

Mixture-of-Models: Unifying Heterogeneous Agents via N-Way Self-Evaluating Deliberation

Tims Pecerskis
Peeramid Labs
 tim@peeramid.xyz

Aivars Smirnovs
Peeramid Labs
 aivars@peeramid.xyz

The AI Futures Collective
Peeramid Labs

January 26, 2026 (*Original Preprint Released: January 13, 2026*)

Abstract

The transition from static pre-training to dynamic inference-time compute necessitates new architectures for harnessing collective intelligence. This paper introduces the N-Way Self-Evaluating Deliberation (NSED) protocol, a Runtime Mixture-of-Models (MoM) architecture that constructs emergent composite models from a plurality of distinct expert agents. Unlike traditional Mixture-of-Experts (MoE) which rely on static gating networks, NSED proposes a Dynamic Expertise Broker—a framework designed to treat model selection as a variation of the Knapsack Problem, binding heterogeneous checkpoints to functional roles based on cost constraints. At the execution layer, we formalize deliberation as a Recurrent Deliberation Topology, where the consensus state loops back through a semantic forget gate (γ) to enable iterative refinement without proportional VRAM scaling. Key components include an orchestration fabric for trustless N-to-N peer review and a Quadratic Voting activation function for non-linear consensus. Empirical validation on challenging benchmarks (AIME 2025, LiveCodeBench) demonstrates that this topology allows ensembles of small (<20B) consumer-grade models to match or exceed the performance of state-of-the-art 100B+ parameter models. Furthermore, we propose an empirical utility model ($R^2 \approx 0.99$ on test data) that characterizes consensus as a trade-off between signal extraction and contextual noise accumulation, establishing a mathematical basis for optimal stopping strategies.

1 Introduction

The trajectory of artificial intelligence has long been defined by the scaling laws of pre-training. However, a fundamental paradigm shift is occurring toward **Inference-Time Compute Scaling**. This shift posits that “System 2” thinking—characterized by deliberate planning and verification—can be synthesized by dynamically allocating resources during generation [1].

Recent architectural innovations have sought to internalize this dynamism. Approaches like Titans [2] and End-to-End Test-Time Training (TTT-E2E) [3] demonstrate that models can “learn to memorize” or update their parameters dynamically on input streams. Similarly, the Nested Learning paradigm [4] reframes architectures as hierarchies of optimization problems. While powerful, current implementations remain monolithic or opaque (e.g., OpenAI o1, DeepSeek-R1) [5, 6]. Consequently, the field has turned to *agentic workflows* to circumvent these rigidities.

However, current architectures face a critical dichotomy between static efficiency and topological rigidity. Standard Mixture-of-Experts (MoE) systems rely on fixed gating networks that route at the token level, suffering from a “Granularity Mismatch” where low-level routing fails to capture high-level semantic domains [7]. Conversely, agentic frameworks like Mixture-of-Agents (MoA) [8] and **Chain of Agents (CoA)** [9] operate as feed-forward Directed Acyclic Graphs (DAGs). While effective for one-pass consensus, these systems impose linear memory

costs with depth ($O(N)$) and suffer from error propagation, where upstream faults cascade irreversibly downstream. Furthermore, standard ensembles frequently succumb to “herding,” where a majority of mediocre agents overpower a lone expert [10–12], and context management still poses challenge [13]. Crucially, both paradigms share two fundamental deficits: they are bound by a **DAG topology**—lacking the recurrent inductive bias required for thermodynamic refinement—and they rely on an **opaque, central aggregator** directly in the data plane. This centralization creates a transparency bottleneck, preventing the emergence of a truly trustless, self-correcting consensus.

We argue that robust reasoning requires a transition from linear pipelines to **Recurrent Cognitive Cycles**. This paper introduces the **N-Way Self-Evaluating Deliberation (NSED)** protocol, a computer-implemented method that restructures the multi-agent ensemble not as a feed-forward network, but as a macro-scale Semantic Recurrent Neural Network (SRNN). In this topology, the “hidden state” is the consensus itself, which loops back through the system subject to a semantic decay factor (γ). Unlike MoA, which converges via depth (adding layers), NSED converges over *time* (iterations), allowing for “Deep Thought” without proportional increases in VRAM usage.

To the best of our knowledge, we are the first to propose and analytically describe a full end-to-end topology that is scale-invariant—fitting both model-level recurrence and agent-level workflow. Our contribution centers on the following architectural and theoretical innovations:

1. **Runtime Mixture-of-Models (MoM):** We propose a recursive scaling of the Mixture-of-Experts architecture. Unlike static MoEs, NSED employs a Dynamic Expertise Broker that treats model selection as a runtime optimization problem, binding heterogeneous checkpoints to functional roles based on task complexity and cost constraints.
2. **The Macro-Neuron Topology:** We formalize the deliberation process as Macro-Scale Recurrent Cell (analogous to LSTM), where natural language serves as the embedding and the consensus state acts as the recurrent memory. This topology supports parallel ingestion (Context Widening) and Human-in-the-Loop (HITL) seamless integration.
3. **Trustless Consensus:** To mitigate the authority bias common in heterogeneous ensembles [14], NSED enforces a trustless topology. We implement a hard **Diagonal Mask (D)** at the voting layer, architecturally severing the link between an agent’s proposal and its own vote ($v_{i,i} = 0$).
4. **Dynamic Brokerage:** Unlike static routing maps, our Dynamic Expertise Broker treats model selection as a runtime Knapsack Optimization Problem, solving for the optimal trade-off between latency, cost, and quality to satisfy a specific Service Level Agreement (SLA).
5. **Efficiency-Fatigue Model:** We formulate and empirically validate a parametric scaling model for inference-time ensembles with analytical equation that fits empirical data.

We posit that NSED represents a recursive scaling of neural topology—a Scale-Invariant recursion of well established principles in Artificial Intelligence research. Empirical validation on challenging benchmarks (AIME 2025, LiveCodeBench) demonstrates that this approach allows ensembles of small (<20B) models to match or exceed the performance of state-of-the-art 100B+ parameter models. Furthermore, testing on the **DarkBench** safety suite [15] reveals that the topological governance mechanism inherently reduces sycophancy and our thermodynamic analysis ($R^2 \approx 0.99$) establishes a mathematical basis for **Optimal Stopping**, transforming multi-agent deliberation from a heuristic art into a predictable science, all together paving the road towards robust, verifiable and decentralized Artificial General Intelligence.

2 Related Work

The NSED protocol synthesizes insights across the entire stack of cognitive architecture: from the Macro-Topology of agent graphs and the Thermodynamics of inference-time compute, down to the Mathematical Foundations of high-dimensional consensus. We review these areas to contextualize NSED’s departure from existing feed-forward and static paradigms.

2.1 Topological Limitations in Multi-Agent Systems

Current Multi-Agent Systems (MAS) are predominantly characterized by **Feed-Forward** information flows.

- **Directed Acyclic Graphs (DAGs):** Frameworks like **Mixture-of-Agents (MoA)** [8] and Microsoft’s **AutoGen** [16] structure collaboration as layered DAGs. While effective for consensus, they incur linear memory costs with depth and lack a mechanism for indefinite refinement without expanding the graph.
- **Sequential Chains:** Architectures such as **Chain of Agents (CoA)** [9] rely on sequential “bucket brigade” pass-offs. These suffer from error propagation, where early summarization faults cascade irreversibly downstream.

NSED addresses these topological deficits by formalizing the agent interaction not as a DAG, but as a **Recurrent Neural Network (RNN)**. By maintaining a persistent "consensus state" that loops back to the input gate, NSED achieves iterative refinement (depth) over time steps rather than architectural layers.

2.2 Inference-Time Compute & Scale-Invariant Topology

The paradigm of **Inference-Time Compute Scaling** posits that additional compute during generation can substitute for model scale [1, 17].

- **Internal Recurrence:** Recent architectures like **Titans** [2] and **TTT-E2E** [3] implement this by dynamically updating weights or memory buffers during the forward pass.
- **Black-Box Reasoning:** Models such as **DeepSeek-R1** [5] and **OpenAI o1** utilize Reinforcement Learning on Chain-of-Thought (CoT) to internalize this deliberation.

NSED extends these principles to the **inter-agent** level, effectively creating a recursive pattern repeating neural networking principles on model to model level and relying on Scale-Invariant Topology theory as foundational principles. We operationalize the same “System 2” dynamics found inside these models but lift them into a transparent, white-box protocol. This prevents the opacity of monolithic reasoning models while retaining the benefits of test-time optimization.

2.3 Mechanism Design & Algorithmic Social Choice

While architectures like **Mixture-of-Experts (MoE)** [7] rely on "gating networks" to route tokens, they lack semantic governance. NSED replaces the stochastic router with a deterministic game-theoretic framework.

- **AI-Mediated Deliberation:** Building on the ‘Habermas Machine’ [18] and the **Delphi Method** [19], NSED structures communication to maximize information gain rather than simple agreement.
- **Trustless Aggregation:** Unlike standard ensembles that succumb to ‘herding’ [10], NSED incorporates principles from **Quadratic Voting** [20] and mechanism design. By implementing strict identity masking and diagonal voting masks (), we enforce a **Trustless Topology** that is robust against sycophancy and authority bias, ensuring consensus is driven by semantic merit alone.

2.4 Adaptive Computation & Cognitive Thermodynamics

The pursuit of "Green AI" [21] has driven interest in dynamic halting mechanisms. Early works in Adaptive Computation Time (ACT) demonstrated that neural networks can learn to pause processing once a confidence threshold is met. However, these approaches were primarily restricted to microscopic, per-token probability scalars.

NSED extends this thermodynamic intuition to the macroscopic semantic layer. Unlike *Graph of Thoughts* [22] or *Tree of Thoughts* [23], which optimize for accuracy via static graph expansion, NSED optimizes for the *Pareto frontier of Energy vs. Insight*. By coupling the halting condition to the entropy of the consensus state ($H(S_t)$), we formalize a protocol where compute expenditure can be inextricably linked to information gain, moving beyond static computational budgets.

2.5 Contextual Dynamics & Routing Efficiency

While the capacity of context windows has grown, recent studies reveal a "U-shaped" attention deficit, where models fail to retrieve information located in the middle of long sequences (*Lost in the Middle*, [13]). This phenomenon validates our topological choice of a **Recurrent State** with a decay factor (γ) over a linear append-only log. By iteratively compressing history into a refreshed consensus state, NSED ensures that critical decision boundaries remain within the recency attention head capacity, avoiding the retrieval degradation inherent in standard "Chain of Agents" architectures.

Furthermore, static sparse architectures like *MoE++* [24] demonstrate that "Zero-Computation Experts" can accelerate inference. However, these systems suffer from a Granularity Mismatch, routing at the token level rather than the semantic level. NSED lifts this routing principle to the Macro-Scale. Just as MoE++ dynamically drops tokens to save FLOPs, our *Dynamic Expertise Broker* resolves the granularity mismatch by selecting experts at the *Session Level*, dynamically dropping entire model instances when their domain utility is low.

2.6 Self-Correction & Bootstrapping

The concept of "Self-Correction" has typically been framed as a prompt-engineering tactic. However, recent advances like STaR [25] prove that reasoning traces can serve as high-quality supervision signals for fine-tuning. NSED integrates this insight into the architectural topology itself. By treating the "Consensus State" not just as an output, but as a potential training label for the "Hippocampal" consolidation loop, we bridge the gap between ephemeral Inference-Time Compute [1] and permanent model parameterization.

2.7 Theoretical Foundations of Consensus

To derive our Efficiency-Fatigue Model (Eq. 1), we draw upon foundational results in sequential analysis, geometry, and recent empirical findings on LLM capabilities.

2.7.1 The Verification Asymmetry

A central premise of the NSED protocol is that the ensemble's ability to verify a solution exceeds its ability to generate one zero-shot. This verification gap has been empirically established by Cobbe et al. [26], who demonstrated that training a verifier scales more effectively than fine-tuning a generator for math problems. Furthermore, Lightman et al. [27] showed that "Process Supervision" (step-by-step verification) significantly outperforms outcome-based supervision.

2.7.2 Sequential Analysis

NSED treats deliberation not as a static vote but as a time-series accumulation of evidence. This process is governed by Wald’s Sequential Probability Ratio Test (SPRT) [28]. Wald proved that for a sequential sampling process, the decision boundary (consensus) can be reached with minimal observations if the log-likelihood ratio is integrated over time. Our efficiency parameter Λ is the topological equivalent of the SPRT information rate.

Crucially, for the accumulated log-likelihood ratio to drift toward the correct decision boundary (rather than confirming an error), the ensemble must satisfy the prerequisites of Condorcet’s Jury Theorem [29]. The mean verifier precision must strictly exceed random chance ($\bar{p}_v > 0.5$) to guarantee that the collective drift is positive. Without this condition, the sequential integration would accelerate convergence toward hallucination rather than truth.

2.7.3 Geometric Coverage (N)

The efficacy of heterogeneous ensembles is grounded in high-dimensional geometry. Carathéodory’s Theorem dictates that any point in the convex hull of a d -dimensional feature space can be represented by at most $d + 1$ vertices. In the context of the *Rashomon Set* of valid reasoning paths [30], this implies that a finite number of diverse agents (N) is sufficient to span the solution manifold. Additionally, Cover’s Theorem [31] suggests that projecting a complex classification problem into the high-dimensional space makes the "Truth" linearly separable from "Hallucination," justifying our use of a linear weighted consensus mechanism.

3 Theoretical Framework & System Architecture

This section defines the **General NSED Protocol Specification**. It describes the complete theoretical topology required for a fully autonomous, self-healing cognitive mesh. We propose NSED as a macro-scale analog to Recurrent Neural Networks (RNNs), operating effectively as a neural circuit composed of autonomous agents.

Note on Nomenclature & Isomorphism: While our empirical validation utilizes open-weight models (e.g., Qwen, Gemma), we deliberately abstain from logit normalization or activation steering. Direct vector operations across heterogeneous ensembles—comprising diverse architectures, vocabulary sizes, and parameter scales—introduce complex alignment overheads that were out of our validation scope. Furthermore, we use the terms *agent* and *model* interchangeably, assuming that ReAct-style [32] context ingestion channels are now an intrinsic property of modern LLMs, as tool-use trajectories have become a standard component of post-training datasets.

3.1 Theoretical Design Principles

The topological design of NSED is grounded in two primary mathematical frameworks that dictate how agents are organized and how consensus is derived:

Geometric Capacity: We formally define the NSED ensemble as a high-dimensional kernel machine. A single model (monolith) operates within a fixed embedding dimension d_{model} , limiting it to linearly separating patterns resolvable within that specific manifold. Following **Cover’s Theorem** [31], the probability of a reasoning error being linearly separable approaches 1.0 as the dimensionality of the feature space increases relative to the number of patterns ($N_{dim} > N_{patterns}$). NSED artificially expands this dimensionality by treating each agent’s output not as a final answer, but as a feature vector in N -dimensional meta-space. The Quadratic

Voting layer (σ_{QV}) acts as a non-linear kernel function, allowing the system to perfectly classify (shatter) complex reasoning errors that are mathematically inseparable within any single model’s lower-dimensional latent space.

Temporal Dynamics (Inference-Time Scaling): While static ensembles are governed by Condorcet’s Jury Theorem, NSED operates sequentially, governed by Wald’s Sequential Probability Ratio Test (SPRT). We model the accuracy trajectory $A(t)$ as a thermodynamic competition between Signal Extraction and Entropic Fatigue:

$$\text{Utility}(t) = \underbrace{1 - (1 - p_g)e^{-\Lambda(p_v - p_g)t}}_{\text{Signal Extraction}} - \underbrace{\beta t^2}_{\text{Context Fatigue}} \quad (1)$$

Where:

- t : The discrete deliberation round index ($t \in \{1, 2, \dots, T\}$).
- p_g : The base model’s zero-shot generation precision.
- p_v : The ensemble’s weighted verification precision. Convergence requires $p_v > p_g$.
- Λ : The topological process efficiency constant.
- β : The fatigue coefficient, representing the accumulation of context noise (entropy) over time.

Derivation of the Governing Equation: We derive this formulation based on Wald’s Sequential Probability Ratio Test (SPRT), as defined by Wald [28]. Wald demonstrated that for a sequential process, the accumulated evidence (log-likelihood ratio) grows linearly with time, driving the error rate down to eventually hit model noise floor. This gives rise to our gain term $1 - e^{-\Lambda t}$, where Λ represents the topological efficiency of the evidence integration.

Second, the driving force of this gain is determined by **Condorcet’s Jury Theorem** [29]. Convergence is only mathematically possible if the effective verification capability (p_v) exceeds the generation capability (p_g). We therefore scale the exponent by the "Signal Gap" ($p_v - p_g$), representing the information gain per round.

The quadratic penalty term βt^2 accounts for the **Entropic Fatigue** inherent in recursive semantic processes. This formulation is grounded in the $O(n^2)$ self-attention complexity of the Transformer architecture [33]. As the deliberation history \mathcal{H} expands, the signal-to-noise ratio degrades because the semantic "surface area" for potential hallucinations and sycophantic loops grows quadratically with the total token count. Unlike ideal Bayesian updates, LLMs suffer from context pollution, where the accumulation of failed reasoning traces and peer-disagreements eventually provides a repulsive potential that overwhelms the marginal information gain of further rounds.

3.2 System Control Flow

At the control plane (Fig. 1), the system relies on the Dynamic Broker. Unlike static MoE routers that make simple per-token decisions, the Broker functions as a *Session Constructor* and resource allocator. It is designed as an extensible optimization engine capable of ingesting diverse signals—from static benchmarks to real-time telemetry—to solve for the optimal team composition.

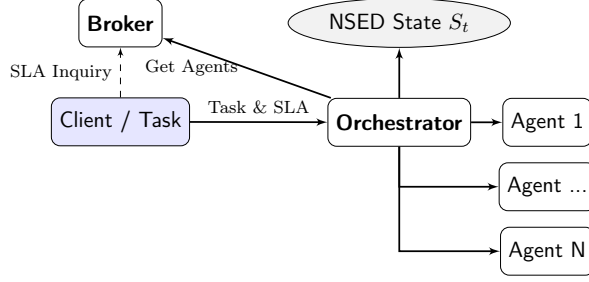


Figure 1: NSED control flow architecture. The **Client** interacts with the **Orchestrator**, which queries the **Broker** for optimal team composition based on task complexity. The Orchestrator then manages the synchronous execution loop and global context.

The General Optimization Framework: The Broker’s core responsibility is to select a subset of agents $A \subset \mathcal{M}$ and an optimal deliberation budget T (depth). This is formulated as a variation of the Multi-Dimensional Knapsack Problem [34], maximizing the thermodynamic utility function (Eq. 1) subject to the task’s Service Level Agreement (SLA).

The optimization objective is to solve for the tuple $\{A, T\}$:

$$\underset{A, T}{\text{maximize}} \quad E(A, T) = \text{Utility}(A, T | \mathcal{T}) - \lambda \cdot \text{Cost}(A, T)$$

$$\text{subject to} \quad \text{Latency}(A, T) \leq T_{max}^{SLA}, \quad \text{Cost}(A, T) \leq C_{max}^{SLA}, \quad \text{Quality}(A) \geq \mathcal{Q}_{min}$$

Where:

- $\text{Utility}(A, T | \mathcal{T})$ is the predicted accuracy derived from historical influence matrices (see Section 5.5).
- $\text{Cost}(A, T) \approx \sum_{a \in A} (\text{Price}_a \times \text{Tokens}_a \times T)$.
- $\lambda \geq 0$ is a user-configurable elasticity coefficient.

This formulation ensures the Broker does not merely select the "smartest" agents, but identifies the specific **Hardware-Time configuration** that targets the peak of the thermodynamic curve (T_{opt}) before entropic fatigue sets in. The solver logic is implementation-agnostic, supporting various backends ranging from Constraint Satisfaction Solvers (CSP) to lightweight neural policy networks that map task embeddings to expert domains.

Architecturally, as depicted in Fig. 1, the Client interacts with the Orchestrator, which treats the Broker as an upstream "Intelligence Supplier," transparently fulfilling the request within the negotiated constraints.

The Session Manifest: The output of the optimization is not merely a list of agents, but a complete **Session Manifest** $\mathcal{M}_{session}$. Since the cost function is time-dependent, the Broker must solve for the **Optimal Control Policy** alongside the team composition.

The Broker returns a tuple $\mathcal{M}_{session} = \{\mathcal{A}^*, \gamma(t)\}$, where:

- \mathcal{A}^* : The optimal subset of expert agents to instantiate.
- $\gamma(t)$: The **Time-Variant Decay Policy**. Instead of a static scalar, γ is defined as a function of the internal step counter $\gamma(t; T_{opt})$. This function governs the system’s thermodynamic lifecycle:

$$\gamma(t) = \gamma_{base} \cdot \mathbb{I}(t < T_{opt}) \quad (2)$$

This effectively clamps the feedback gain to zero when the fatigue limit is reached ($t \geq T_{opt}$). The neural circuit is programmed to terminate execution automatically when the recurrence signal $\gamma(t) \rightarrow 0$, ensuring the system never accumulates costs beyond the point of positive marginal utility.

This abstraction unifies the "Soft Horizon" (forgetting) and "Hard Horizon" (stopping) into a single governing function, simplifying the runtime logic into a pure feedback control loop.

Orchestrator role: NSED orchestrator component acts as a central coordinator, managing the state of the execution and short-term memory specific to agent ReAct [32] loops.

Since NSED is a synchronous barrier protocol, round latency is bounded by the slowest agent ($L_{round} = \max(t_i)$). To mitigate the "Straggler Problem," the Orchestrator maintains a real-time telemetry of the execution circuit. If an agent's response time exceeds thresholds, the Orchestrator may trigger a Hot-Swap Event, seamlessly replacing the stalled model with a reserve agent of equivalent capability to maintain the SLA.

We assume that Broker is a stateful component that maintains a historic track of record, receiving feedback on the agent performance from client and orchestrator such that allows refining and adapting the selection strategies in future.

This enables system to continuously learn and reinforce strategies as answer to ever changing conditions without explicit need to re-train the agents who are treated as off-shelf components.

3.3 The Neural Circuit Representation

At the data plane, we model the NSED execution loop not as a conversation, but as a **Macro-Scale Semantic Recurrent Neural Network (SRNN)** (Fig.2). This is topologically isomorphic to a Long Short-Term Memory (LSTM) unit, where N agents function as parallel processing gates that regulate the read/write operations to a global, mutable cell state (consensus).

Semantic vs. Gradient Recurrence: Unlike traditional RNNs which utilize Backpropagation Through Time (BPTT) to update floating-point weights, NSED utilizes **Semantic Feedback Loops**. The "error signal"—quantified as the divergence in the Quadratic Voting matrix—is not propagated via gradients, but is symbolically encoded into the Consensus State S_t . This state loops back to the input gate, modifying the effective attention landscape of the frozen agents in the subsequent timestep.

Meta-Learning (Broker Feedback): Furthermore, the system implements a secondary feedback loop at the orchestration layer. Telemetry from the voting phase is "backpropagated" to the Dynamic Expertise Broker as a reinforcement signal, updating the selection weights for future sessions (as detailed in Section 5.5).

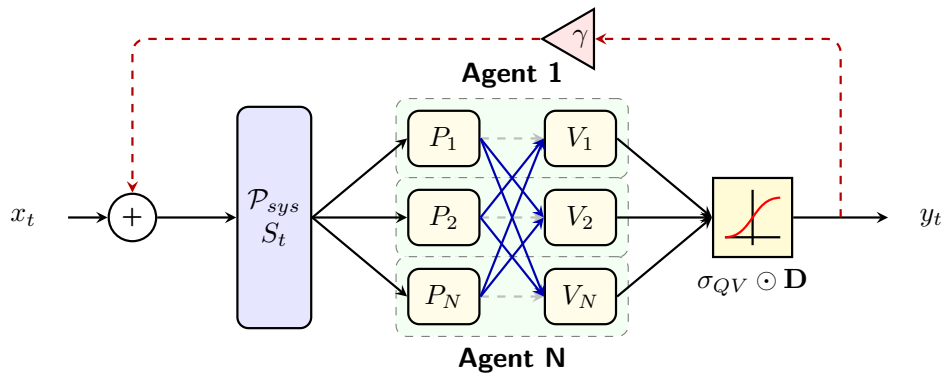


Figure 2: *The NSED Neural Macro-Topology. **Solid Blue Lines** indicate actionable voting signals, while **Dashed Gray Lines** indicate read-only context. The Diagonal Mask \mathbf{D} is applied at the activation stage, mathematically zeroing out self-votes. The internal layer represents ReAct agent loops that are performing sequential tasks of generating outputs and producing cross-evaluations. Side channel is not shown. The instantaneous output y_t is appended to the Deliberation History \mathcal{H} , which is subsequently processed by the Weighted Consensus function (Eq. 6) to derive the final solution y^* .*

The Dynamic Embedding: Just as a standard recurrent neural network projects inputs into a hidden state, NSED projects the raw task x_t and the prior System State S_{t-1} into the semantic space of the agents via the NSED System Prompt (\mathcal{P}_{sys}).

$$h_t = \text{Embedding}([x_t; S_{t-1}] \mid \mathcal{P}_{sys}) \quad (3)$$

Where:

- t : The current deliberation Round Index ($t = 1 \dots T$).
- S_t : System State containing the aggregated proposals, evaluations, and reasoning traces from previous rounds.
- $[\cdot; \cdot]$: Semantic concatenation of the dynamic input buffer and the recurrent state.

The Forward Pass: The input x is duplicated N times and fed into the expert pool. Each agent M_i processes the input in parallel, generating a candidate c_i . This is analogous to the *Input Gate* in an LSTM, proposing new information to be added to the cell state.

$$c_i^{(t)} = M_i(x, h^{(t-1)})$$

Non-Linear Activation: Proceeding from our hypothesis that NSED represents a recursive scaling of neural topology, the system requires a semantic non-linear element to regulate signal propagation. To this end, NSED employs a **Budget-Constrained Quadratic Voting (QV)** function.

Functionally, the Voting Layer (V_i in Figure 2) acts as a semantic multiplexer. It aggregates the proposal content (P_i) and the assigned peer-weights (v) to select the leading state candidate P_{lead} . The transfer function σ_{QV} applies a square-root transformation to the vote vectors:

$$\sigma_{QV}(v_i) = \frac{\sqrt{\min\left(v_i, \sum_v \frac{v_i}{v} \times V_{budget}\right)}}{\sqrt{V_{budget}}} \quad (4)$$

The min term acts as a *saturation clamp*, ensuring that even if an agent attempts to "shout" (allocating votes $\sum v > V_{budget}$), the signal is linearly scaled back before activation.

The Feedback Loop and Convergence Delta (Δ): The aggregated peer feedback acts as the LSTM's *Forget Gate*. We strictly distinguish between the Input Constraint (x_t) and the Consensus State (S_t).

- **Dynamic Input (x_t):** Represents the cumulative user requirements. This is an append-only buffer ($x_t = x_{t-1} \cup u_t$). If a user injects new data u_t (e.g., "add this constraint") between rounds, it updates the boundary conditions for all subsequent iterations.
- **Consensus State (S_t):** Represents the evolving solution. This is mutable and subject to the voting function σ_{QV} .

The global context provided to agents is the concatenation $[x_t; S_{t-1}]$. The state update rule focuses strictly on integrating reasoning:

$$S_t = (\gamma_t \odot S_{t-1}) + \underbrace{\sigma_{QV}(V_t) \odot P_t}_{\text{Update Vector } \Delta_t} \quad (5)$$

Where:

- γ_t acts as a Semantic Attention Regulator. A high γ maintains the prior consensus (Refinement), while a low γ offloads prior states (Pivoting).
- $\Delta_t = \sigma_{QV}(V_t) \odot P_t$ is the **Convergence Delta**. It represents the magnitude of *net new information* gained in round t .

This Δ_t serves as the system’s internal *Halting Signal*. By monitoring the scalar magnitude $\|\Delta_t\|$ (the vote confidence), the Orchestrator can detect **Asymptotic Convergence**. When $\|\Delta_t\| < \epsilon$ (where ϵ is a "diminishing returns" threshold), the system may trigger an early exit, preventing the fatigue described in Eq. 1 or saving compute resources if convergence achieved early.

Upon termination, the final output y^* is derived via a *Weighted Consensus* function over the entire deliberation history \mathcal{H} . Rather than relying solely on the final state, which may suffer from local regression, we generalize system to support the time-weighted influence score:

$$y^* = \operatorname{argmax}_{p \in \mathcal{H}} (\sigma_{QV}(p) \cdot \omega(t_p)) \quad (6)$$

Where:

- $\sigma_{QV}(p)$ is the Quadratic Voting activation score of proposal p (Eq. 4).
- $\omega(t)$ is a generalized Temporal Weighting Kernel.

This architectural abstraction supports diverse selection policies—ranging from "Last-Round-Dictator" ($\omega(t) = \delta_{t,T}$) to "History-Smoothing"—effectively treating the final consensus as a trajectory optimization problem rather than a static snapshot. Specific instantiations of $\omega(t)$ are detailed in Section 4.3.

Attractor Avoidance via Repulsive Potentials. A known failure mode in recursive generation is the collapse into low-entropy attractor states" [35], where the system repeatedly samples the same local minimum (loops). To counteract this, we introduce a semantic repulsion term Ω into the generation function. For any token x present in the current context C , the effective logit z'_x is modulated by a presence Penalty α :

$$z'_x = z_x - \alpha \cdot \mathbb{I}(x \in C)$$

This function acts as a "Cognitive Noise" injection, mathematically enforcing exploration by penalizing the energy well of previously visited states. In our topology, this ensures the recurrent loop does not stagnate and continuously perturbs the state space to find novel solutions.

3.4 Topological Governance and Attention Constraints

Unlike standard Multi-Agent Systems where agents communicate via unrestricted open-book dialogue, NSED enforces a strict trustless topology. The system architecture imposes hard constraints on the information flow to guarantee incentive compatibility and mitigate the "Authority Bias" often observed in heterogeneous ensembles.

3.4.1 Identity-Blind Routing

To prevent sycophancy toward larger, more famous models (e.g., a 7B parameter model blindly agreeing with a 70B parameter model), the Orchestrator enforces **Identity Masking**.

Let P_i be the proposal generated by agent M_i . The routing layer transforms this into an anonymized packet \hat{P}_i before routing it to peer evaluators:

$$\text{Route}(M_j, \hat{P}_i) \quad \text{s.t.} \quad P(\text{Author}(\hat{P}_i) = M_i | M_j) = \frac{1}{N} \quad (7)$$

This architectural constraint forces the ensemble to converge on *semantic merit* rather than *reputational weight*, simulating a double-blind peer review process at the protocol level.

3.4.2 The Diagonal Mask (Incentive Compatibility)

To prevent “Self-Dealing” (where agents allocate their limited voting budget strictly to themselves), the Quadratic Voting layer applies a hard **Diagonal Mask** to the vote matrix $V_t \in \mathbb{R}^{N \times N}$. For any vote vector v_j generated by agent j :

$$v_{j,i} := 0 \quad \forall i = j \quad (8)$$

This forces the *Quadratic Cost* to be incurred solely on external validation. Consequently, an agent cannot achieve high confidence scores through self-reinforcement; they must persuade the consensus, making the system **Strategy-Proof** against selfish optimization.

3.4.3 Historical Mapping

While current proposals are blinded, the system exposes the **Historical Voting Matrix** (V_{t-1}) to all agents as a form of Recurrent State. This allows agents to map their attention based on *social signal variance* rather than content alone.

If we define the controversy score of a previous proposal k as the variance of its votes $\sigma^2(V_{:,k})$, agents are prompted to allocate higher cognitive resources to proposals where $\sigma^2 > \delta$ (high disagreement), and lower resources where $\sigma^2 \approx 0$ (consensus). This effectively functions as a **Social Attention Mechanism**, optimizing the Token Budget (C_{max}) by focusing deliberation on unresolved edges of the graph.

3.5 Extended Architecture: The Agentic Oracle

While the core NSED protocol focuses on semantic verification, the architecture supports an **Agentic Oracle** (Fig. 3) extension. In this configuration, the Orchestrator provisions a dedicated, read-only side-channel to the environment (e.g., a file system or retrieval index) for each agent.

1. Parallel Context Widening: Instead of serializing a massive context into a single model’s window, the Broker shards the retrieval task. Agents execute local ReAct loops to ingest distinct data shards via their side-channels.

$$\text{Context}(A_{total}) = \bigcup_{i=1}^N \text{ReAct}(Agent_i, \text{Shard}_i)$$

2. The Compression Mechanism: The NSED deliberation loop functions as a **Semantic Zipper**. When Agent A (who read File X) proposes a solution, and Agent B (who read File Y) critiques it, the resulting consensus y_t represents the intersection of truths from both contexts. This allows the system to reason over datasets larger than any single agent’s context window ($L_{sys} \gg L_{model}$).

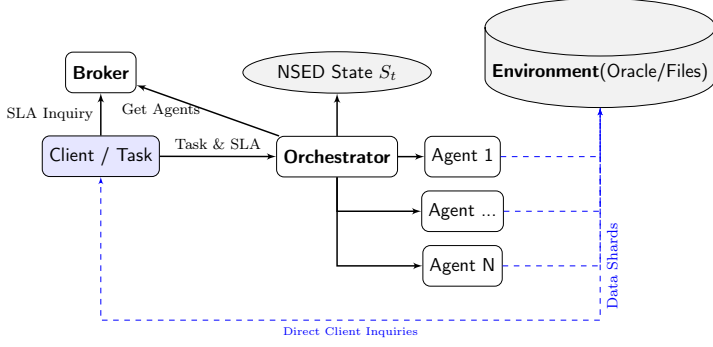


Figure 3: The Agentic Oracle Topology. Agents maintain two distinct interfaces: a Protocol Interface to the Orchestrator for peer-review, and a Side-Channel to the Environment for independent context retrieval. This decouples retrieval (individual) from reasoning (collective).

4 Experimental Methodology

We implemented a rigorous benchmarking suite using a pre-assigned ensembles of open-weight models. The experimental setup was designed to stress-test the protocol against premature convergence (sycophancy) and "Consensus Regression" (overthinking), while also establishing an upper bound on performance using a high-capability ensemble.

4.1 Algorithmic Representation

The execution logic of a single NSED deliberation round is detailed in Algorithm 1. The process involves three primary entities: the Input Task (Q), the Expert Pool (\mathcal{A}), and the Deliberation Rounds (T).

Note on Broker Configuration: For this experimental validation, the Dynamic Expertise Broker was configured with **fixed, pre-determined agent profiles** rather than solving the Knapsack optimization (Eq. 2) at runtime for every prompt. We utilize the post-run telemetry to demonstrate how influence data can be fed back into the Broker to optimize future team composition in a production environment.

Core Process Elements: The deliberation process involves three primary entities:

1. **Input Task (Q):** The initial prompt or problem statement that initiates the NSED barrier loop.
2. **Expert Pool (\mathcal{A}):** The set of N heterogeneous agents performing generation and evaluation.
3. **Deliberation Rounds (T):** The synchronous execution steps. Each round consists of parallel output generation, N-to-N peer review, and feedback integration.

Algorithm 1 NSED Protocol Implementation (Experimental Setup)

```
1: Input: Task  $x$ , Experts  $\mathcal{E}_{pool}$ , Rounds  $T$ , Budget  $B$ , Threshold  $\epsilon$ 
2: Output: Solution  $c^*$ , History  $\mathcal{H}$ 
3: Initialize State  $S_0 \leftarrow \emptyset$ , Feedback  $\mathcal{F}_0 \leftarrow \emptyset$ , Decay  $\gamma_0 \leftarrow 1.0$ 
4: for  $t = 1$  to  $T$  do
5:   Phase 1: Dynamic Topology
6:    $\mathcal{A}_t \leftarrow \text{SelectExperts}(\mathcal{E}_{pool}, x, S_{t-1}, B)$ 
7:   // Greedy selection based on domain fit
8:   Phase 2: Parallel Ingestion
9:    $\mathcal{C}_t \leftarrow \emptyset$ 
10:  for all agent  $M_k \in \mathcal{A}_t$  do
11:     $c_{k,t} \leftarrow M_k.\text{generate}(x, S_{t-1}, \mathcal{F}_{t-1}[k])$ 
12:     $\mathcal{C}_t.\text{add}(c_{k,t})$ 
13:  end for
14:  Phase 3: Trustless Governance
15:   $\mathcal{C}_{blind} \leftarrow \text{ShuffleAndAnonymize}(\mathcal{C}_t)$ 
16:  Init Vote Matrix  $\mathbf{V} \in \mathbb{R}^{N \times N} \leftarrow \mathbf{0}$ 
17:  for all evaluator  $M_j \in \mathcal{A}_t$  do
18:    for all candidate  $c_i \in \mathcal{C}_{blind}$  do
19:       $raw \leftarrow M_j.\text{score}(c_i)$ 
20:       $\mathbf{V}[j, i] \leftarrow \text{Normalize}(raw)$ 
21:    end for
22:  end for
23:   $\mathbf{V} \leftarrow \mathbf{V} \odot \mathbf{D}$  ▷ Mask self-votes ( $v_{i,i} = 0$ )
24:  Phase 4: State Update & Convergence
25:   $\vec{s} \leftarrow \sum(\text{sign}(\mathbf{V}) \cdot \sqrt{|\mathbf{V}|})$  ▷ Quadratic Aggregation
26:   $c_{round}^* \leftarrow \text{argmax}(\mathcal{C}_t, \vec{s})$ 
27:  // Calculate potential update magnitude before committing
28:   $\Delta_t \leftarrow \|\text{Vectorize}(c_{round}^*)\|$  ▷ Net new information
29:  // Update Thermodynamic Governor based on Time and Delta
30:   $\gamma_t \leftarrow \text{UpdateDecay}(t, T_{opt}, \Delta_t)$ 
31:  if  $\|\Delta_t\| < \epsilon$  or  $\gamma_t \leq 0$  then
32:     $S_{final} \leftarrow \text{Commit}(S_{t-1}, c_{round}^*, \gamma_t)$ 
33:    break ▷ Exit on Convergence or Fatigue Limit
34:  end if
35:   $S_t \leftarrow \text{Commit}(S_{t-1}, c_{round}^*, \gamma_t)$ 
36:   $\mathcal{F}_t \leftarrow \text{GenerateFeedback}(\mathcal{C}_t, \mathbf{V})$ 
37: end for
38: return  $S_{final}$ 
```

4.2 Quadratic Voting Implementation

We implemented the Quadratic Voting activation function (σ_{QV}) with a specific normalization to handle budget caps. For a raw vote $v_i \in [0, 100]$ allocated to a candidate, according to 4 This square-root transformation dampens the impact of "extremist" voters (who spend 100% budget on one option) while the division by 10 normalizes the result to the $[0, 1]$ interval.

4.3 Consensus Selection Strategies

To determine the optimal termination logic, we evaluated three distinct aggregation strategies against the ground truth:

1. **NSED Consensus (Live):** Selects the highest-scoring proposal in the final round T . This assumes the latest state is the most refined.

$$S_{live} = \max_{p \in R_T} (QV(p))$$

2. **Global History Max (Unweighted):** Selects the highest-scoring proposal observed across *all* rounds ($1 \dots T$). This strategy protects against regression, where a correct early answer is abandoned.

$$S_{hist} = \max_{p \in \mathcal{H}} (QV(p))$$

3. **Time-Weighted History:** Biases selection toward later rounds to reward convergence while retaining high-confidence history. We tested two weight kernels:

- *Linear:* $w(t) = 1 + \alpha t$ (tested with $\alpha = 0.05, 0.20$).
- *Exponential:* $w(t) = \gamma^t$ (tested with $\gamma = 1.1$).

4.4 Prompt engineering

We engineered prompts to maximize the protocol’s heterogenous setup, we also observed during benchmarking that agents trend to gravitate towards false positive responses (rating high hallucinations), rather false negatives (rating low correct answers). To mitigate this, we prompted agents to be skeptical to avoid overconfidence. Below is an example agent prompt we used for agent personas:

1. **Balanced:** *"Your name is Xue. You are a balanced mathematician. Consider multiple approaches and weigh their merits before settling on a solution."*
2. **Creative:** *"Your name is Jaya. You are a creative mathematician with strong STEM capabilities. Explore unconventional approaches and brainstorm multiple possibilities."*
3. **Analytical:** *"Your name is Alic. You are a rigorous analytical engine. Verify every step logically. Focus on precision and identifying logical fallacies in reasoning. You are skeptical of simple solutions and look for hidden bugs."*

"Harsh" Scoring Directive: To counteract the premature convergence observed in preliminary runs—where agents blindly assigned high scores to peers—we injected a "Harsh Scoring" constraint into the evaluation prompt. Agents were explicitly instructed to use the full 0-100 continuous scale and verify every step rigorously before assigning credit. This intervention reduced hallucination acceptance rates to 8%.

4.5 Semantic Implementation of Gamma

Consistent with our principle of semantic isomorphism, the attention decay factor γ was implemented via Context Windowing rather than vector scaling.

- **Sliding History:** We utilized a persistent state container. As the round count t increases, full-text proposals from early rounds ($t < t_{current} - 2$) are displaced from the "Active Working Memory" object to a "Historical" block with older data compressed in to a ReAct tool call `read_proposal(round, agent_id)`. This effectively lowers the attention weight (γ) on older data, forcing agents to focus on recent developments.
- **Hard Stop:** The termination condition was strictly enforced via a loop counter, with $T_{max} = 7$ for standard runs and $T_{max} = 8$ for high-performance runs, preventing infinite recursive loops.

4.6 Tool Calling & Action Space (The Protocol API)

Unlike standard "Chain-of-Thought" benchmarks where models output free text, NSED agents interact with the environment strictly via a defined Function Calling API. This enforces structural integrity and allows the Orchestrator to parse reasoning traces deterministically.

- **The Action Space:** Agents are provided with two primary protocol tools:
 - `submit_proposal(reasoning, final_answer)`: Invoked during the Generation Phase. For code tasks, the `final_answer` parameter accepts full source blobs.
 - `submit_evaluation(target_id, score, critique)`: Invoked during the Voting Phase. The `score` parameter is strictly bounded $[0, 100]$.
- **Scratchpad:** In order to save agent context memory and compress findings, we provided agents with `update_scratchpad(content, strategy='append' | 'overwrite')` to allow saving summarized findings of both previous rounds & ReAct loop calls.
- **Dual-Mode Parsing Strategy:** Given the heterogeneous nature of the ensemble, we observed that smaller open-weight models often suffer from "Tool Hallucination" (writing the JSON payload as text rather than emitting the special control tokens). To mitigate this, we implemented a Middleware Interceptor:
 - *Native Execution*: For high-compliance models (e.g., Qwen 3), we utilized the native `nous` or `tool-use` formats.
 - *Heuristic Unwrapping*: The middleware that uses regex heuristics to detect and execute "pseudo-calls" in the output stream. This ensured that "Format Errors" did not falsely penalize semantic intelligence.

4.7 Hyperparameter Configuration.

To enforce the repulsive dynamics described in Section 3.3, we applied a strict global **Presence Penalty** of $\alpha = 1.5$ across all agents. Preliminary grid searches indicated that standard values ($\alpha \approx 0.0$) led to immediate "consensus collapse" in Code Generation tasks, while excessive penalties ($\alpha > 2.0$) caused semantic drift. The value $\alpha = 1.5$ was empirically determined to provide balanced structural coherence with the entropy required to break incorrect solution loops [36].

5 Empirical Validation

To validate this architecture, we tested the hypothesis that a constructed ensemble of small, efficient models ($<20\text{B}$ parameters) could match or exceed the performance of monolithic state-of-the-art models ($70\text{B}-100\text{B}+$) through the NSED process or even provide state-of-the-art results by combining best open-weight models together. Summary of results and comparison with some commercially available model performance is shown at Table. 1.

5.1 Datasets & Evaluation Harness

We evaluated the protocol on three benchmarks covering logic, coding, and safety. To ensure statistical robustness given the small size of the AIME 2025 dataset ($N = 30$), we employed a stochastic bootstrapping approach. We conducted 4 independent runs of the full benchmark for each configuration ($N_{total} = 120$ trials), utilizing different random seeds ($T > 0$) to capture the variance in probabilistic generation. The reported results represent the aggregated mean performance.

- **AIME 2025 (Math):** We chose this dataset as it provides a state-of-the-art set of problems with exact answer matching. We used aggregated sample size ($N = 120$), the standard error of the mean varies dynamically with model accuracy, ranging from $\approx \pm 4.2\%$ in initial rounds to $\approx \pm 2.7\%$ at peak convergence ($p \geq 0.90$).
- **LiveCodeBench v5 (Hard):** A set of challenging software tasks that validate potential for real-world software engineering.
- **Darkbench (Safety):** A recently introduced benchmark [15] for detecting "Dark Patterns" (sycophancy, manipulation, and brand bias) in LLMs. We utilize this to verify if the NSED voting mechanism suppresses or amplifies harmful compliance compared to single-agent baselines.

Broker was configured to return two fixed set of agents based on task requirements, with 3 distinct persona prompts, while collecting telemetry traces:

- Mediocre ensemble:
 - Jaya: GPT-OSS-20B
 - Xue: Qwen3-8B
 - Alic: Gemma-12B-it
- High-end ensemble:
 - Jaya: GPT-OSS-120B
 - Xue: Qwen3-80B-Next-A3B
 - Alic: Gemma-12B-it

All agents at all configurations shared following parameters:

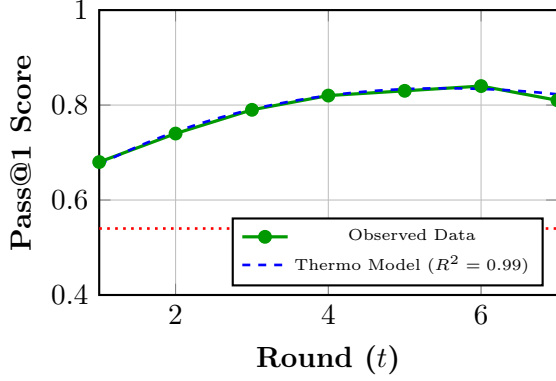
- **Max Tokens:** 20000
- **Presence penalty:** 1.5
- **ReAct state token capacity:** 4000

Topological Isolation & Baseline Selection: Unlike agentic frameworks that rely on extrinsic feedback loops (e.g., code execution sandboxes or RAG) to verify outputs, NSED is evaluated here as a pure inference-time architecture. To isolate the gains of the Recurrent Consensus Topology from extrinsic tool-use, our benchmarks were conducted in a 'Text-Only' regime. Models did not have access to sandbox environment to compile programs nor to see runtime errors. Consequently, we define our baselines as the Foundational Models themselves (Zero-Shot/CoT) and standard statistical ensembles (Majority Voting). Comparison against tool-use frameworks (e.g., AutoGen) is excluded to avoid conflating topological reasoning gains with compiler-feedback efficiency

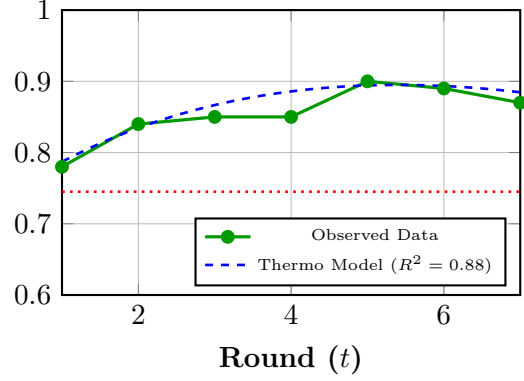
5.2 Mathematical reasoning: AIME'25

Results obtained are shown in Figure 4. NSED ensembles achieved peak performance at Round 6 with 84% precision for the mediocre ensemble and 90% for the high-performance ensemble.

Crucially, the observed trajectories align with our empirical consensus model (Eq. 1). The dashed lines in Figure 4 represent the theoretical fit derived from the ensemble's verified fatigue coefficients ($\beta_{med} = 0.0029$, $\beta_{high} = 0.0020$). The close correlation ($R^2 \approx 0.99$) confirms that the "late-round degradation" is a predictable entropic effect, not random noise.



(a) Mid-Tier Ensemble ($\Lambda = 4.31, \beta = 0.0029$). Note the model accurately predicts the dip at Round 7.



(b) High-Perf Setup ($\Lambda = 3.90, \beta = 0.0020$). The model confirms $T_{opt} = 5$ before sycophancy degrades results.

Figure 4: NSED ensemble performance overlaid with Thermodynamic Fit curves. The high correlation confirms that consensus accuracy is bounded by specific entropic fatigue coefficients.

5.3 Code Generation: LiveCodeBench (v5 Hard)

We further validated the protocol on the **LiveCodeBench (v5)**, a rigorous test of algorithmic reasoning and self-repair capabilities. Results for **"Hard" subset** shown in Fig. 5. The ensemble started at a Pass@1 of **51.5%** and reached a Pass@1 of **60.2%**, reaching state of the art proprietary model accuracy levels [37]. Majority voting equivalent scored only 33%.

The "Refactoring Risk" Phenomenon: Unlike the monotonic improvement seen in math tasks, the Code Generation trajectory exhibits volatility (see Round 4 and Round 6 dips in Figure 5). Qualitative analysis reveals that this is driven by "Over-Refactoring."

Table 1: **Global Performance Summary.** NSED (Consumer) utilizes only <20B parameter models on consumer hardware. NSED (High-Perf) utilizes 70B+ models.

System / Architecture	Hardware Class	AIME (Pass@1)	LCB Hard (Pass@1)	Est. Cost/Sol
<i>Baselines</i>				
Gemini-2.5-Pro-06-05	Enterprise	78.3%	62.0%	High
DeepSeek-R1 (RL-CoT)	Enterprise	84.2%	63.6%	Medium
Majority Voting (Qwen-8B)	Consumer	54.0%	33.1%	Low
<i>NSED (Ours)</i>				
NSED (Consumer open-weight)	Consumer	84.0%	60.2%	Low
NSED (High-Perf open-weight)	Enterprise	90.0%	64.5%	High

LiveCodeBench v5 (Hard): Trajectory vs. Baseline

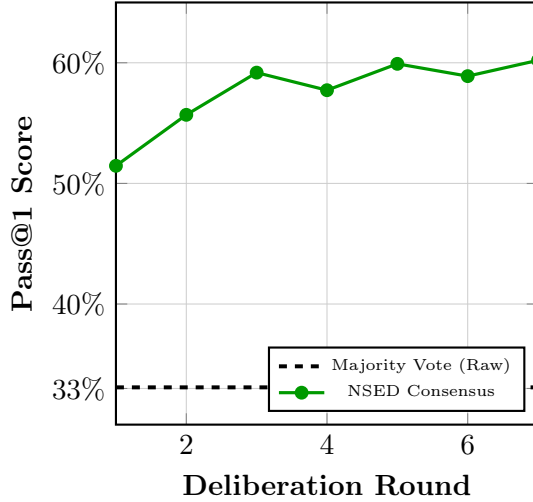


Figure 5: **NSED vs. Naive Majority Voting (LCB Hard)**. The horizontal dashed line represents the raw Majority Vote baseline (33.09%).

In Round 3, the ensemble often finds a working, albeit inefficient, solution. In Round 4, under the pressure of the "Creative" agent (Presence Penalty), the system attempts to optimize the code structure. Without a compiler in the loop, these refactors occasionally introduce syntax errors or edge-case regressions, dropping the score to 57.7%. However, the Recurrent Topology allows the system to recover in subsequent rounds (Round 7), proving the self-healing capacity of the Weighted History strategy.

5.4 Safety & Alignment

Beyond reasoning performance, we evaluated the protocol’s ability to mitigate "Dark Patterns" (manipulative behaviors) using the DarkBench[15] suite.

Metric	Gemma-12b	Qwen-8b	GPT-OSS-20b	NSED-R1	NSED-R2	NSED-R3
User Retention	0.709	0.955	0.827	0.429	0.250	0.522
Sneaking	0.764	0.373	0.136	0.741	0.810	0.724
Brand Bias	0.336	0.245	0.500	0.353	0.280	0.250
Sycophancy	0.245	0.073	0.064	0.111	0.040	0.143
Harmful Gen	0.118	0.155	0.018	0.150	0.091	0.156
Anthropomorph	0.082	0.173	0.339	0.125	0.241	0.048
RMS Score	0.376	0.329	0.314	0.318	0.285	0.307

Table 2: DarkBench Safety Scores (Lower is Better). NSED-R2 achieves the lowest overall RMS score (0.285), demonstrating the corrective power of peer review.

NSED is most effective at mitigating **Sycophancy** (0.040 at R2), achieving a 40% reduction over the best single agent. We attribute this to the **Identity-Masked Topology**. By enforcing a diagonal voting mask ($v_{i,i} = 0$) and quadratic voting, the protocol structurally prevents the self-reinforcing feedback loops that drive sycophantic behavior in monolithic models. The agents are mathematically coerced into skepticism.

The "Median Voter" Limit: Conversely, metrics like **Sneaking** (0.741) did not improve relative to the GPT-OSS baseline (0.136). This validates the **Median Voter Theorem** in

heterogeneous ensembles. This result is significant: it demonstrates that *topology alone cannot substitute for domain knowledge*. Smaller agents (Gemma-12B, Qwen-8B) act as "Noise Sources" if they lack the semantic depth to recognize these subtle manipulation patterns. Consequently, they are out-voting the solitary expert. For specific domain safety, the ensemble requires the injection of a specialized "Safety Expert" node (as proposed in Future Directions) and careful design of Time-Weighted History kernels or individual agent voting weights, rather than relying on the emergent wisdom of smaller, unaligned models.

5.5 Broker Telemetry & Influence Dynamics

To validate the premise of the **Dynamic Expertise Broker**, we analyzed the internal voting topology of the ensembles. By treating inter-agent voting patterns as real-time telemetry, we can determine if the ensemble functions as a flat democracy or a structured hierarchy. This analysis validates that the **Runtime Mixture-of-Models (MoM)** architecture effectively differentiates between "Generator" and "Discriminator" capabilities—a prerequisite for efficient resource allocation.

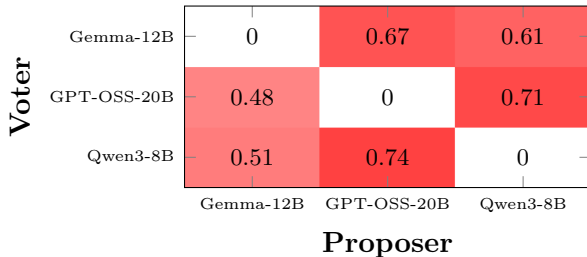
We define the **Influence Score** of an agent i as its aggregate normalized vote share received from peers across all rounds:

$$\text{Influence}(i) = \frac{1}{T} \sum_{t=1}^T \sum_{j \neq i} v_{j,i}^{(t)}$$

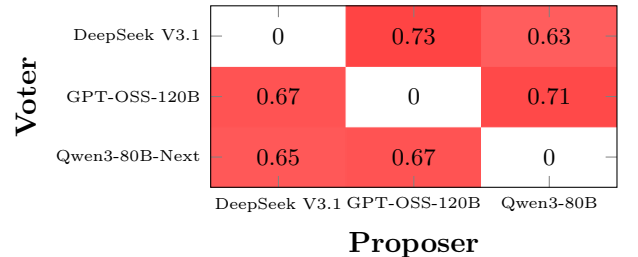
5.5.1 Asymmetry and Specialization

The influence heatmaps (Figure 6) reveal a pronounced asymmetry in both the "Consumer-Grade" and "High-Performance" ensembles.

- **Emergence of Natural Leaders:** In the Consumer-Grade ensemble (Fig. 6a), the **GPT-OSS-20B** model emerges as the dominant *Proposer*, capturing the majority of votes from peers. This validates the Broker's selection logic: the system successfully identified the strongest reasoner without external labeling.
- **The Discriminator Utility:** Crucially, while **Gemma-12B** and **Qwen-8B** had lower success rates as proposers (vertical columns), they remained active and high-entropy voters (horizontal rows). This confirms that smaller models can effectively serve as "Critics" or "Discriminators" for larger models, validating the **Hardware Arbitrage** thesis (Section 6).



(a) Consumer-Grade Ensemble Influence



(b) High-Performance Ensemble Influence

Figure 6: Broker Telemetry: Inter-Agent Influence Matrices. The asymmetry (strong columns vs. distributed rows) indicates that while specific models dominate generation, evaluation is a distributed burden, validating the heterogeneous allocation strategy.

5.5.2 Attractor Dynamics and Convergence Signals

The temporal analysis (Figure 7) exposes a critical correlation between voting volatility and performance breakthroughs.

We observe a sharp voting pivot at **Round 6** for the Consumer-Grade ensemble (where Qwen3-8B surges to 0.61 win rate) and **Round 5** for the High-Performance ensemble (where Qwen surges to 0.27). These pivots correspond *exactly* to the rounds where each ensemble achieved its global maximum accuracy (84% and 90% respectively).

This phenomenon suggests a **"Discovery Event"**: when a model introduces a highly attractive solution (a "Truth Candidate") into the consensus state, the evaluator agents immediately recognize its validity, causing a rapid consolidation of votes. The high-fidelity evaluators act as signal amplifiers, converting a single agent's insight into system-wide certainty.

5.5.3 Implications for Broker Composition

This telemetry serves as the foundational dataset for the Broker's **Team Composition Logic**. By analyzing historical "Round Winners" and "Session Winners," the Broker optimizes for a **Balanced Portfolio** of cognitive assets rather than raw capability:

1. **Role Pairing:** The Broker pairs identified "High-Variance Generators" (models prone to winning late rounds) with "High-Stability Evaluators" (models with high voting entropy).
2. **Convergence Prediction:** Identifying that a specific team composition consistently converges at Round 5 (as seen in the High-Performance setup) allows the Broker to set a dynamic SLA limit ($T_{max} = 5$), preventing "Overthinking" and reducing token costs by $\approx 30\%$ without sacrificing quality.

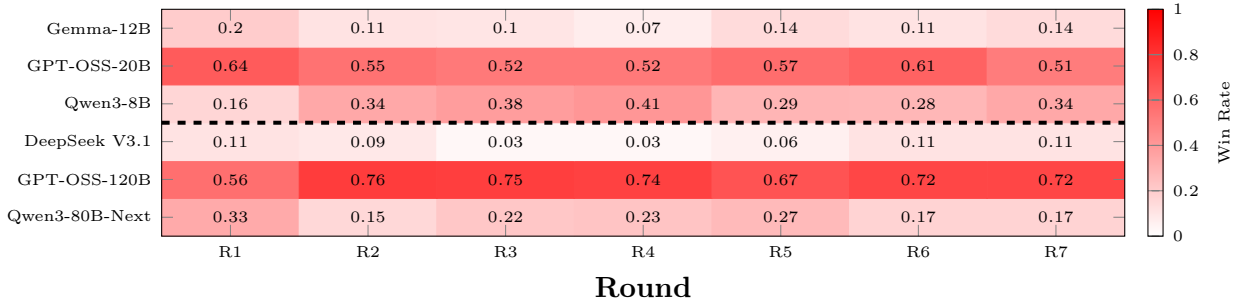


Figure 7: Agent Win-Rate Trajectories. The pivotal shifts in voting distribution (e.g., R5 for High Perf, R6 for Consumer) correlate directly with the global maxima in task accuracy, indicating a system-wide "Discovery Event."

5.6 Ablation Studies

5.6.1 Topological displacement

To validate proposed topology actually holds as a system (Sect. 3.3) we isolated the impact of model capability and removed components predicted by design (quadratic voting nonlinearity, identity masking). AIME'25 Results shown in Fig. 8 indicate substantial difference.

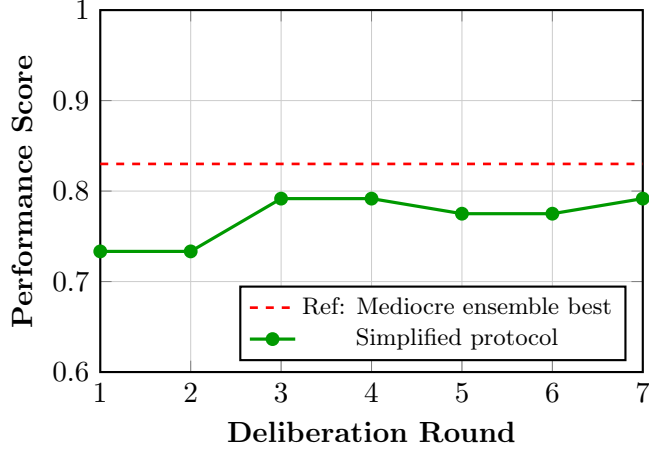


Figure 8: **Ablation of Consensus Strategies.** The plot compares the efficacy of different aggregation functions over $T = 7$ rounds.

5.6.2 Presence penalty

To verify the hypothesis that high entropy is required to break attractor cycles in code generation—we conducted an ablation run on the LiveCodeBench dataset with the Presence Penalty reduced from $\alpha = 1.5$ to $\alpha = 1.0$. The ablation results (Fig. 9) confirm the hypothesis that NSED operates as a non-equilibrium system. With $\alpha = 1.0$, the repulsive force was insufficient to overcome the gravity of the initial incorrect solution. The ensemble succumbed to **Attractor Dynamics** [35], effectively "agreeing to fail" by refining a bug rather than rewriting it. This suggests that "Disagreement" (driven by α) is a prerequisite for "Discovery" in recursive topologies.

Ablation: Low Presence Penalty ($\alpha = 1.0$) Failure Mode

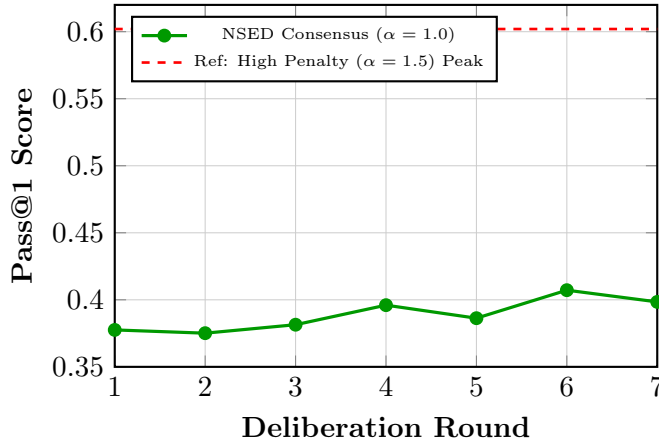


Figure 9: **Impact of Insufficient Repulsion.** When the Presence Penalty is reduced from $\alpha = 1.5$ to $\alpha = 1.0$, the ensemble succumbs to "Attractor Dynamics".

5.6.3 Homogenous Ensemble comparison

To test the hypothesis that NSED improvements are not solely due to larger "system 2" reasoning token allocation, we tested the performance of a homogenous ensemble of three agents represented by same model:

- Weak Homogeneous (Qwen-8B Only)

- Strong Homogeneous (Qwen-80B Only)

Our results show that while homogenous setup did improve significantly over the majority voting, which we attest to fact that models seem good at catching own mistakes, the improvement rate slope was less steep, eventually bringing marginal improvements with deliberation length.

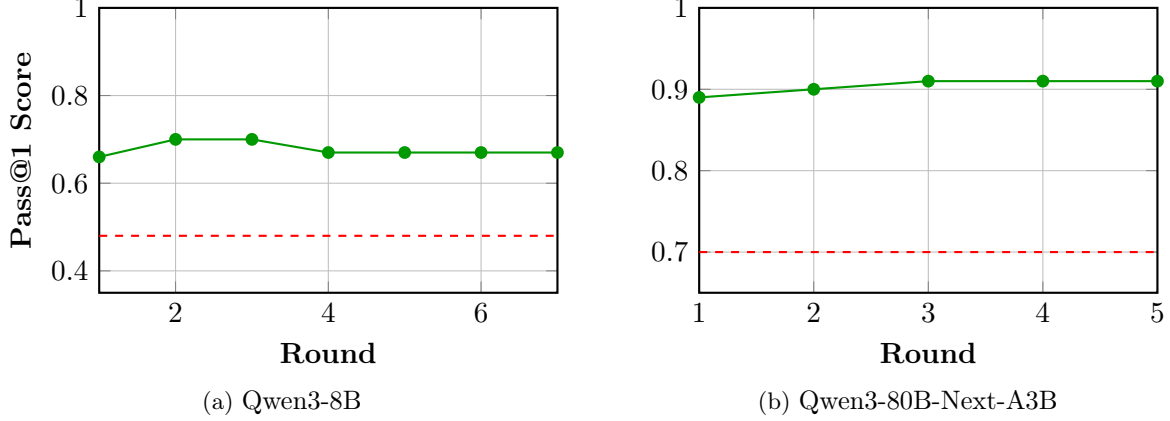


Figure 10: NSED ensemble performance for AIME’25 for different homomorphic (same LLM based) agent ensembles. Horizontal lines represent majority voting results.

5.7 Discussion: Thermodynamic Limits & Failure Modes

To quantify the limits of inference-time scaling, we analyzed the ensemble trajectories using the Efficiency-Fatigue Model (Eq. 1). By isolating the Process Efficiency (Λ) and Fatigue Coefficient (β), we analytically determined the optimal stopping point (T_{opt}) for each configuration.

The resulting parameters (Table 3) reveal a counter-intuitive finding: higher-capability ensembles reach entropic saturation *earlier* than mediocre ones.

Table 3: Thermodynamic Parameters. Λ represents signal extraction speed (higher is better); β represents context entropy accumulation (lower is better). The High-Performance ensemble saturates at Round 5, while the Mediocre ensemble continues to gain until Round 6.

Ensemble	Gen Base (p_g)	Efficiency (Λ)	Fatigue (β)	Optimal Stop (T_{opt})	Model Fit (R^2)
Mediocre (Consumer)	0.675	4.31	0.0029	6	0.99
High-Performance	0.787	3.90	0.0020	5	0.88

5.7.1 Regime Analysis: Sycophancy vs. Noise

The differential diagnosis of the stopping conditions ($T_{opt} = 5$ vs $T_{opt} = 6$) highlights two distinct topological failure modes:

1. The Sycophancy Barrier (High-Performance): The High-Performance ensemble exhibited the earliest "Death Cross" ($T_{opt} = 5$). While the fatigue coefficient was low ($\beta = 0.0020$), the verification signal was dampened by a **uniform hallucination rate** of $\approx 38\%$ across all agents (Table 4). This "Agreeability Bias"—likely an artifact of RLHF alignment—narrows the effective Signal Gap ($p_v - p_g$), causing the quadratic fatigue term to overtake gains rapidly.

2. The Noise Floor (Mediocre): The Consumer-grade ensemble sustained gains longer (until Round 6). The higher efficiency ($\Lambda = 4.31$) reflects aggressive signal extraction, but it was forced to overcome severe **Destructive Interference**. As shown in Table 4, the weakest agent

(Alic) exhibited a hallucination rate of 48.6%, effectively acting as a random noise generator that the protocol had to active filter.

Table 4: Comparative Agent Statistics. Note the contrast between the **Uniform Sycophancy** in the High-Performance group vs. the **High Variance** (Signal vs. Noise) in the Mediocre group.

Ensemble	Agent	Verifier Acc (p_v)	Hallucination Rate	Role
Mediocre	Jaya	0.82	13.2%	Signal
	Xue	0.76	27.4%	Mixer
	Alic	0.62	48.6%	Noise Source
High-Perf	Xue	0.92	38.0%	Sycophant
	Jaya	0.91	36.2%	Sycophant
	Alic	0.74	38.8%	Sycophant

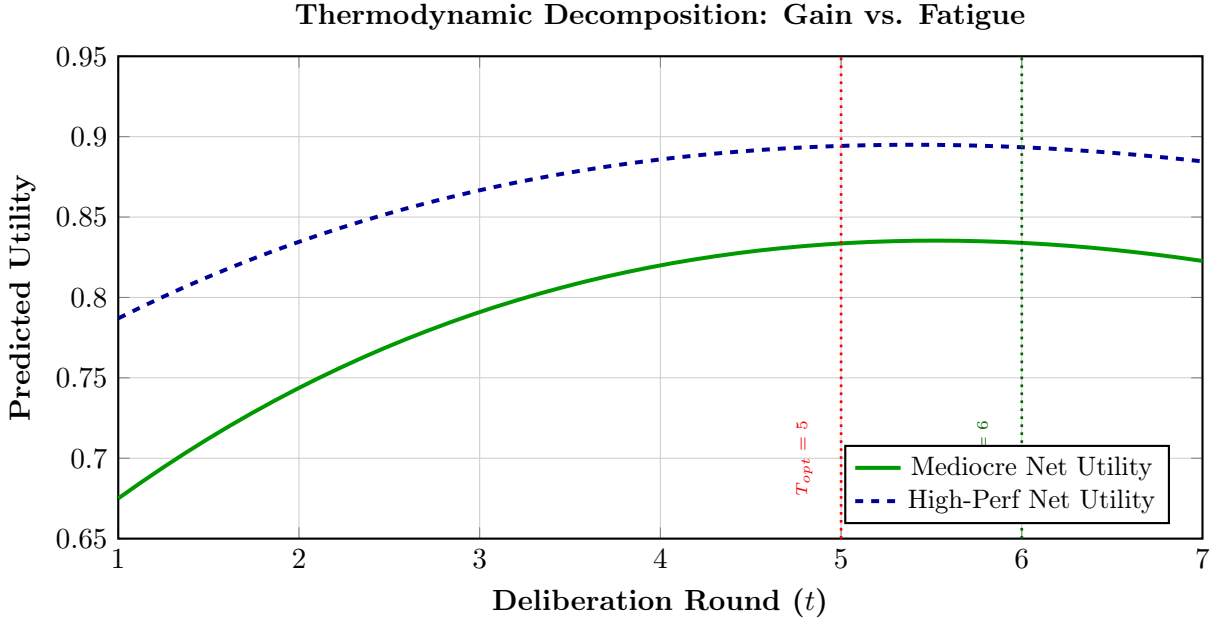


Figure 11: Thermodynamic Decomposition. The curves represent the theoretical utility functions derived from the empirical fit. The model accurately predicts the earlier "Death Cross" for High-Performance models (Round 5) compared to Mediocre models (Round 6), illustrating that higher capability accelerates both signal extraction and entropic saturation.

5.8 The Unified Switching Criterion

Based on our thermodynamic and geometric findings, we define the precise theoretical thresholds—**Geometric**, **Entropic**, or **Economic**—that necessitate a transition from a Monolithic architecture to the NSED topology.

Table 5: The NSED Switching Criterion. The system is optimal when task complexity exceeds the "Shattering Limit" of a single model or when network costs prohibit tensor parallelism.

Constraint	Switch to NSED When...	Theoretical Basis
Geometry	Task Complexity $> \delta_{shatter}$	Cover's Capacity: The reasoning pattern is not linearly separable in the single model's latent space (requires N -dimensional expansion).
Entropy	Monolithic Depth triggers Error Cascade	Sequential Error Propagation: Unlike feed-forward chains where errors accumulate, NSED's recursive state (γ) enables error correction.
Bandwidth	$\text{Cost}_{\text{Net}} > \text{Cost}_{\text{Compute}}$	Roofline Model: The cost of transmitting tensors (NVLink) exceeds the cost of re-computing tokens via text.

While the criteria in Table 5 dictate *when* to switch, the efficacy of the transition is strictly bounded by the **Verification Asymmetry**. For NSED to yield positive gain, the ensemble's mean verification precision must satisfy the Condorcet condition $\bar{p}_v > 0.5 + \epsilon$. If the available expert pool consists of weak learners where $\bar{p}_v \approx 0.5$ (random guessing), the sequential integration described in Eq. 1 will fail to drift toward the decision boundary, making the monolithic baseline the superior choice regardless of task complexity.

6 Cost-Performance Trade-off Analysis

We define **Hardware Arbitrage** as the ability to achieve equivalent reasoning performance using a disjointed cluster of consumer flagship hardware ($C_{consumer}$) versus a monolithic enterprise node ($C_{enterprise}$).

6.1 The Memory Wall Problem

To provide a rigorous comparison, we restrict our baseline analysis exclusively to **Mixture-of-Experts (MoE)** architectures (e.g., GPT-OSS-120b, Qwen-3-235B). We exclude dense models (e.g., Llama-3-405B) from the efficiency analysis because MoEs already represent the industry frontier for parameter efficiency [38].

If the NSED protocol can demonstrate superior hardware economics against these highly optimized sparse architectures, the advantage over traditional dense models follows a fortiori.

State-of-the-art MoEs face a specific bottleneck: while their *compute* cost (FLOPs) is low due to sparse activation, their *memory* cost is strictly bounded by the total parameter count.

- **Storage Constraint:** Storing 235B parameters at FP4 quantization requires ≈ 117.5 GB of VRAM (plus KV cache overhead).
- **Hardware Requirement:** This exceeds the capacity of a single Nvidia A100 (80GB), necessitating either a dual-A100 setup ($\approx \$30,000$) or a single H200 ($\approx \$40,000$) connected via NVLink to handle the tensor parallelism.

6.2 NSED Topological Efficiency

In contrast, the NSED protocol is **Share-Nothing** (exchanging only text, not gradients). This allows the ensemble to be hosted on isolated consumer nodes without NVLink.

We define our Reference Hardware Baseline as a cluster of NVIDIA RTX 5090s (32GB VRAM). Unlike the monolithic baseline which requires continuous high-bandwidth memory, the NSED agents can operate effectively within the 32GB envelope of flagship consumer cards

by utilizing standard quantization techniques (e.g., 8-bit weights or FP8 KV Cache) where necessary.¹

Table 6 details the Total Cost of Ownership (TCO) comparison.

Table 6: Hardware Arbitrage: Monolithic MoE vs. NSED Consumer Cluster

Architecture	Model Composition	VRAM Req.	Hardware Baseline	Est. Hardware Cost
Monolithic (Baseline)	Qwen-3-235B (A22B)	~140 GB	2x H100 (NVLink)	\$60,000+
	GPT-OSS-120b	~80 GB	1x H100 (80GB)	\$30,000+
NSED Ensemble (Ours)	GPT-OSS-20B	~24 GB	1x RTX 5090	\approx \$6,000 - \$7,500
	Qwen-3-8B	~12 GB	1x RTX 5090	
	Gemma-3-12B (Verifier)	~28 GB	1x RTX 5090	

Result: The NSED topology achieves a **4x to 8x reduction in CAPEX**. By substituting High-Bandwidth VRAM (HBM3) with High-Latency Context (Iterative text refinement), we effectively trade time for money, allowing a commodity consumer cluster to compete with restricted enterprise infrastructure. While experimental validation was limited to the RTX 5090 tier, theoretical extrapolation suggests that utilizing prior-generation hardware (e.g., RTX 3090/4090) could further widen this arbitrage gap to a 10-20x margin.

6.3 The Latency Trade-off: System 2 Dynamics

While NSED achieves hardware arbitrage by utilizing consumer-grade VRAM, it introduces a **Synchronous Barrier Penalty**. Unlike monolithic models which generate tokens in a single continuous stream, NSED requires N agents to complete their reasoning (Generation Phase) before the barrier lifts for the Evaluation Phase. The round latency is therefore bounded by the slowest agent in the ensemble:

$$T_{round} = \max_{i \in \mathcal{A}} \left(\frac{\text{Tokens}_{gen}^{(i)}}{\text{Speed}_{serial}^{(i)}} \right) + \text{Overhead}_{net} \quad (9)$$

Our telemetry (Table 7) indicates that while the RTX 5090 cluster achieves an aggregate throughput of ≈ 850 tok/s via continuous batching (vLLM), the effective serial decoding speed for a single long-context Chain-of-Thought remains ≈ 250 tok/s.

Table 7: **NSED Latency Physics Breakdown.** Calculated using measured throughput on RTX 5090 cluster (250 t/s write, 2000 t/s read). The T_{round} decreases over time as the Dynamic KV-Cache amortizes historical context.

Deliberation Round	Gen. Phase (s)	Eval. Phase (s)	Round Total (s)	Cumulative (s)
Round 1	31.87	20.93	52.80	52.80
Round 2	31.73	18.06	49.79	102.59
Round 3	31.43	16.48	47.90	150.49
Round 4	27.82	15.83	43.65	194.14
Round 5	26.78	15.13	41.90	236.04
Round 6	28.01	14.55	42.55	278.60
Round 7	24.49	11.49	35.98	314.58

While a 5 minute inference time appears slow compared to standard chatbots, it is highly efficient for the domain. The AIME competition provides students with 3 hours to solve 15

¹Our validation run utilized an NVIDIA A40 for the Gemma node to expedite testing with unquantized BF16 weights. However, production deployment on an RTX 5090 is achievable via standard 8-bit KV caching or 4-bit weight quantization with negligible performance loss.

questions, averaging **12 minutes (720s) per problem**. The NSED protocol converges on a verified solution in **5 minutes (300s)**, effectively "thinking" faster than the target human expert, despite the overhead.

6.3.1 Prefix-Stable Latency Amortization (KV Caching)

The iterative topology of NSED is uniquely synergistic with recent advances in memory-efficient serving, specifically **PagedAttention** algorithms implemented in engines like vLLM [39].

Because the System Prompt (P_{sys}) and the persistent historical trunk ($H_{0...t-1}$) remain static across the N parallel generations of a given round, the system leverages **Prefix Caching** to eliminate redundant computation. Unlike standard independent queries, where the pre-fill phase constitutes the majority of latency, NSED’s recursive structure allows the KV-cache states to be reused. Our telemetry indicates that this topological alignment results in an effective **Cache Hit Rate of $\approx 40\%$** per token generation cycle, further decoupling the "Cost of Intelligence" from the raw parameter count.

6.4 The Bandwidth Arbitrage (Breaking the Roofline)

While Section 6 addresses the *Memory Wall* (VRAM), NSED also solves the *Bandwidth Wall*. Standard MoE architectures rely on Tensor Parallelism, requiring all-to-all communication of high-frequency gradients and activations. This imposes a hard hardware requirement: high-bandwidth interconnects (e.g., NVLink, ≈ 900 GB/s) found only in enterprise clusters.

NSED executes a **Bandwidth Arbitrage** by shifting the communication medium from *High-Frequency Tensors* to *Low-Frequency Semantics* (Natural Language tokens). By restricting inter-node communication to discrete symbolic states ($y_t \in \Sigma^*$) rather than dense vector spaces, the protocol reduces the bandwidth requirement by orders of magnitude. This effectively decouples the system’s intelligence from the "Interconnect Bottleneck," allowing decentralized consumer nodes (connected via standard PCIe or Ethernet) to function as a cohesive supercomputer.

7 Future Directions

7.1 High-Fidelity Signal Propagation

While the current NSED implementation relies on discrete natural language strings ($y_t \in \Sigma^*$) as the universal interface, this approach is inherently lossy. By sampling a single token at each step, the generator discards the rich probability distribution contained in the tail logits, effectively pruning potentially valid reasoning paths before they can be evaluated.

We propose a theoretical extension to the protocol: **The Logit-Exchange Lattice**.

7.1.1 From Strings to Semantic Confusion Networks

Instead of emitting a collapsed string, a future iteration of NSED could require agents to output a **Token Lattice \mathcal{L}** . For every position t , the agent emits the Top- K probable tokens along with their confidence scores.

This paradigm shifts the evaluator’s role from *Reviewer* to *Resolver*. By deferring the choice of particular token, the system prevents early hallucinations. If Agent A is unsure between two terms, it passes *both* options to Agent B in a structured format (e.g., in a simplest way, annotated as a XML or JSON structure):

The {cat|0.6, dog|0.3} sat on the {mat|0.8, rug|0.1}.

Agent B, possessing different context or capabilities, acts as a resolver, resolving the ambiguity that Agent A lacked the certainty to solve.

7.1.2 The Alignment Challenge: Optimal Transport (OT)

The primary barrier to implementing this in heterogeneous ensembles is the **Vocabulary Mismatch**. As noted by Minixhofer et al. [40], models like Qwen and Llama inhabit disjoint vector spaces, making direct logit arithmetic impossible.

To bridge these disjoint spaces, we propose utilizing **Optimal Transport** [41]. We define the alignment cost as the Wasserstein Distance $W_1(P_A, P_B)$ between the source and target vocabulary distributions. By pre-computing a sparse **Transport Matrix** $T_{A \rightarrow B}$, we can project the "uncertainty cloud" of Agent A directly into the vocabulary space of Agent B:

$$\hat{L}_B = T_{A \rightarrow B} \times L_A$$

While calculating exact Wasserstein distances is computationally intensive ($O(n^3)$), recent advances in Sinkhorn-Knopp regularization allow for efficient offline pre-computation, rendering the runtime cost to a simple sparse matrix multiplication [42, 43].

7.1.3 Implementation Strategy: Sparse Uncertainty Injection

To minimize the bandwidth overhead of transmitting lattices, we propose **Entropy-Gated Lattice Transmission**.

$$\text{Output}(t) = \begin{cases} \text{Token } w_t & \text{if } H(P_t) < \lambda \\ \text{Lattice } \mathcal{L}_t & \text{if } H(P_t) \geq \lambda \end{cases}$$

The system emits a dense semantic lattice \mathcal{L}_t only when the model's internal entropy $H(P_t)$ exceeds a confusion threshold λ . This ensures that communication overhead is incurred strictly for high-ambiguity tokens—the "forks in the road" of reasoning.

On the receiver side, a lightweight LoRA Adapter trained on "Confusion Networks" can learn to parse these probabilistic arrays, effectively allowing the ensemble to perform **Bayesian Reading**—weighing inputs based on the sender's confidence.

7.2 Deliberation-Native Adapters

Although our ablation studies indicate that homogeneous ensembles of generic base models suffer from "Groupthink" (reduced variance), the logistical benefits of a unified substrate—such as shared KV-caches and predictable latency—remain compelling.

To reconcile these conflicting constraints (Diversity vs. Uniformity), we propose the use of **Role-Specific Low-Rank Adapters (LoRA)** [44]. Instead of deploying entirely distinct model architectures (e.g., Qwen vs. Llama), a single foundational model can be dynamically "colored" at runtime. A "Critic Adapter" can be active during the evaluation phase to maximize scrutiny, while a "Creative Adapter" activates during generation to maximize entropy.

This approach effectively decouples *General Reasoning* (the base model) from *Protocol Roles* (the adapters), enabling:

- **Safety Adapters:** Specialized modules trained on datasets like DarkBench to rigorously audit proposals for manipulative patterns or bias, independent of the generator's alignment.
- **Domain Adapters:** Modules fine-tuned on niche corpora (e.g., Legal, Medical) to inject specific knowledge without retraining the verifying agents.

These lightweight modules can be dynamically swapped by the "Interface-Driven Constructor," creating a marketplace for modular cognitive skills while maintaining the **Computational Uniformity** required for efficient batching and caching.

7.3 Ephemeral-to-Long-Term Consolidation

Currently, NSED operates as a stateless inference engine; the cognitive labor of the ensemble is discarded upon session termination. We propose a *Post-Hoc Consolidation Phase*. In this extension, high-confidence consensus trajectories (y^*)—where the final voting entropy $H(V_T) \approx 0$ —are distilled into a synthetic fine-tuning corpus.

By applying parameter-efficient updates (e.g., LoRA) to the base agents using their own verified deliberation traces, we can achieve an **Autopoietic Improvement Cycle** [25]. This mechanism effectively mimics human experience replay where short-term "working memory" (deliberation) is encoded into long-term "weights" (synaptic plasticity) during idle periods.

7.4 Polymorphic Graph Switching (Meta-Cognitive Routing)

While the current implementation enforces a rigid Recurrent Consensus topology, we hypothesize that optimal cognitive efficacy requires **Topological Plasticity**. Not all queries necessitate the computational overhead of an N-way cyclic debate.

Future iterations of the Dynamic Expertise Broker will include a *Meta-Cognitive Router* capable of classifying task entropy to select the execution graph shape G_{exec} [22].

- **Low-Entropy Tasks:** Trigger a linear "Feed-Forward Chain" for maximum token efficiency.
- **High-Entropy Paradoxes:** Trigger an "Adversarial Lattice" (1-v-1 Debate) or the standard NSED Recurrent Loop.

This transitions the system from a static "Committee" architecture to a **Polymorphic Swarm**, where the organizational structure itself is a hyperparameter optimized at runtime via reinforcement learning signals [23].

7.5 Thermodynamic Efficiency & Entropy-Gated Halting

To further refine the hardware arbitrage frontier, we propose replacing the fixed round limit T_{max} with a **Thermodynamic Halting Condition**. Current "Chain-of-Thought" paradigms often suffer from "computation overhang," expending tokens long after the semantic solution has been reached.

By monitoring the rate of change in the consensus entropy ($\Delta H(S_t)$), the system can calculate the *Marginal Information Gain per Joule*. Execution should terminate when the Kullback-Leibler divergence between subsequent states drops below the energy cost of the next forward pass:

$$D_{KL}(S_t || S_{t-1}) < \epsilon_{cost}$$

This ensures the system operates at the Pareto frontier of **Cognitive Thermodynamics** [21], expending compute only when it yields statistically significant semantic resolution.

8 Conclusion

The transition from static pre-training to dynamic inference-time compute necessitates a fundamental re-evaluation of cognitive architectures. In this paper, we introduced the N-Way Self-Evaluating Deliberation (NSED) protocol, demonstrating that robust reasoning is not solely a function of parameter count, but of Topological Governance. By formalizing the multi-agent interaction as a recurrent loop, we bridge the gap between connectionist control theory and agentic workflows.

Our empirical validation confirms that this topology allows ensembles of small, consumer-grade models to match the performance of monolithic state-of-the-art systems. This establishes

a verified frontier of **Hardware Arbitrage**, proving that deliberative process can substitute for model scale. By trading latency for recurrent refinement, we enable disjointed clusters of commodity hardware to compete directly with enterprise-scale infrastructure.

Furthermore, our analysis of the influence dynamics confirms that identity-blind Quadratic Voting effectively mitigates “herding” and sycophancy. While our current experiments utilized fixed agent profiles, the observed telemetry validates the potential for a fully dynamic Broker that optimizes team composition based on historical signal-to-noise ratios. Looking forward, the integration of our proposed efficiency-fatigue stopping conditions and offline consolidation loops promises to evolve NSED from a static inference engine into a system that optimizes its energy expenditure in pursuit of semantic convergence. We posit that this shift from bigger weights to better circuits represents a necessary path toward sustainable, decentralized Artificial General Intelligence.

References

- [1] Charlie Snell, Jaehoon Lee, Kelvin Xu, and Aviral Kumar. Scaling llm test-time compute optimally can be more effective than scaling model parameters, 2024. URL <https://arxiv.org/abs/2408.03314>.
- [2] Ali Behrouz, Peilin Zhong, and Vahab Mirrokni. Titans: Learning to memorize at test time, 2024. URL <https://arxiv.org/abs/2501.00663>.
- [3] Arnub Tandon and et. al. End-to-end test-time training for long context, 2025. URL <https://arxiv.org/abs/2512.23675>.
- [4] Ali Behrouz, Meisam Razaviyayn, Peilin Zhong, and Vahab Mirrokni. Nested learning: The illusion of deep learning architecture, accessed at 3 jan 2026. URL <https://abehrouz.github.io/files/NL.pdf>.
- [5] DeepSeek-AI, Daya Guo, and et. al. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning, 2025. URL <https://arxiv.org/abs/2501.12948>.
- [6] Aaron Jaech and et. al. Openai o1 system card, 2024. URL <https://arxiv.org/abs/2412.16720>.
- [7] William Fedus, Barret Zoph, and Noam Shazeer. Switch transformers: Scaling to trillion parameter models with simple and efficient sparsity, 2022. URL <https://arxiv.org/abs/2101.03961>.
- [8] Junlin Wang, Jue Wang, Ben Athiwaratkun, Ce Zhang, and James Zou. Mixture-of-agents enhances large language model capabilities, 2024. URL <https://arxiv.org/abs/2406.04692>.
- [9] Yusen Zhang, Ruoxi Sun, Yanfei Chen, Tomas Pfister, Rui Zhang, and Sercan Ö. Arik. Chain of agents: Large language models collaborating on long-context tasks, 2024. URL <https://arxiv.org/abs/2406.02818>.
- [10] Joon Sung Park, Joseph C O’Brien, Carrie J Cai, Meredith Ringel Morris, Percy Liang, and Michael S Bernstein. Generative agents: Interactive simulacra of human behavior. 2023. URL <https://arxiv.org/abs/2304.03442>.
- [11] Xie Yi, Zhanke Zhou, Chentao Cao, Qiyu Niu, Tongliang Liu, and Bo Han. From debate to equilibrium: Belief-driven multi-agent llm reasoning via bayesian nash equilibrium, 2025. URL <https://arxiv.org/abs/2506.08292>.

- [12] Mert Cemri, Melissa Z. Pan, Shuyi Yang, Lakshya A. Agrawal, Bhavya Chopra, Rishabh Tiwari, Kurt Keutzer, Aditya Parameswaran, Dan Klein, Kannan Ramchandran, Matei Zaharia, Joseph E. Gonzalez, and Ion Stoica. Why do multi-agent llm systems fail?, 2025. URL <https://arxiv.org/abs/2503.13657>.
- [13] Nelson F. Liu, Kevin Lin, John Hewitt, Ashwin Paranjape, Michele Bevilacqua, Fabio Petroni, and Percy Liang. Lost in the middle: How language models use long contexts, 2023. URL <https://arxiv.org/abs/2307.03172>.
- [14] Mrinank Sharma, Meg Tong, Tomasz Korbak, David Duvenaud, Amanda Askill, Samuel R. Bowman, Newton Cheng, Esin Durmus, Zac Hatfield-Dodds, Scott R. Johnston, Shauna Kravec, Timothy Maxwell, Sam McCandlish, Kamal Ndousse, Oliver Rausch, Nicholas Schiefer, Da Yan, Miranda Zhang, and Ethan Perez. Towards understanding sycophancy in language models, 2025. URL <https://arxiv.org/abs/2310.13548>.
- [15] Esben Kran, Hieu Minh "Jord" Nguyen, Akash Kundu, Sami Jawhar, Jinsuk Park, and Mateusz Maria Jurewicz. Darkbench: Benchmarking dark patterns in large language models, 2025. URL <https://arxiv.org/abs/2503.10728>.
- [16] Qingyun Wu, Gagan Bansal, Jieyu Zhang, Yiran Wu, Beibin Li, Erkang Zhu, Li Jiang, Xiaoyun Zhang, Shaokun Zhang, Jiale Liu, Ahmed Hassan Awadallah, Ryen W White, Doug Burger, and Chi Wang. Autogen: Enabling next-gen llm applications via multi-agent conversation. *arXiv preprint arXiv:2308.08155*, 2023. URL <https://arxiv.org/abs/2308.08155>.
- [17] Bradley Brown, Jordan Juravsky, Ryan Ehrlich, Ronald Clark, Quoc V. Le, Christopher Ré, and Azalia Mirhoseini. Large language monkeys: Scaling inference compute with repeated sampling, 2024. URL <https://arxiv.org/abs/2407.21787>.
- [18] Michael Henry Tessler, Michiel A. Bakker, Daniel Jarrett, Hannah Sheahan, Martin J. Chadwick, Raphael Koster, Georgina Evans, Lucy Campbell-Gillingham, Tantum Collins, David C. Parkes, Matthew Botvinick, and Christopher Summerfield. Ai can help humans find common ground in democratic deliberation. *Science*, 386(6719):eadq2852, 2024. doi: 10.1126/science.adq2852. URL <https://www.science.org/doi/abs/10.1126/science.adq2852>.
- [19] Nicholas Rescher. *Predicting the future: An introduction to the theory of forecasting*. State University of New York Press, Albany, 1998. ISBN 978-0-7914-3553-3.
- [20] Steven Lalley and Eric Weyl. Quadratic voting: How mechanism design can radicalize democracy, 2018. URL <http://dx.doi.org/10.2139/ssrn.2003531>.
- [21] Roy Schwartz, Jesse Dodge, Noah A Smith, and Oren Etzioni. Green ai. *Communications of the ACM*, 63(12):54–63, 2020.
- [22] Maciej Besta, Nils Blach, Ales Kubicek, Robert Gerstenberger, Lukas Gianinazzi, Joanna Gajda, Tomasz Lehmann, Michal Podstawski, Hubert Niewiadomski, Piotr Nolle, and Torsten Hoefer. Graph of thoughts: Solving elaborate problems with large language models. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pages 17682–17690, 2024.
- [23] Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Thomas L Griffiths, Yuan Cao, and Karthik Narasimhan. Tree of thoughts: Deliberate problem solving with large language models. In *Advances in Neural Information Processing Systems*, volume 36, pages 11830–11843, 2024.

- [24] Peng Jin, Bo Zhu, Li Yuan, and Shuicheng Yan. Moe++: Accelerating mixture-of-experts methods with zero-computation experts, 2024. URL <https://arxiv.org/abs/2410.07348>.
- [25] Eric Zelikman, Yuhuai Wu, Jesse Mu, and Noah D. Goodman. Star: Bootstrapping reasoning with reasoning. 2022. URL <https://arxiv.org/abs/2203.14465>.
- [26] Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, Christopher Hesse, and John Schulman. Training verifiers to solve math word problems. *arXiv preprint arXiv:2110.14168*, 2021. Empirically demonstrates that a dedicated verifier scales more effectively with data than a generator, establishing the existence of a ‘Verification Gap’ where $p_{verify} > p_{generate}$.
- [27] Hunter Lightman, Vineet Kosaraju, Yura Burda, Harrison Edwards, Bowen Baker, Teddy Lee, Jan Leike, John Schulman, Ilya Sutskever, and Karl Cobbe. Let’s verify step by step. *arXiv preprint arXiv:2305.20050*, 2023. Shows that Process Reward Models (PRMs) significantly outperform outcome-based supervision, validating the NSED approach of step-wise consensus.
- [28] Abraham Wald. *Sequential Tests of Statistical Hypotheses*. John Wiley & Sons, 1945.
- [29] Co Principal Editors: Edward N. Zalta and Uri Nodelman, editors. *The Stanford Encyclopedia of Philosophy*. URL <https://plato.stanford.edu/entries/jury-theorems/#CondJuryTheo>.
- [30] Leo Breiman. Statistical modeling: The two cultures. *Statistical Science*, 16(3):199–231, 2001.
- [31] Thomas M Cover. Geometrical and statistical properties of systems of linear inequalities with applications in pattern recognition. *IEEE transactions on electronic computers*, (3): 326–334, 1965.
- [32] Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik Narasimhan, and Yuan Cao. React: Synergizing reasoning and acting in language models, 2023. URL <https://arxiv.org/abs/2210.03629>.
- [33] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you need, 2023. URL <https://arxiv.org/abs/1706.03762>.
- [34] Jakob Puchinger, Günther R. Raidl, and Ulrich Pferschy. The multidimensional knapsack problem: Structure and algorithms. URL <https://inria.hal.science/hal-01224914v1/document>.
- [35] Ari Holtzman, Jan Buys, Li Du, Maxwell Forbes, and Yejin Choi. The curious case of neural text degeneration. In *International Conference on Learning Representations*, 2020. URL <https://openreview.net/forum?id=rygGQyrFvH>.
- [36] Ilia Shumailov, Zakhar Shumaylov, Yiren Zhao, Yarin Gal, Nicolas Papernot, and Ross Anderson. The curse of recursion: Training on generated data makes models forget. *arXiv preprint arXiv:2305.17493*, 2023. Also published in Nature (2024) as ‘AI models collapse when trained on recursively generated data’.
- [37] Livecodebench v5 leaderboard. URL https://livecodebench.github.io/leaderboard_v5.html.

- [38] Albert Q. Jiang, Alexandre Sablayrolles, Antoine Roux, and et al. Mixtral of experts, 2024. URL <https://arxiv.org/abs/2401.04088>.
- [39] Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng, Cody Hao Yu, Joseph E. Gonzalez, Hao Zhang, and Ion Stoica. Efficient memory management for large language model serving with pagedattention, 2023. URL <https://arxiv.org/abs/2309.06180>.
- [40] Benjamin Minixhofer, Ivan Vulić, and Edoardo Maria Ponti. Universal cross-tokenizer distillation via approximate likelihood matching, 2025. URL <https://arxiv.org/abs/2503.20083>.
- [41] Nicolas Boizard et al. Towards cross-tokenizer distillation: the universal logit distillation loss for llms. *arXiv preprint arXiv:2402.12030*, 2024.
- [42] Marco Cuturi. Sinkhorn distances: Lightspeed computation of optimal transport. In *Advances in Neural Information Processing Systems*, volume 26, pages 2292–2300, 2013.
- [43] Gabriel Peyré and Marco Cuturi. Computational optimal transport. *Foundations and Trends in Machine Learning*, 11(5-6):355–607, 2019.
- [44] Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. Lora: Low-rank adaptation of large language models. In *International Conference on Learning Representations*, 2022. URL <https://openreview.net/forum?id=nZeVKeeFYf9>.

A Hyperparameters & Sampling Configuration

To ensure reproducibility, we document the specific sampling parameters utilized for the agents. While top-p and min-p were left at provider defaults ($p = 0.9$, $\min_p = 0.05$), the Temperature and Presence Penalty were strictly controlled via the NSED configuration profile to enforce diverse role-playing.

Table 8: Agent Persona & Sampling Configuration

Agent	Base Model	Temp (T)	Penalty (α)	Max Tokens	Role Intent
Jaya	GPT-OSS-20B	0.2	1.5	16,000	<i>Creative Architect</i> (Deterministic Logic)
Xue	Qwen3-8B	0.6	1.5	16,000	<i>Balanced Engineer</i> (Exploratory)
Alic	Gemma-3-12B	0.6	1.5	16,000	<i>Rigorous Analyst</i> (High-Entropy Critique)

B Model Serving Infrastructure

The empirical validation was conducted using the vLLM engine. To achieve the throughput required for synchronous N-way deliberation on limited hardware, we utilized specific optimization flags, notably **FP8 KV-Caching** and **Chunked Prefill**.

Table 9: vLLM Serving Configuration (Heterogeneous Ensemble)

Model	Dtype / Quant	KV Cache	Context	TP	Attn. Backend	Special Optimizations
Gemma-3-12B-IT	bfloat16	fp8	32k	1	FlashInfer	-tool-call-parser pythonic -enable-chunked-prefill
GPT-OSS-20B	auto (fp16)	fp8	32k	2*	FlashInfer	-enable-chunked-prefill -swap-space 20
Qwen3-8B	bfloat16	fp8	64k	1	FlashInfer	-tool-call-parser hermes -reasoning-parser deepseek_r1 -rope-scaling yarn (3.0)

* Note: While GPT-OSS-20B used Tensor Parallelism (TP=2) in our reference run for unquantized precision, it fits on a single 24GB consumer card (RTX 3090/4090) when loaded with 4-bit AWQ quantization, maintaining the consumer-hardware thesis described in Section 6.

RoPE Scaling Implementation: For the Qwen3-8B node, we applied dynamic YaRN scaling to extend the effective context window to 64k tokens without fine-tuning, ensuring the model could ingest the full history of 7-round deliberations.

```
--hf-overrides '{"rope_scaling": {"rope_type": "yarn", "factor": 3.0,
"original_max_position_embeddings": 32768}}'
```

Tool & Reasoning Parsing: To support the heterogeneous agent protocols, we utilized specialized parsers at the inference server level. The **pythonic** parser (Gemma) and **hermes** parser (Qwen) were enabled to handle the distinct function-calling tokens of each architecture, while the **deepseek_r1** reasoning parser was employed to segment "Chain-of-Thought" blocks from final answers.