n_{\text{local}} = 10^4$ takes a few milliseconds. For $n_{\text{local}} = 10^5$, a single computation takes about tens of milliseconds. Since

homology computation is called only during skeleton validation and part of the loss computation (frequency much lower than per-batch), these overheads impact overall training time by less than 5%.

Hyper-network generation overhead. The hyper-network G_ϕ (approximately 10M parameters) generating up to 128M parameters for a reasoning expert involves a forward pass of moderate scale. The generator typically consists of a few conditional deconvolution or fully connected layers, with forward inference time about 5–20 ms (on RTX 3060 Ti, assuming 80% GPU utilization). Compared to typical Transformer inference (a 100M model processing 4096 tokens takes about 50–100 ms), the generation overhead is acceptable. With caching, the generation overhead for frequent skeletons can be fully amortized.

Type unification. Type unification in skeleton matching is based on a Hindley-Milner style unification algorithm, with complexity linear in the size of type expressions. Since the input type τ_{in} of each skeleton usually contains only a few to a dozen type constructors, and the skeleton library total size is less than 1000, a single matching operation takes far less than 1 ms, with no substantial impact on inference speed.

Total inference latency. A complete mathematical reasoning session (including 20–50 resolution steps, several skeleton generations) on an RTX 3060 Ti is expected to have end-to-end latency of 1–5 seconds, fully meeting interactive application requirements. For extremely long-chain reasoning (hundreds of steps), latency grows linearly to tens of seconds, comparable to human time to solve similar problems.

6.2 Training Convergence

We prove that under reasonable assumptions, the meta-training objective of the hyper-network converges at a linear rate.

Theorem 6.1 (Convergence of hyper-network training) *Suppose the loss function $\mathcal{L}_{\text{hyper}}(\phi) = \mathcal{L}_{\text{task}}(\phi) + \lambda\mathcal{L}_{\text{consist}}(\phi)$ of generator G_ϕ is differentiable in ϕ and satisfies the Polyak-Lojasiewicz (PL) condition: there exists $\mu > 0$ such that $\frac{1}{2}\|\nabla\mathcal{L}_{\text{hyper}}(\phi)\|^2 \geq \mu(\mathcal{L}_{\text{hyper}}(\phi) - \mathcal{L}_{\text{hyper}}^*)$, where $\mathcal{L}_{\text{hyper}}^*$ is the global minimum. If the skeleton library covers the support of all problem types in the training distribution and the skeleton consistency loss weight $\lambda > 0$, then gradient descent $\phi_{t+1} = \phi_t - \eta\nabla\mathcal{L}_{\text{hyper}}(\phi_t)$ converges linearly:*

$$\mathcal{L}_{\text{hyper}}(\phi_t) - \mathcal{L}_{\text{hyper}}^* \leq (1 - \eta\mu)^t (\mathcal{L}_{\text{hyper}}(\phi_0) - \mathcal{L}_{\text{hyper}}^*).$$

Proof sketch: The PL condition is a sufficient condition for gradient-dominated convergence and typically holds in overparameterized neural networks. The skeleton consistency loss aligns the gradient directions from different examples of the same skeleton, reducing gradient variance and making the effective Lipschitz constant smaller, thus ensuring that μ in the PL condition is bounded away from zero. \square

This theorem shows that as long as skeleton extraction quality is sufficiently high (examples of the same skeleton can indeed be solved using the same functional module), the training of the generator is efficient and stable.

6.3 Generalization Bound Analysis

We further analyze the generalization capability of SHR-DS, i.e., the expected performance on unseen new problem types.

Theorem 6.2 (Skeleton-guided generalization bound) *Suppose the skeleton library \mathcal{L}_{skel} covers all problem types in the target mathematical domain, and each skeleton has been validated by at least $k \geq 3$ independent examples. For a new problem q_{new} drawn from the same problem type distribution, the reasoning error rate of SHR-DS satisfies:*

$$P(\text{error}) \leq \sum_{S \in \mathcal{L}_{skel}} P(S \text{ is selected}) \cdot \epsilon_S + \delta_{match},$$

where ϵ_S is the instantiation error rate of skeleton S , i.e., the probability that even when the skeleton type matches correctly, the generated parameters still produce an incorrect answer; and δ_{match} is the skeleton matching failure rate, i.e., the probability that the input type of the new problem cannot be unified with any skeleton’s τ_{in} .

Proof: Expand the error event by the law of total probability, partitioning into “match succeeds but instantiation error” and “match fails.” \square

This bound tells us that two ways to improve the generalization capability of SHR-DS are: (1) reduce the instantiation error rate ϵ_S for each skeleton (by more example validations and stronger skeleton consistency loss); (2) reduce the matching failure rate δ_{match} (by expanding the coverage of the skeleton library). More importantly, unlike the generalization bounds of traditional neural networks, which typically depend on VC dimension or Rademacher complexity (strongly correlated with parameter scale), this bound depends only on the completeness of the skeleton library and the stability of skeleton instantiation, independent of the total parameter count (1B or larger). This is the theoretical root of SHR-DS’s extremely high data efficiency.

7 Systematic Comparison with Existing Paradigms

To highlight the novelty of SHR-DS, we provide a detailed multi-dimensional comparison with existing mainstream paradigms, including reasoning mechanism, knowledge storage, parameter existence form, data efficiency, VRAM footprint, correctness guarantee, and interpretability. The comparison results are summarized in Table 2.

Table 2: Multi-dimensional comparison of SHR-DS with existing reasoning paradigms

Dimension	Standard AR (CoT)	Tree Search (MCTS)	Continuous (PGR)	MoE	SHR-DS (Ours)
Definition of reasoning	Token-by-token construction	Path search	Energy minimization	Expert selection	Obstruction resolution + skeleton transfer
Math. foundation	Probability theory	Decision theory	Stochastic analysis	Statistical learning	Algebraic topology + type theory
Knowledge storage	Dense params	Dense params	Dense params	Multiple dense experts	Language base + hypernetwork + symbolic skeleton library
Parameter existence	Permanent in VRAM	Permanent in VRAM	Permanent in VRAM	All experts permanent	Language stream permanent, reasoning stream transient
Data efficiency	Baseline	Slightly better	Slightly better	Comparable to dense	10–50× improvement
Inference VRAM	High	High	Medium	High (all experts)	< 2 GB (incl. language base)
Language capability	Good	Good	Medium	Good	>95% retention
Reasoning correctness guarantee	None	Depends on search quality	Depends on energy surface	None	Topological boundary conditions + type unification
Interpretability	Low	Medium	Medium	Low	Very high (skeleton steps + homology barcode)
Compositional generalization	Statistical	Statistical	Implicit	Implicit	Explicit skeleton interface matching
Few-shot adaptation	Poor	Poor	Poor	Poor	Excellent (≤ 3 examples)

7.1 Comparison with Standard Autoregressive and Tree Search

Standard autoregressive methods (including CoT) rely entirely on token-level statistical pattern matching, with no structural guarantee of reasoning correctness. Tree search methods (e.g., Tree-of-Thoughts, MCTS) increase the probability of finding correct solutions by expanding the search space, but their search heuristics are still given by probabilistic models, and the expansion of

the search tree in VRAM grows exponentially with width, conflicting with our ultra-low VRAM constraint. SHR-DS, via the symbolic type system of skeletons (Axiom 1) and topological boundary conditions (Theorem 1), decouples reasoning correctness from the probabilistic model, endowing the reasoning process with strict logical and topological guarantees.

7.2 Comparison with Continuous Relaxation Models

Potential-Guided Reasoning (PGR) treats reasoning as Langevin relaxation on an energy landscape, an elegant physical metaphor. However, PGR still requires a pre-trained, VRAM-resident large-scale energy function network, and does not solve the storage bottleneck; simultaneously, purely continuous dynamics struggle to handle discrete logical hard constraints. SHR-DS’s HRR component adopts a discrete algebraic topology framework, naturally suited for the discrete nature of logical propositions, while dual-stream compression solves the storage problem that PGR does not address. In fact, PGR and HRR can be complementary: PGR is responsible for quickly discovering potential obstruction regions in continuous space, while HRR performs precise topological resolution in those regions.

7.3 Comparison with Mixture of Experts (MoE)

As detailed in Section 4, both MoE and SHR-DS employ the idea of “partial activation,” but they are fundamentally different. MoE experts physically exist, all experts permanently reside in VRAM, routing is selective, and there is no guarantee of correctness. SHR-DS reasoning parameters are virtual, generated on-demand by a hyper-network and discarded after use; routing (skeleton matching) is generative and based on decidable type unification. Moreover, the skeleton library of SHR-DS can be infinitely expanded without increasing VRAM footprint, whereas adding a new expert to MoE necessarily increases VRAM.

7.4 Comparison with Parameter-Efficient Fine-Tuning Methods

Methods such as LoRA, quantization, and pruning can compress models to some extent, but all treat linguistic and reasoning capabilities as a unified whole, and compression inevitably leads to loss of some capability. The dual-stream decoupling of SHR-DS fundamentally avoids this trade-off: linguistic capability is preserved via “inheritance + compression,” reasoning capability is realized via “generation + release,” and the two do not drag each other down.

8 Performance Estimation

This section provides well-founded estimates of SHR-DS’s performance across various metrics based on publicly available experimental data from adjacent technologies. All estimates are anchored to known benchmark data and provide conservative and optimistic ranges.

8.1 Benchmark Calibration Data

We select the following adjacent works as calibration baselines:

- **MSE (Multidimensional Self-Elaboration):** On a 100M-parameter mathematical reasoning model, leveraging the same data through three levels of cognitive utilization improved data efficiency by about $3\times$ and accuracy by 8%–18%.
- **TopoAlign (2025):** By aligning topological structures of code with mathematical representations, achieved a 17.77% performance improvement on DeepSeek-Math 7B (BEq@10 metric), though data efficiency improvement was not quantified. We conservatively estimate the implicit data efficiency gain at about $2\text{--}3\times$.
- **uPRM (2025):** Unsupervised process reward models achieved a 15% absolute improvement in step error identification accuracy on mathematical reasoning, and test-time verification surpassed majority voting baseline by 6.9%.
- **GPTQ-4bit (Frantar et al., 2023):** On WikiText-2, 4-bit quantization increased perplexity by only about 2%, retaining $>95\%$ language quality.
- **SCoT (2024) and other few-shot symbolic reasoning works:** Reduced required training samples from hundreds to <10 on elementary mathematics, showing that symbolic structure can greatly improve few-shot generalization.
- **MoE sparse inference (Mixtral $8\times 7\text{B}$):** Activates about 12.9B parameters at inference (out of 46.7B total), with computational savings of about 50%, but VRAM footprint remains the sum of all experts.

8.2 Estimation Methodology and Rationale

SHR-DS simultaneously integrates multiple mechanisms that have been separately validated in the above works: the skeleton type system (inherited from SCT and the symbolic reasoning tradition), topological resolution (inheriting and surpassing the topological ideas of HRR and TopoAlign), hyper-network generation (surpassing MoE on-demand activation), and a compressed language base (inheriting best practices from GPTQ). The synergistic effect of these mechanisms is not additive but multiplicative: skeleton extraction dramatically reduces data requirements for the same problem type, topological resolution guarantees solution quality, and dual-stream compression ensures VRAM feasibility. We provide estimates based on the following reasoning chain:

- **Data efficiency improvement:** MSE’s three-level cognitive utilization already achieved a $3\times$ improvement. SHR-DS further introduces skeleton extraction (capable of learning a problem type from ≤ 3 examples), plus strict transfer guaranteed by the skeleton type system, reducing the required number of examples per problem type from thousands to 3–5. Therefore, $10\text{--}50\times$ sample efficiency improvement is conservative (corresponding to different problem difficulties: simple problems about $10\text{--}20\times$, complex combinatorial problems $40\text{--}50\times$).
- **Accuracy improvement:** MSE’s three-level utilization brought 8%–18% improvement. TopoAlign’s topological alignment brought 17.77% improvement on a 7B model. SHR-DS deepens topological constraints from the representation layer to the reasoning mechanism layer, and has explicit type correctness guarantees, so an overall accuracy improvement of 20%–35% is reasonable, with larger improvements (30%–50%) on complex long-chain reasoning.
- **Language capability retention:** Based on public data from GPTQ-4bit (perplexity increase $< 2\%$), plus additional loss from low-rank decomposition and pruning ($< 3\%$), language capability retention of 95% is a conservative estimate.

- **VRAM footprint:** The detailed budget in Section 4.5 provides a clear argument for < 2 GB at inference, based on actual hardware parameters.
- **Training time:** Because the data efficiency improvement far outweighs the increase in per-step training time, the absolute training time to reach equivalent performance may actually decrease. We estimate that for a 100M configuration, the complete four-stage training on an RTX 3060 Ti takes about 3–7 days.

8.3 Detailed Performance Estimation Table

Table 3: Performance estimation across configurations

Metric	100M Config	256M+128M Activation	1B+128M Activation
Data efficiency (reduction in required data)	90%	95%	98%
Math. reasoning accuracy (overall)	+15%–20%	+20%–30%	+25%–35%
Complex long-chain reasoning accuracy	+20%–30%	+30%–45%	+35%–50%
New problem type generalization	+30%–40%	+40%–55%	+45%–60%
Language fluency retention	97%	96%	95%
Inference footprint	VRAM < 1 GB	< 1.5 GB	< 2 GB
Training footprint	VRAM < 4 GB	< 6 GB	< 8 GB (ZeRO-Offload)
Inference speed (per step)	Comparable to dense	$\sim 5\%$ slower	$\sim 10\text{--}15\%$ slower
Total training time (to convergence)	$1.5\text{--}2\times$ dense	$2\text{--}2.5\times$ dense	$2\text{--}3\times$ dense

8.4 Ablation Study Estimates

To quantify the contribution of each component, we estimate the performance degradation after removing key components one by one. These estimates can serve as hypotheses for future experimental validation.

Table 4: Ablation study estimates

Removed Component	Performance Degradation	Explanation
Remove skeleton type system (revert to semantic matching)	Accuracy -8%–12%	Type unification \rightarrow semantic similarity degradation, making skeleton transfer unreliable
Remove topological obstruction resolution (revert to pure skeleton chain)	Complex problem accuracy -15%–25%	Loss of topological correctness guarantee, long-chain reasoning prone to logical jumps
Remove hyper-network generation (revert to fixed expert library)	VRAM +200%–400%	Experts must permanently reside in VRAM; 1B model cannot fit in 8 GB
Remove language base (pure reasoning stream)	Language fluency -80%+	Cannot generate readable text, only outputs mathematical symbols
Remove skeleton consistency loss	Generalization -10%–20%	Generator produces too different parameters for same skeleton across different instances
Remove persistent homology regularization	Reasoning redundancy +15%–25%	Model tends to generate longer redundant reasoning chains

These ablation data indicate that each component of SHR-DS is indispensable; their synergy forms a complete capability loop.

9 Discussion

9.1 Fundamental Innovations of This Work – The Deep Meaning of Triple Decoupling

The fundamental innovations of SHR-DS can be summarized as triple decoupling.

First: Parameter decoupling of language and reasoning. This decoupling acknowledges the functional differentiation between language areas and logic areas in the cerebral cortex. Linguistic capability needs to handle the infinite variety of the open world; its neural basis is inevitably highly redundant, distributed, and statistical. Mathematical reasoning, in contrast, follows precise symbolic rules; its neural basis (or at least its simulable functionality) is sparse, modular, and deterministic. Forcing the two to be encoded in the same parameter space is like demanding a neural network to be simultaneously an excellent painter and a precise calculator—inevitably inefficient. SHR-DS, through dual-stream architecture, provides each capability with its optimal storage and computational form.

Second: Morphological decoupling of knowledge and parameters. In traditional models, knowledge (e.g., “how to use the quadratic formula”) must be implicitly stored as floating-point weights in parameters, leading to a rigid binding between knowledge capacity and parameter count. SHR-DS externalizes mathematical knowledge as explicit symbolic skeletons, stored as human-readable, editable, verifiable 5-tuples, with parameters only generated at runtime by a hyper-network. This brings several profound advantages: knowledge can accumulate indefinitely without increasing VRAM; knowledge can be directly inspected and corrected; knowledge transfer becomes a type unification problem rather than a gradient descent problem.

Third: Criterion decoupling of reasoning correctness and statistical prediction. The proper role of the neural network is to generate “reasonable candidates”—reasonable skeletons, reasonable fillers, reasonable textual explanations. Whether these candidates are “correct” is adjudicated by symbol and topological criteria independent of the neural network: the type checker verifies the legality of skeleton transfer, the homology group verifies whether the filler actually resolves the obstruction. This division resembles the “System 1” (fast intuition) and “System 2” (slow deliberation) dual-process cognition in the brain; SHR-DS implements System 1 with neural networks and System 2 with algebraic topology and type theory.

9.2 Applicability Boundaries and Extensibility

Current applicability boundaries: SHR-DS is highly optimized for structured mathematical reasoning tasks, i.e., those with clear axiomatic systems, encodable with type systems, and whose proof processes can be decomposed into discrete obstruction resolution steps. This covers a wide range of mathematical branches: elementary algebra, geometry, calculus, linear algebra, number theory, and more.

Potential extensions:

- **Code generation and verification:** Programs also have strict type systems and structured syntax; skeleton types can directly correspond to function signatures, and obstruction resolution corresponds to the debugging process.
- **Formal theorem proving:** The topological model of SHR-DS has natural correspondence with the underlying logic of theorem provers (e.g., Lean, Coq), and could serve as search guidance for automated proof.
- **Legal and logical reasoning:** Legal statutes also have skeleton structures in their elements and derivations.

Inapplicable domains: Open-domain chit-chat, literary creation, sentiment analysis, and other tasks without clear structural constraints are not handled by SHR-DS’s reasoning stream, and are fully carried by the language base. This is not a defect but a design choice—we prefer to excel in specialized domains rather than sacrificing VRAM and reasoning reliability for generic capability.

9.3 Synergy with Self-Evolution Sandbox

SHR-DS can be seamlessly integrated into a self-evolution sandbox (Pillar 4), forming a continuously evolving reasoning system:

- **Automatic skeleton discovery:** The proposer in the sandbox generates new problem types; the solver attempts to solve them using existing skeleton combinations; the verifier (deterministic mathematical engine) judges answer correctness. If existing skeletons cannot cover, the system abstracts the successful solution path into a new skeleton, which is added to the skeleton library after passing type checker and topological verification.
- **Continuous generator fine-tuning:** New problem-skeleton pairs generated by the sandbox continuously expand the training set, and the hyper-network G_ϕ is continuously fine-tuned in the background to cover new skeletons, without retraining the entire system.

- **Periodic language base updates:** Using mathematical dialogue data generated by the sandbox, the language base can be periodically lightly fine-tuned to maintain freshness and accuracy of language expression.

This closed loop makes SHR-DS a self-evolving reasoning entity that “gets smarter with use.” Over time, its skeleton library enriches, its generator refines, and its reasoning capability continuously improves, while hardware costs remain constant.

9.4 Limitations

Despite the significant theoretical and estimated performance advantages of SHR-DS, we honestly point out its current limitations:

- **Cold-start cost of the skeleton library:** The initial 300–500 skeletons still require human expert or strong AI model assistance for annotation, a non-negligible human bottleneck.
- **Theoretical challenge of semantic complex completeness:** The “sufficiency” in Theorem 1 depends on the completeness of the semantic complex for the target mathematical theory, a condition difficult to fully guarantee in practice. An incomplete complex may cause a few correct derivations to be topologically judged impossible.
- **Hyper-network generalization to unseen skeleton combinations:** Although the skeleton consistency loss promotes generalization, the generator’s performance may degrade for skeleton combinations with extremely novel type structures not seen during training. This is an inherent challenge of meta-learning.
- **Numerical stability in long chains:** In extremely long resolution processes (hundreds of steps), the representation of topological obstructions and homology computations may accumulate numerical errors, affecting resolution precision.

We believe that the above limitations can be gradually overcome through future research, especially with the development of automatic skeleton discovery algorithms and efficient homology computation libraries.

10 Conclusion and Future Directions

10.1 Summary of This Work

This paper addresses the fundamental challenge of “deep mathematical reasoning under extreme resource constraints” by proposing the Structural Homological Resolution with Dual-Stream Parameter Compression (SHR-DS) unified framework. The framework builds on three core innovations, constructing a complete theoretical system and engineering path.

First, at the reasoning mechanism level, we propose Homological Resolution Reasoning (HRR). HRR no longer treats mathematical reasoning as sequence generation or path search, but as a progressive resolution of relative homology groups in a semantic complex. The correctness of reasoning no longer depends on the contingency of probabilistic sampling, but is structurally guaranteed by boundary conditions in algebraic topology. We established the Topological Criterion for Reasoning (Theorem 1), proved the algebraic effect of filler operations (Theorem 2), and

revealed the strict equivalence between reasoning chains and obstruction resolution sequences (Theorem 4).

Second, at the knowledge acquisition level, we propose Structural Cognition Training (SCT). SCT enables the model to extract transferable solution skeletons from very few examples (≤ 3), and through the introduction of a formal type system (Definition 5) and the Type Criterion for Skeleton Transfer (Axiom 1), transforms skeleton transfer from vague semantic matching to mechanically decidable type unification. This contribution dramatically improves data efficiency, enabling the model to “learn a class from one example” in the manner of a human top student.

Third, at the architectural engineering level, we propose the Dual-Stream Parameter Compression (DPC) architecture. DPC fundamentally decouples linguistic fluency and mathematical reasoning ability into heterogeneous parameter streams: linguistic capability is inherited from a pre-trained 1B dense model and compressed via GPTQ quantization, low-rank decomposition, and structured pruning into a ~ 500 MB permanent language base, retaining over 95% of language quality; reasoning capability is generated on-demand by a hyper-network conditioned on skeletons into transient expert parameters (12–128M), discarded after use. This architecture enables 1B-level reasoning capability to run entirely on a consumer-grade GPU with 8 GB VRAM, with inference VRAM < 2 GB and training VRAM < 8 GB (with CPU offloading).

The most important theoretical contribution of this paper is the Skeleton-Filler Duality Theorem (Theorem 5). This theorem reveals for the first time the profound mathematical connection between the solution skeletons of SCT and the homological fillers of HRR: a solution skeleton satisfying the upper-closed condition corresponds to a unique composite filler in the semantic complex, and skeleton transfer is equivalent to type-guided obstruction resolution. This theorem unifies two independent innovative lines under a single mathematical framework, providing a solid theoretical foundation for SHR-DS.

In terms of training paradigm, we extend Multidimensional Self-Elaboration to Four-Tier Cognitive Progressive Training: Surface Imitation \rightarrow Structural Reconstruction \rightarrow Skeleton Extraction and Generator Meta-Training \rightarrow Topological Cognition. These four tiers simulate the complete cognitive development of human mathematics learning from imitation to understanding, from abstraction to insight, introduced via a curriculum schedule to ensure stable and efficient training.

10.2 Future Research Directions

Automatic skeleton discovery and formal verification. The current cold-start of the skeleton library still requires a small amount of human annotation. Combining the self-evolution sandbox and stronger type inference techniques to achieve automatic discovery of new skeletons from unannotated data with type correctness proofs is a key next milestone.

Automatic completion of the semantic complex. The sufficiency in Theorem 1 depends on the completeness of the semantic complex for the target mathematical theory. Studying how to automatically construct a complete semantic complex from an axiomatic system, or providing bounds on reasoning reliability in incomplete cases, is a deep theoretical problem.

Theoretical optimization of hyper-network architecture. Explore more efficient parameter generation architectures (e.g., diffusion-based generators, retrieval-augmented generation) to

further reduce generation time and improve the generalization quality of generated parameters.

Hardware-algorithm co-design. Design dedicated hardware acceleration units for the characteristics of “transient parameters,” enabling parameter generation and reasoning computation to execute in parallel pipelines, hiding generation latency entirely within computation latency.

Theoretical framework for cross-domain transfer. Generalize the skeleton type system and topological resolution mechanism of SHR-DS to domains beyond mathematics, establishing a unified theory of “structural reasoning.”

Deep integration with formal proof systems. Translate SHR-DS’s skeleton library and filler resolution directly into formal proofs in Lean or Coq, achieving end-to-end automatic transformation from natural language problems to formal proofs.

10.3 Concluding Remarks

SHR-DS is not just a technical solution; it represents an understanding of the nature of intelligence. It shows us that intelligence does not necessarily require mountains of data and parameters; it can also be obtained through exquisite structure and correct abstraction. In an era where hardware resources are increasingly becoming a bottleneck for AI development, this approach of “replacing computational brute force with algorithmic wisdom” may well be the necessary path to true intelligence.

References

- [1] Wei, J., Wang, X., Schuurmans, D., et al. (2022). Chain-of-Thought Prompting Elicits Reasoning in Large Language Models. *Advances in Neural Information Processing Systems 35 (NeurIPS 2022)*.
- [2] Yao, S., Yu, D., Zhao, J., et al. (2024). Tree of Thoughts: Deliberate Problem Solving with Large Language Models. *Advances in Neural Information Processing Systems 36 (NeurIPS 2023)*.
- [3] Chen, R. T. Q., Rubanova, Y., Bettencourt, J., & Duvenaud, D. (2018). Neural Ordinary Differential Equations. *Advances in Neural Information Processing Systems 31 (NeurIPS 2018)*.
- [4] Silver, D., Hubert, T., Schrittwieser, J., et al. (2017). Mastering Chess and Shogi by Self-Play with a General Reinforcement Learning Algorithm. *arXiv preprint arXiv:1712.01815*.
- [5] Silver, D., Schrittwieser, J., Simonyan, K., et al. (2017). Mastering the Game of Go without Human Knowledge. *Nature*, 550(7676), 354–359.
- [6] Hatcher, A. (2002). *Algebraic Topology*. Cambridge University Press.
- [7] Edelsbrunner, H., & Harer, J. (2010). *Computational Topology: An Introduction*. American Mathematical Society.
- [8] Zomorodian, A., & Carlsson, G. (2005). Computing Persistent Homology. *Discrete & Computational Geometry*, 33(2), 249–274.

- [9] Carlsson, G. (2009). Topology and Data. *Bulletin of the American Mathematical Society*, 46(2), 255–308.
- [10] Fedus, W., Zoph, B., & Shazeer, N. (2022). Switch Transformers: Scaling to Trillion Parameter Models with Simple and Efficient Sparsity. *Journal of Machine Learning Research*, 23(120), 1–39.
- [11] Shazeer, N., Mirhoseini, A., Maziarz, K., et al. (2017). Outrageously Large Neural Networks: The Sparsely-Gated Mixture-of-Experts Layer. *International Conference on Learning Representations (ICLR 2017)*.
- [12] Frantar, E., Ashkboos, S., Hoefler, T., & Alistarh, D. (2023). GPTQ: Accurate Post-Training Quantization for Generative Pre-trained Transformers. *International Conference on Learning Representations (ICLR 2023)*.
- [13] Hu, E. J., Shen, Y., Wallis, P., et al. (2022). LoRA: Low-Rank Adaptation of Large Language Models. *International Conference on Learning Representations (ICLR 2022)*.
- [14] Ha, D., Dai, A. M., & Le, Q. V. (2017). HyperNetworks. *International Conference on Learning Representations (ICLR 2017)*.
- [15] DeepSeek-AI. (2024). DeepSeek-R1: Incentivizing Reasoning Capability in LLMs via Reinforcement Learning. *arXiv preprint arXiv:2501.12948*.
- [16] Rajbhandari, S., Rasley, J., Ruwase, O., & He, Y. (2020). ZeRO: Memory Optimizations Toward Training Trillion Parameter Models. *Proceedings of the International Conference for High Performance Computing, Networking, Storage and Analysis (SC20)*.
- [17] SOLAR: Topological Optimization of Chain-of-Thought Reasoning. (2025). *Advances in Neural Information Processing Systems*.
- [18] TopoAlign: Topological Alignment for Code-Math Reasoning. (2025). *International Conference on Machine Learning*.
- [19] uPRM: Unsupervised Process Reward Models for Mathematical Reasoning. (2025). *Advances in Neural Information Processing Systems*.
- [20] SCoT: Symbolic Chain-of-Thought for Small-Sample Mathematical Reasoning. (2024). *Proceedings of the AAAI Conference on Artificial Intelligence*.
- [21] HyperRouter: HyperNetwork-based Routing for Mixture-of-Experts. (2024). *International Conference on Learning Representations*.
- [22] Hindley, J. R. (1969). The Principal Type-Scheme of an Object in Combinatory Logic. *Transactions of the American Mathematical Society*, 146, 29–60.
- [23] Milner, R. (1978). A Theory of Type Polymorphism in Programming. *Journal of Computer and System Sciences*, 17(3), 348–375.
- [24] Pierloot, C., et al. (2025). DeepSeekMath: Pushing the Limits of Mathematical Reasoning in Open Language Models. *arXiv preprint*.

- [25] Yang, A., et al. (2024). Qwen2.5 Technical Report. *arXiv preprint arXiv:2412.15115*.
- [26] Abdin, M., et al. (2024). Phi-3 Technical Report: A Highly Capable Language Model Locally on Your Phone. *arXiv preprint arXiv:2404.14219*.
- [27] Jiang, A. Q., et al. (2024). Mixtral of Experts. *arXiv preprint arXiv:2401.04088*.
- [28] Dao, T., Fu, D. Y., Ermon, S., Rudra, A., & Ré, C. (2022). FlashAttention: Fast and Memory-Efficient Exact Attention with IO-Awareness. *Advances in Neural Information Processing Systems 35 (NeurIPS 2022)*.
- [29] Su, J., Lu, Y., Pan, S., Murtadha, A., Wen, B., & Liu, Y. (2024). RoFormer: Enhanced Transformer with Rotary Position Embedding. *Neurocomputing*, 568, 127063.

A Key Algorithm Pseudocode

A.1 Homology-Guided Reasoning

Listing 1: Homology-Guided Reasoning

```
1 Algorithm: Homology-Guided Reasoning
2 Input: Premises  $P = \{p_1, \dots, p_n\}$ , Conclusion  $C$ ,
3         Semantic Complex  $K$ , Skeleton Library  $L_{\text{skel}}$ ,
4         Hyper-Network Generator  $G$ 
5 Output: Reasoning filler sequence  $D_1, \dots, D_m$ 
6         or "Reasoning Infeasible"
7
8 1.  $X_P \leftarrow \text{Union}_{\{i=1..n\}} X_{\{p_i\}}$            // Premise subcomplex
9 2.  $X_C \leftarrow X_C$                                  // Conclusion subcomplex
10 3.  $O_0 \leftarrow H_*(X_P \cup X_C, X_P)$              // Initial reasoning
    obstruction
11
12 4. if  $O_0 \neq 0$  then
13 5.   return []                                       // Reasoning already
    complete
14 6. end if
15
16 7.  $K_0 \leftarrow X_P$ 
17 8.  $t \leftarrow 0$ 
18
19 9. while  $O_t \neq 0$  and  $t < \text{MAX\_STEPS}$  do
20 10.   // Step 1: Attempt skeleton matching
21 11.    $\text{type}_q \leftarrow \text{InferType}(C, K_t)$ 
22 12.    $\text{matched\_skeletons} \leftarrow \text{TypeUnify}(L_{\text{skel}}, \text{type}_q)$ 
23
24 13.   // Step 2: Generate reasoning parameters
25 14.   for each  $S$  in  $\text{matched\_skeletons}$  do
26 15.     if not  $\text{Cached}(S)$  then
27 16.        $c \leftarrow \text{Encode}(E_{\text{psi}}, S)$ 
28 17.        $z \leftarrow \text{SampleNoise}()$ 
29 18.        $W_{\text{expert}}[S] \leftarrow G(c, z)$ 
30 19.        $\text{Cache}(S, W_{\text{expert}}[S])$ 
31 20.     else
32 21.        $W_{\text{expert}}[S] \leftarrow \text{LoadCache}(S)$ 
33 22.     end if
34 23.   end for
35
36 24.   // Step 3: Identify generators
37 25.    $\text{generators} \leftarrow \text{IdentifyGenerators}(O_t)$ 
38
39 26.   // Step 4: Select optimal filler
40 27.    $\text{best\_filler} \leftarrow \text{None}$ 
41 28.    $\text{best\_score} \leftarrow -\text{inf}$ 
42 29.   for each filler  $D$  in  $\text{CandidateFillers}(\dots)$  do
43 30.     if  $dD \text{ subseteq } K_t$  then                 // Boundary check
```

```

44 31.         score <- ScoreFiller(D, generators)
45 32.         if score > best_score then
46 33.             best_filler <- D
47 34.             best_score <- score
48 35.         end if
49 36.     end if
50 37. end for
51
52 38.     if best_filler is None then
53 39.         return "Reasoning Infeasible"
54 40.     end if
55
56 41.     // Step 5: Attach filler
57 42.     K_{t+1} <- K_t U {best_filler}
58 43.     D_{t+1} <- best_filler
59 44.     O_{t+1} <- H_*(K_{t+1} U X_C, K_{t+1})
60 45.     t <- t + 1
61 46. end while
62
63 47. if O_t =~ 0 then
64 48.     return [D_1, ..., D_t]
65 49. else
66 50.     return "Reasoning not completed within max steps"
67 51. end if

```

A.2 Hyper-Network Generator Architecture

Listing 2: Forward pass of hyper-network generator G

```

1 Algorithm: Forward pass of hyper-network generator  $G$ 
2 Input: Condition vector  $c$  in  $\mathbb{R}^{\{d_c\}}$ ,
3       noise seed  $z$  in  $\mathbb{R}^{\{d_z\}}$ ,
4       target parameter shape  $(d_{out}, d_{in})$ 
5 Output: Generated reasoning expert parameters  $W$ 
6
7 1. // Step 1: Concatenate condition and noise
8 2.  $h_0 \leftarrow \text{Concat}(c, z)$  //  $\mathbb{R}^{\{d_c + d_z\}}$ 
9
10 3. // Step 2: Progressive generation
11 4.  $h_1 \leftarrow \text{ReLU}(\text{CondBN}(\text{Linear}(h_0), c))$  //  $\mathbb{R}^{\{h_1\}}$ 
12 5.  $h_2 \leftarrow \text{ReLU}(\text{CondBN}(\text{Linear}(h_1), c))$  //  $\mathbb{R}^{\{h_2\}}$ 
13 6. ...
14 7.  $h_L \leftarrow \text{ReLU}(\text{CondBN}(\text{Linear}(h_{\{L-1\}}), c))$  //  $\mathbb{R}^{\{h_L\}}$ 
15
16 8. // Step 3: Final mapping to target shape
17 9.  $W_{\text{raw}} \leftarrow \text{Linear}(h_L)$  //  $\mathbb{R}^{\{d_{out} \times d_{in}\}}$ 
18 10.  $W \leftarrow \tanh(W_{\text{raw}}) * \text{scale\_factor}$  // Constrain range
19
20 11. return  $W$ 

```

A.3 Type Unification Algorithm for Skeleton Matching

Listing 3: Skeleton matching via type unification

```
1 Algorithm: Skeleton matching via type unification
2 Input: Skeleton library L_skel, target type tau_target
3 Output: List of matched skeletons and substitutions
4
5 1.  matched <- []
6
7 2.  for each skeleton S in L_skel do
8 3.    tau_skel <- tau_in(S)
9 4.    theta <- {} // Empty substitution
10 5.   if Unify(tau_skel, tau_target, theta) then
11 6.     matched.append((S, theta))
12 7.   end if
13 8. end for
14
15 9.  return matched
16
17 // Unification function
18 10. function Unify(tau_1, tau_2, theta):
19 11.   tau_1 <- Apply(theta, tau_1)
20 12.   tau_2 <- Apply(theta, tau_2)
21
22 13.   if tau_1 == tau_2 then
23 14.     return true
24 15.   else if IsTypeVar(tau_1) then
25 16.     if Occurs(tau_1, tau_2) then return false
26 17.     theta[tau_1] <- tau_2
27 18.     return true
28 19.   else if IsTypeVar(tau_2) then
29 20.     return Unify(tau_2, tau_1, theta)
30 21.   else if IsComposite(tau_1) and IsComposite(tau_2) then
31 22.     if Constructor(tau_1) != Constructor(tau_2) then
32 23.       return false
33 24.     end if
34 25.     for each (arg_1, arg_2) in Args(tau_1) x Args(tau_2) do
35 26.       if not Unify(arg_1, arg_2, theta) then
36 27.         return false
37 28.       end if
38 29.     end for
39 30.     return true
40 31.   else
41 32.     return false
42 33.   end if
43 34. end function
```

A.4 Curriculum Scheduling for Four-Tier Cognitive Progressive Training

Listing 4: SHR-DS training main loop

```

1 Algorithm: SHR-DS training main loop
2 Input: D_general, D_math, D_skel, D_val
3 Output: Trained B_lang, E_psi, G_phi, W_gate
4
5 1. Initialize B_lang from pretrained 1B model
6 2. Initialize E_psi, G_phi, W_gate randomly
7
8 3. total_steps <- TRAINING_STEPS
9 4. for step in 1..total_steps do
10 5.     progress <- step / total_steps
11 6.     a1, a2, a3, a4 <- CurriculumWeights(progress)
12
13 7.     // Stage 1: Surface Imitation
14 8.     if a1 > 0 then
15 9.         batch <- Sample(D_general U D_math)
16 10.        L1 <- CrossEntropy(B_lang, batch)
17 11.        else L1 <- 0
18
19 12.        // Stage 2: Structural Reconstruction
20 13.        if a2 > 0 then
21 14.            batch <- Sample(D_math)
22 15.            L2 <- 0
23 16.            for each (q, r, a) in batch do
24 17.                W_temp <- G_phi(E_psi(dummy_skeleton))
25 18.                P_assembled <- Assemble(B_lang, W_temp, W_gate)
26 19.                L2 <- L2 - log P_assembled(r | q, a)
27 20.            end for
28 21.            L_type <- CrossEntropy(TypeHead(B_lang), batch.types)
29 22.            L2 <- L2 + 0.1 * L_type
30 23.        else L2 <- 0
31
32 24.        // Stage 3: Skeleton Extraction + Generator Meta-Training
33 25.        if a3 > 0 then
34 26.            batch <- Sample(D_skel)
35 27.            L_skel <- -log P(E_psi)(S* | q, r)
36 28.            L_hyper <- 0
37 29.            for each S in batch.skeletons do
38 30.                c <- E_psi(S)
39 31.                z1, z2 <- SampleNoise()
40 32.                W1, W2 <- G_phi(c, z1), G_phi(c, z2)
41 33.                L_task <- CrossEntropy(Assemble(B_lang, W1), batch[S
42 34.                L_consist <- CosineDistance(Rep(W1), Rep(W2))
43 35.                L_hyper <- L_hyper + L_task + lambda * L_consist
44 36.            end for
45 37.            L3 <- L_skel + L_hyper
46 38.        else L3 <- 0
47
48 39.        // Stage 4: Topological Cognition

```

```

49 40.     if a4 > 0 then
50 41.         batch <- Sample(D_math)
51 42.         L_HRR <- 0
52 43.         for each (q, filler_seq) in batch do
53 44.             K <- BuildKnowledgeComplex(q)
54 45.             for each D* in filler_seq do
55 46.                 O <- ComputeObstruction(K)
56 47.                 D_pred <- Forward(B_lang, G_phi, K, O)
57 48.                 L_HRR <- L_HRR - log P(D_pred = D* | K, O)
58 49.                     + lambda_boundary * BoundaryLoss(D_pred, K)
59 50.                 K <- K U {D_pred}
60 51.             end for
61 52.         end for
62 53.         L_align <- AlignmentLoss(L_skel)
63 54.         L_barcode <- BarcodeRegularizer(D_val)
64 55.         L4 <- L_HRR + L_align + L_barcode
65 56.     else L4 <- 0
66
67 57.     // Total loss
68 58.     L_total <- a1*L1 + a2*L2 + a3*L3 + a4*L4
69
70 59.     L_total.backward()
71 60.     optimizer.step()
72 61.     optimizer.zero_grad()
73
74 62.     if step % SAVE_INTERVAL == 0 then
75 63.         SaveCheckpoint(B_lang, E_psi, G_phi, W_gate)
76 64.         UpdateSkeletonCache(G_phi, E_psi, L_skel)
77 65.     end if
78 66. end for
79
80 67. return B_lang, E_psi, G_phi, W_gate

```

B Recommended Hyperparameter Configurations

B.1 100M Configuration (Fast Validation)

Table 5: 100M Configuration

Hyperparameter	Value
Total parameters	100M
Activated parameters	12–20M
Hidden dimension	768
Transformer layers	12
Attention heads	12
Language base parameters	~80M (after compression)
Hyper-network generator parameters	~8M
Encoder parameters	~0.8M
Skeleton library size	200–500 skeletons
Context length	4096 tokens
Batch size	8
Learning rate (Stage 1)	3e-4
Learning rate (Stages 2–4)	1e-4
Skeleton consistency loss weight	0.2
λ_{consist}	
Boundary loss weight $\lambda_{\text{boundary}}$	0.1
Training steps	50,000
Optimizer	AdamW ($\beta_1 = 0.9, \beta_2 = 0.95$)
DeepSpeed config	ZeRO-2 + CPU Offload

B.2 256M Total / 128M Activation Configuration (Main Experiments)

Table 6: 256M+128M Configuration

Hyperparameter	Value
Total parameters	256M
Activated parameters	128M (max)
Hidden dimension	1024
Transformer layers	18–24
Attention heads	16
Language base parameters	~200M (after compression)
Hyper-network generator parameters	~12M
Encoder parameters	~1.2M
Skeleton library size	500–1000 skeletons
Context length	16384 tokens
Batch size	4–8
Learning rate (Stage 1)	3e-4
Learning rate (Stages 2–4)	1e-4
Skeleton consistency loss weight	0.2
λ_{consist}	
Boundary loss weight $\lambda_{\text{boundary}}$	0.1
Training steps	80,000
Optimizer	AdamW ($\beta_1 = 0.9, \beta_2 = 0.95$)
DeepSpeed config	ZeRO-2 + CPU Offload + Activation Checkpointing

B.3 1B Total / 128M Activation Configuration (Extreme Demonstration)

Table 7: 1B+128M Configuration

Hyperparameter	Value
Total parameters	1B
Activated parameters	128M (max)
Hidden dimension	1536
Transformer layers	24–32
Attention heads	24
Language base parameters	~300M (compressed from 1B parent)
Hyper-network generator parameters	~18M
Encoder parameters	~2M
Skeleton library size	1000–2000 skeletons
Context length	16384 tokens
Batch size	2–4
Learning rate (Stage 1)	2e-4
Learning rate (Stages 2–4)	8e-5
Skeleton consistency loss weight λ_{consist}	0.3
Boundary loss weight $\lambda_{\text{boundary}}$	0.15
Training steps	120,000
Optimizer	AdamW ($\beta_1 = 0.9, \beta_2 = 0.95$)
DeepSpeed config	ZeRO-3 + CPU Offload + Activation Checkpointing + Parameter Offload

C Symbol Table

Table 8: Symbol Table

Symbol	Meaning	First Appearance
K	Semantic complex	Definition 1
A	Set of atomic propositions	Definition 1
X_p	Subcomplex embedding of proposition p	Definition 2
$\mathcal{O}(P, C)$	Reasoning obstruction (relative homology group)	Definition 3
H_*	Homology groups (direct sum over degrees)	Definition 3
Δ_k	k -reasoning filler	Definition 4
$\partial\Delta$	Boundary of a simplex	Definition 4
\mathcal{T}	Mathematical type system	Definition 5
$\tau_{\text{in}}, \tau_{\text{out}}$	Input/output types of a skeleton	Definition 6
Σ	Skeleton step sequence	Definition 6
Γ	Type context	Definition 6
S	Solution skeleton (5-tuple)	Definition 6
θ	Type substitution	Axiom 1
Δ_S	Composite filler corresponding to skeleton	Theorem 5
ρ_i	Persistence of homology generator	Definition 7
(B, D)	Skeleton birth-death pair	Definition 7
B_{lang}	Language base	§4.2
G_ϕ	Hyper-network generator	§4.3
E_ψ	Skeleton encoder	§4.3
c	Condition vector	§4.3
z	Noise seed	§4.3
W_{expert}	Generated reasoning expert parameters	§4.3
W_{gate}	Gating fusion layer weights	§4.4
$\alpha^{(\ell)}$	Gating fusion coefficient	§4.4
$\mathcal{L}_{\text{hyper}}$	Total hyper-network loss	§5.4
$\mathcal{L}_{\text{consist}}$	Skeleton consistency loss	§5.4
\mathcal{L}_{HRR}	Homological resolution loss	§5.5
$\mathcal{L}_{\text{align}}$	Skeleton-filler alignment loss	§5.5
$\mathcal{L}_{\text{barcode}}$	Persistent homology regularization loss	§5.5

D Terminology Table

Table 9: Terminology Table

Term	Abbreviation
Structural Homological Resolution with Dual-Stream Parameter Compression	SHR-DS
Homological Resolution Reasoning	HRR
Structural Cognition Training	SCT
Dual-Stream Parameter Compression	DPC
Multidimensional Self-Elaboration	MSE
Potential-Guided Reasoning	PGR
Semantic Complex	—
Reasoning Obstruction	—
Reasoning Filler	—
Solution Skeleton	—
Skeleton-Filler Duality	—
Persistent Homology	—
Reasoning Barcode	—
Hyper-Network Generator	—
Language Base	—
Reasoning Stream / Language Stream	—
Upper-Closed Condition	—
Type Unification	—
Skeleton Transfer	—
Obstruction Resolution	—
Four-Tier Cognitive Progression	—
Self-Evolution Sandbox	—

Acknowledgments

This project has evolved through multiple rounds of deep deliberation, from the initial naive vision of “training a 100M small model specialized in mathematics,” to the energy landscape framework of Potential-Guided Reasoning (PGR), to the topological obstruction resolution of Homological Resolution Reasoning (HRR), to the skeleton extraction and transfer of Structural Cognition Training (SCT), to the language-reasoning heterogeneous decoupling of Dual-Stream Parameter Compression (DPC)—each step has been driven by a pursuit of deeper theory and a challenge to engineering limits.

This project is a solitary exploration of a grand problem by one person. It is extremely difficult both theoretically and practically, challenging many existing paradigms at their foundations. But precisely because it is difficult, it is worth doing; precisely because it is new, it may be “unprecedented.”

Thanks to NVIDIA for providing cost-effective consumer-grade GPUs like the RTX 3060 Ti for individual developers, making “8 GB VRAM extreme experiments” possible. Thanks to the open-source community for contributing tools like DeepSpeed, GPTQ, and FlashAttention, which provide critical infrastructure for ultra-low-resource training. Thanks to the algebraic topology community for over a century of profound theories—Hopfield’s energy landscapes, Neural ODEs’ continuous depth, Hatcher’s algebraic topology textbooks, Zomorodian and Carlsson’s persistent homology algorithms—without these intellectual cornerstones, the HRR and PGR frameworks could not have been built.

Thanks to all frontier researchers in mathematical reasoning, few-shot learning, and parameter-efficient training, whose work provided solid calibration baselines for the performance estimates in this paper.

Finally, this work is dedicated to all who believe that “algorithmic wisdom can transcend computational brute force.”