

Distributed Intelligence: Architecting Multi-Agent Swarms and Decentralized Expert Flows for Scalable Large Language Model Inference

Authors: John Villwock and Aubra Taylor, Ph.D.

Institution: ReEnvision AI

Abstract

The trajectory of generative artificial intelligence has historically been dominated by the scaling hypothesis, relying on massive monolithic models deployed in centralized clusters. This paper presents a novel alternative: coordinating smaller, specialized Large Language Models (LLMs) in a decentralized Mixture-of-Experts (dMoE) and Mixture-of-Agents (MoA) swarm. By introducing the concept of Speculative Orchestration—where specialized draft models predict complex reasoning paths verified by a central orchestrator—we demonstrate a paradigm that mathematically optimizes the cost-latency tradeoff. Combined with Dynamic Speculative Planning (DSP) and no-regret drafter selection routing, the proposed ReEnvision AI architecture achieves intelligence parity with frontier monolithic models at unprecedented execution speeds across completely decentralized hardware topologies.

The Paradigm Shift Toward Decentralized Multi-Agent Architectures

The trajectory of generative artificial intelligence has historically been dominated by the scaling hypothesis, which posits that exponentially increasing the parameter count of monolithic Large Language Models (LLMs) inherently yields superior cognitive and reasoning capabilities. However, this centralized approach presents profound bottlenecks in compute resource availability, energy consumption, and logical degradation. Massive models deployed on centralized clusters face issues such as the "Lost-in-the-Middle" effect, wherein vast context windows fail to accurately retrieve or process deeply embedded reasoning steps during complex execution sequences.¹ Recent empirical research demonstrates a formidable, structurally divergent alternative: coordinating smaller, specialized LLMs—such as 7B to 25B parameter models—in highly parallelized, decentralized swarms to achieve intelligence parity with, or superiority over, massive monolithic structures.²

This architecture effectively mimics a sophisticated organizational hierarchy. In such a framework, a moderately sized model (for instance, a 25B parameter orchestrator) serves as a product manager, dynamically decomposing complex tasks and delegating highly specific

workloads to a fleet of smaller (e.g., 7B parameter) expert models.³ These experts, highly specialized in discrete domains such as mathematical reasoning, code execution, or literature research, operate simultaneously across entirely separate physical machines.⁵ By distributing the workload across decentralized hardware, each agent maintains a pristine, independent context window, preventing the logical contamination and token exhaustion that plagues singular massive models.⁶ Subsequently, an independent layer of verifier models synthesizes the disparate outputs, engaging in parallel checking, debating, and majority voting to determine the optimal result before returning the final execution to the user.⁷ The second-order implication of this paradigm is a fundamental shift in AI infrastructure economics, moving the industry away from reliance on centralized, capital-intensive supercomputers toward decentralized, asynchronous edge networks utilizing consumer-grade hardware and specialized routing mechanisms.⁹

Foundational Frameworks: Mixture-of-Experts versus Mixture-of-Agents

To accurately contextualize multi-agent swarms operating across discrete machines, it is imperative to deeply delineate the architectural and operational distinctions between the traditional Mixture-of-Experts (MoE) methodology and the rapidly emerging Mixture-of-Agents (MoA) paradigm. While both seek to leverage specialization to enhance output quality and optimize compute, their execution layers, network topologies, and hardware requirements represent fundamentally different approaches to artificial intelligence design.

Internal Token Routing: Mixture of Experts (MoE)

Mixture-of-Experts operates strictly within the internal neural architecture of a single forward pass during inference.¹² The core premise behind MoE models originates from foundational research conducted in 1991 concerning adaptive mixtures of local experts, which proposed training an AI system composed of separate networks specializing in different subsets of training cases.¹³ In contemporary LLMs, such as Mistral's Mixtral 8x7B, a routing or gating network evaluates the incoming input data and selectively routes specific tokens to distinct, highly specialized sub-networks, or "experts," embedded directly within the model's multi-layer perceptron architecture.¹³

This internal routing mechanism allows a model containing hundreds of billions of parameters to selectively activate only a minor fraction of its neural network for any given task, drastically reducing the computational cost per token and accelerating both pre-training and inference.¹³ While highly efficient in terms of sparse activation, traditional MoE remains a centralized, monolithic process. It is constrained by a single, shared context window and a unified attention mechanism. Furthermore, because the experts are sheer in size and tightly coupled, they are typically bound to massive, centralized GPU clusters utilizing tensor parallelism, where weights must be sharded across high-bandwidth memory architectures.¹⁶ This centralization fundamentally limits the ability to isolate context or deploy individual reasoning agents across

truly decentralized hardware topologies.

External Collaborative Synthesis: Mixture-of-Agents (MoA)

In stark contrast to the internal, token-level routing of MoE, Mixture-of-Agents (MoA) operates entirely at the application and inference orchestration layer, treating multiple, complete, and distinct LLMs as discrete entities that collaborate iteratively over a network.¹² The MoA framework constructs a layered pipeline comprising multiple "Proposer" agents and subsequent "Aggregator" agents.¹⁹ The methodology thrives on model diversity, explicitly combining varied architectures to eliminate the shared weaknesses inherent in relying on a single structural design.¹⁹

During execution, multiple independent open-source models—such as Qwen, Llama 3, Mixtral, and DBRX—receive the identical initial prompt as Proposer agents.²¹ Because each model possesses unique architectural weights, distinct pre-training corpora, and specialized fine-tuning, they generate highly diverse candidate responses.¹⁹ The outputs from all Proposer agents are subsequently fed into the next layer of Aggregator models. These Aggregator models are explicitly prompted to synthesize the diverse perspectives, discard hallucinated logic, and merge the valid components into a single, highly refined response.²⁰ The entire system functions via natural language communication and prompt engineering, requiring absolutely no fine-tuning of the underlying model weights.¹⁹

The empirical success of the MoA architecture demonstrates the inherent collaborativeness of language models. Research reveals that LLMs generate substantially higher-quality text when they have access to the varied outputs of their peers as auxiliary information, even if some of those individual peer outputs are objectively lower in quality.² Implementations of the Together AI MoA architecture have achieved state-of-the-art results, scoring 65.1% on the AlpacaEval 2.0 benchmark using exclusively open-source models, significantly surpassing the 57.5% score achieved by the heavily capitalized GPT-4 Omni.¹

Architecture Feature	Mixture-of-Experts (MoE)	Mixture-of-Agents (MoA)
Execution Layer	Model internals (Neural Network weights)	Orchestration layer (API / Inference generation)
Routing Mechanism	Token-level internal gating mechanism	Full-response aggregation and explicit synthesis

Context Window	Shared context across all internal experts	Independent context windows for every agent
Inference Latency	Highly optimized, parallelized internally	Higher latency due to multi-layered sequential routing
Hardware Topology	Tightly coupled, centralized GPU clusters	Distributed, hardware-agnostic, decentralized networking
Model Diversity	Single architectural framework	Combines multiple distinct architectures natively

The Mathematics of Ensembles: Scaling Laws and Prior Probability

The mathematical and statistical foundation validating the efficacy of multi-agent swarms is rigorously articulated in research identifying that simply increasing the number of instantiated agents profoundly scales the intelligence of the overall system.²⁴ The research presents a comprehensive exploration of raw agent scaling, proving that deploying a swarm of models executing a parallel sampling-and-voting protocol yields consistent performance improvements across diverse benchmarks spanning reasoning and generation tasks.²⁶

Sampling, Voting, and Task Difficulty Optimization

When deploying a homogeneous swarm of smaller models—such as instantiating ten parallel versions of a 7B parameter expert—the system relies on generating independent reasoning paths concurrently. The swarm approaches the prompt, generates diverse outputs by leveraging temperature scaling and inference variations, and then aggregates these paths through a majority voting mechanism to determine the most logically sound output.²⁴ Researchers have determined that this raw ensemble methodology is orthogonal and highly complementary to existing complex reasoning enhancements, such as Chain-of-Thought (CoT) prompting or Debate mechanisms.²⁶ By stacking multi-agent ensembling on top of CoT, the system achieves synergistic accuracy gains.

A critical third-order insight derived from these ensemble dynamics is the inverse relationship between individual agent capability and the scaling benefits of the swarm. Smaller, comparatively weaker LLMs gain disproportionately higher performance improvements from ensembling than massive, frontier models.²⁷ Furthermore, the degree of enhancement correlates directly with the inherent difficulty of the task. Highly complex problems, such as multi-step mathematical reasoning or intricate code generation, exhibit dramatic performance spikes when subjected to agent swarms, whereas trivial tasks plateau quickly.²⁶ This implies a profound economic shift: for highly complex enterprise logic or rigorous mathematical validation, deploying a massive swarm of cost-effective 7B models executing in parallel will mathematically outperform a single, computationally expensive monolithic model.²⁴

Hierarchical Step-Wise Ensembling and Prior Probability

To mathematically optimize this probabilistic phenomenon, advanced implementations discard simple end-to-end voting in favor of step-wise sampling-and-voting methods. Instead of the swarm voting solely on the final output string, the orchestrator decomposes the task into a sequence of discrete logical steps.²⁶ At each individual step, the swarm of agents generates multiple intermediate states, votes on the optimal intermediate progression, and updates the task state before proceeding to the subsequent logical hurdle.²⁶ Experimental data reveals that this step-wise hierarchical aggregation yields accuracy gains of 15% to 42% over simple end-to-end voting, effectively eliminating cascading hallucination errors early in the generation cycle.²⁶

Furthermore, system architects can manipulate the prior probability of success by tuning parameters related to the number of agents and the difficulty of the subtasks. By decomposing tasks hierarchically into high-probability subtasks and addressing them with homogeneous or heterogeneous combinations of models (e.g., mixing smaller localized open-source models with larger API-based models), the system maximizes the statistical likelihood of reaching the correct conclusion.²⁸ This approach not only enhances absolute performance but introduces critical cost savings by dynamically employing simpler, faster models for easier subtasks while reserving dense computational swarms for the most complex logical nodes.²⁸

Hierarchical Orchestration: The Product Manager and Specialized Worker Paradigm

Translating raw mathematical agent ensembling into functional, enterprise-grade software engineering requires highly sophisticated hierarchical orchestration. The most effective topology emerging in the field directly mirrors human corporate structures: a primary reasoning model acts as an intelligent "Product Manager" or Orchestrator, dynamically governing a highly specialized pool of "Worker" models.³

The ParaManager and Agent-as-Tool Architecture

The "Agent-as-Tool" paradigm effectively unifies autonomous reasoning agents and traditional

external software APIs into a single, standardized execution space.⁴ In frameworks such as ParaManager, a lightweight orchestrator (e.g., a 25B parameter model) is explicitly trained on system-level objectives, focusing entirely on planning, subtask decomposition, state-aware delegation, and asynchronous execution.⁴ This architecture completely decouples the strategic planning phase from the tactical execution phase.

When presented with a highly complex user query, the Orchestrator evaluates the multi-faceted requirements and spawns parallel sub-agents.³⁰ Each sub-agent is initialized with a highly specific system prompt defining its persona and limitations, such as a "Mathematical Verifier," a "Code Synthesizer," or a "Literature Researcher".³¹ The critical mechanism driving the success of this hierarchy is absolute context isolation.⁴ Because the Orchestrator handles the overarching meta-logic and state tracking, the specialized 7B worker models operate with pristine, uncontaminated independent context windows.⁶ They are not burdened by the token-heavy instructions of the overall master plan; rather, they receive only the highly specific parameters required for their designated micro-task.³

This isolation drastically curtails the risk of error propagation. In massive monolithic models, an attention failure or hallucination early in a 100,000-token context window can irrevocably corrupt all subsequent logic.⁴ In an orchestrator-worker swarm, if the localized "Mathematical Verifier" agent produces an illogical output or execution fault, the Orchestrator can simply terminate that specific node, utilize explicitly trained recovery mechanisms to identify the failure state, and spawn a new verifier agent without jeopardizing the stability of the entire systemic workflow.⁴

Real-World Implementation: Swarm Mode and Git Worktree Isolation

These theoretical orchestration paradigms are increasingly reflected in cutting-edge production environments. An illustrative example is the native "Swarm Mode" orchestration layers integrated into advanced AI coding assistants, which utilize deep internal coordination protocols such as the TeammateTool architecture developed alongside Claude Sonnet.³ In these highly structured implementations, the lead agent generates sub-agents for specialized engineering roles—frontend development, backend architecture, unit testing, and documentation—and provisions them with independent communication channels.³³

To facilitate this architecture physically and prevent destructive interference between agents, robust systems rely on isolation mechanisms like Git Worktrees.³ This allows multiple autonomous coding agents to execute concurrently in isolated filesystem environments. They manipulate codebases, run specific test suites, and generate modular features entirely in parallel.³ The Orchestrator monitors the independent progress of each workspace and automatically merges the code only after the specialized testing agents confirm that all isolated implementations pass their respective verification suites.³ Anthropic's internal research utilizing the BrowseComp evaluation confirmed that managing token usage and context window distribution across multiple specialized agents accounts for 80% of the performance variance in multi-agent systems, validating that physical and logical isolation enables a 5x to 10x

improvement in complex software engineering workflows compared to sequential, single-agent generation.³

Parallel Execution and Systemic Speed Optimization Strategies

A fundamental critique of deploying layered Mixture-of-Agents or orchestrating massive voting swarms is the geometric escalation of wall-clock latency.³⁵ Generating dozens of parallel responses and routing them sequentially through subsequent aggregator models inherently consumes significantly more temporal and computational resources than a single monolithic inference pass.¹² This latency multiplication can render multi-agent systems slower than monolithic alternatives if not managed with rigorous architectural precision.³⁵ Thus, production-grade swarms mandate advanced parallel execution frameworks.

M1-Parallel and Asynchronous Event-Driven Architectures

To directly address the latency deficits inherent in sequential multi-step task execution, researchers have developed frameworks explicitly optimized for speed, such as M1-Parallel built upon the Magentic-One system.³⁶ Rather than running a single team of agents sequentially through a complex logical maze, M1-Parallel utilizes an event-driven communication model with asynchronous messaging to concurrently spawn multiple separate multi-agent teams.³⁶ Each distinct team attempts to solve the identical overall problem by exploring diverse execution plans simultaneously across parallel hardware nodes.³⁶

The critical innovation enabling massive latency reduction in this architecture is the implementation of an "early termination" protocol. Because the parallel multi-agent swarms are executing asynchronously on separate machines, they progress at different rates based on the complexity of their chosen reasoning path.³⁶ The first swarm to successfully reach a verified, accurate conclusion immediately broadcasts a termination signal across the network, instantly halting the computation of all other exploratory swarms.³⁶ Through raw parallel execution and aggressive culling of slower nodes, this methodology forces the systemic latency to match the speed of the single fastest valid reasoning path.³⁶ Empirical evaluations on highly complex reasoning tasks demonstrate that M1-Parallel with early termination achieves up to a 2.2x speedup in end-to-end wall-clock execution time compared to traditional sequential multi-agent execution, while rigorously preserving task accuracy.³⁶

Speculative Orchestration: Enhancing dMoE through Predictive Action

To further eliminate the latency penalty of iterative agent communication, our proposed framework incorporates **Speculative Orchestration**—applying the mechanics of speculative decoding directly to the multi-agent routing layer. Instead of the 25B orchestrator waiting sequentially for tools or sub-agents to compute, lightweight 7B "Speculator" models predict the

exact sequence of reasoning steps, API calls, or expert outputs ahead of time. The 25B "Actor" orchestrator then verifies these drafted actions in massive parallel batches.

Mathematical Modeling of Speculative Speedup and Dynamic Planning

The theoretical speedup of speculative execution can be modeled by evaluating the expected number of accepted tokens against the relative cost of the draft model. Let

$c = T_{draft} / T_{target}$ represent the relative latency ratio between the fast 7B draft model and

the 25B target model, and let γ represent the number of tokens or actions drafted per speculation cycle. The mathematical speedup factor over standard autoregressive decoding is

approximated by the function $\frac{1-\alpha^{\gamma+1}}{(1-\alpha)(\gamma c+1)}$, where α is the effective acceptance probability of the speculative sequence.

To optimize this mathematical ratio in a dynamic swarm without incurring excessive token waste on incorrect predictions, we employ **Dynamic Speculative Planning (DSP)**. DSP treats the optimal speculation depth k as a real-time value-estimation problem, utilizing an online reinforcement learning framework powered by $TD(\lambda)$ to dynamically adjust how many steps the draft model predicts on the fly. This effectively steers the swarm to make aggressive, high- γ predictions for simple, high-confidence subtasks, while forcing careful, step-by-step verification on dense logic tasks. This mathematical control mechanism yields comparable execution acceleration to fixed-depth methods, while simultaneously reducing unnecessary computation costs by up to 60%.

Provably No-Regret Drafter Routing and Expert Budgeting

Maximizing the acceptance probability (α) requires the orchestration layer to select the perfect speculative draft model for each specific subtask. We achieve this via adaptive routing protocols like **HedgeSpec** and **SpecRouter**, which dynamically profile incoming prompts. HedgeSpec utilizes a full-information online learning feedback loop to evaluate the performance of all available 7B expert drafters in real-time. This provides a mathematically provable, no-regret guarantee that the task is routed to the optimal specialized drafter, dramatically outperforming static routing heuristics—especially on long reasoning chains.

Furthermore, parallel speculative paths across decentralized agents threaten to trigger an immense number of unique expert activations, creating a severe memory bandwidth bottleneck across network nodes. To counteract this, the architecture applies **MoE-Spec**, a mathematical budgeting algorithm that enforces a strict upper limit on the number of experts loaded into VRAM simultaneously. By evaluating aggregate routing probabilities and loading only the top-scoring experts mathematically most likely to contribute to the verification phase, MoE-Spec decouples speculation depth from memory bandwidth constraints, yielding 10% to

30% higher throughput than standard speculative decoding baselines.

Decision Protocols: Verifiers, Consensus, and Multi-Agent Debate Dynamics

Once specialized expert models generate parallel responses, the multi-agent system must possess a resilient, mathematically sound protocol for evaluating truth, discarding hallucinations, and synthesizing the final output. The literature identifies several primary modalities for this checking phase: Explicit Verifier Validation, Unanimous Consensus Debate, and Pure Majority Voting.³⁸ The selection of the decision protocol dictates both the ultimate intelligence of the swarm and its execution speed.

Explicit Verifier Agents and Sample Set Aggregators

The most structurally straightforward quality control mechanism involves deploying a dedicated subset of LLMs purely as "Validators" or "Checkers." While worker agents utilize generative sampling methodologies to explore reasoning paths, the Verifier agent operates on a strict evaluative system prompt.³⁸ A robust verifier explicitly checks for question relevance, mathematical soundness, execution safety, and strict adherence to the orchestrator's initial constraints, evaluating whether the proposed sequence directly solves the core premise without drift.⁴⁰

Recent advancements in scaling test-time compute involve training highly compact LLMs specifically as Sample Set Aggregators (SSA). These specialized models take a concatenated sequence of multiple candidate reasoning paths generated by the worker swarm, rank them dynamically, and output the final verified answer, achieving immense performance gains explicitly in mathematical and coding domains where objective ground-truth can be verified algorithmically.²³

Multi-Agent Debate and The Vulnerabilities of Consensus

In high-stakes scenarios lacking clear, objective ground truth (such as complex strategic analysis or nuanced literature review), systems frequently employ Multi-Agent Debate mechanisms.⁷ Rather than relying on a static aggregator, agents actively present their initial findings to one another across multiple communication rounds, attempting to persuade or correct their peers.²³

However, rigorous research systematically assessing different decision protocols reveals highly nuanced trade-offs between "Consensus" protocols (requiring unanimous agreement among all debating agents) and simple "Voting" protocols.³⁹ Data demonstrates that consensus-based protocols improve performance by 2.8% specifically in knowledge-retrieval tasks, as the iterative multi-round debate allows agents to correct factual inaccuracies via peer review.³⁹ Yet, in complex logical reasoning and mathematical deduction, demanding consensus actively harms the output. In these logical domains, simple independent majority voting improves

performance by a staggering 13.2% compared to consensus debate.³⁹

This severe discrepancy arises from the inherent conformity bias programmed into language models during alignment training. During extended multi-round debates, agents that initially produced mathematically correct reasoning are often swayed by the confident, verbose, yet entirely hallucinatory assertions of their peers, leading to cascading error propagation and a rapid degradation of the overall swarm intelligence.⁸

Free-MAD: Consensus-Free Debate and Anti-Conformity

To systematically mitigate this conformity collapse while preserving the intellectual benefits of multi-agent interaction, advanced architectures employ frameworks such as Consensus-Free Multi-Agent Debate (Free-MAD).⁸ Free-MAD fundamentally eliminates the requirement for agents to reach an agreement. Instead, the framework reconstructs the debate phase by introducing explicit "anti-conformity" prompting instructions, forcing individual agents to stubbornly defend their distinct logical branches and resist the influence of the majority.⁸

The final decision in the Free-MAD framework is managed not by the debating agents themselves, but by an external orchestrator utilizing a novel score-based mechanism. This orchestrator evaluates the entirety of the debate trajectory, tracking precisely how each agent's reasoning evolves over time rather than simply taking the last round's output.⁸ This maintains extreme intellectual diversity within the swarm right up until the final aggregation step, preventing the premature convergence on suboptimal solutions that plagues traditional consensus protocols.⁸ Furthermore, evaluating these debate dynamics utilizing a time-varying mixture of Beta-Binomial distributions and employing adaptive stopping criteria based on the Kolmogorov-Smirnov statistic allows the system to halt debates at the exact mathematical moment accuracy peaks, maintaining computational efficiency.⁷

Intelligent Query Routing: Managing Compute and Economic Scale

Deploying massive multi-agent swarms inherently introduces substantial computational expenditure. Generating dozens of parallel expert responses and passing them through debate and aggregation layers uses significant GPU resources. Consequently, it is economically and computationally unviable to trigger the full multi-agent swarm for every user query.³⁵

Production-grade enterprise systems must implement intelligent query routing cascades to act as the primary gatekeepers for inference.

RouteLLM and FrugalGPT Cascades

Frameworks such as FrugalGPT and RouteLLM implement hierarchical cascades that evaluate prompts before any generation occurs.⁴² Utilizing lightweight classification algorithms—such as similarity matching with a Bradley-Terry model or BERT classifiers—the routing layer evaluates the semantic complexity of the incoming prompt and predicts the necessary computational

weight required to satisfy it.⁴²

If the routing model determines the task is relatively trivial or falls within a high-confidence threshold, it directs the query immediately to an inexpensive, high-speed single monolithic model.⁴³ The query is only escalated to the expensive, decentralized 25B orchestrator and its 7B multi-agent swarm if the lightweight model fails a predefined reliability threshold or explicitly triggers a complexity flag.⁴³

The benchmark numbers validating these routing strategies are striking. Empirical deployment of RouteLLM methodologies, particularly utilizing matrix factorization routers, demonstrates that enterprises can achieve over 85% reductions in API and compute costs while rigorously maintaining 95% of the performance capability of frontier models like GPT-4.⁴² Specifically, on the MMLU and GSM8K reasoning benchmarks, intelligent routing reduced inference costs by 45% and 35%, respectively, requiring only 14% to 26% of calls to actually hit the most expensive model tiers.⁴² By treating AI model selection as a dynamic, real-time "air traffic control" system, the architecture ensures that the massive computational expenditure and latency of a parallel swarm are reserved exclusively for high-stakes, logically dense problem sets that demand deep intelligence.⁴²

Decentralized Physical Topologies: Hosting Models on Separate Machines

Realizing the vision of a multi-agent swarm where specialized 7B models and 25B orchestrators run in parallel requires abandoning centralized cloud infrastructure. Deploying dozens of concurrent LLMs using traditional data center GPU clusters poses severe bottlenecks regarding API quota limits, massive capital expenditure, and thermal density limits. Consequently, achieving a true, high-speed swarm necessitates distributing the inference workload across discrete, physically separate machines connected via local area networks (LAN) or the open internet.⁹

BitTorrent-Style Inference and Swarm Parallelism

Decentralized frameworks, most notably Petals, have pioneered the use of Distributed Hash Tables (DHT) to coordinate inference across massive pools of volunteer edge hardware.⁴⁷ In this model, the foundational parameter layers of a massive model (such as Llama 3.1 405B, Mixtral 8x22B, or BLOOM 176B) are intricately partitioned and sharded across disparate consumer-grade GPUs spread globally across the internet.⁴⁷

When a prompt is injected into the Petals network, it executes a pipeline of peer-to-peer computations, passing the hidden states sequentially from one consumer node to the next.⁴⁷ The Decentralized Mixture-of-Experts (dMoE) layer was explicitly designed for training and inferencing with vast amounts of highly unreliable consumer-grade devices.¹¹ Utilizing PyTorch frameworks and algorithms like DeDLOC and All-Reduce SGD, the dMoE layer accommodates node failure gracefully; if a volunteer's GPU drops offline, the DHT requires at most

$O(k \log N)$ queries to locate alternate experts and seamlessly continue processing.¹¹ While this successfully democratizes access to frontier-class models, running purely sequential layers across public internet networks introduces substantial latency due to bandwidth constraints, capping generation at roughly 4 to 6 tokens per second.⁹

Local AI Supercomputers: The Exo Architecture

To support the high-speed requirements of a multi-agent orchestrator managing multiple sub-agents simultaneously, enterprise edge architectures are shifting toward localized, high-bandwidth cluster models, exemplified by the open-source Exo framework.¹⁰ Exo facilitates the deployment of models comprising hundreds of billions of parameters (capable of running DeepSeek 671B models) across heterogeneous consumer hardware, seamlessly clustering Apple Silicon MacBooks, Linux workstations, and even ARM-based Raspberry Pis within a single local network.¹⁰

The system completely bypasses the need for a rigid master-worker topology; devices automatically discover one another via peer-to-peer protocols and partition model shards dynamically based on the available Unified Memory or VRAM of each physical node.⁵² Utilizing the MLX distributed framework optimized for Apple Silicon, Exo processes data using a ring topology.⁵² To circumvent the latency constraints that plague internet-based decentralized networks, advanced Exo edge clusters employ RDMA (Remote Direct Memory Access) over Thunderbolt 5 or high-bandwidth ethernet connections.¹⁰ By enabling direct memory-to-memory data transfers without taxing the host CPU, these localized setups achieve a 99% reduction in transmission latency compared to traditional TCP/IP sockets, enabling parallel agents hosted on separate machines to communicate and generate tokens at speeds rivaling centralized cloud environments.¹⁰

Wireless Distributed Mixture-of-Experts (WDMoE)

The furthest extrapolation of this decentralized topology is the deployment of independent agent swarms over wireless Edge and Internet-of-Things (IoT) infrastructure.⁵⁴ Deep research into Wireless Distributed Mixture-of-Experts (WDMoE) establishes a paradigm where the 25B orchestrator (or the MoE attention mechanism) resides on a central base station (BS), while the independent 7B expert modules are hosted dynamically on mobile devices or edge nodes—such as NVIDIA Jetson hardware kits—communicating via standard WiFi or 5G channels.⁵⁵

In the WDMoE framework, as an orchestrator query arrives, a highly optimized, latency-aware bandwidth allocation algorithm jointly optimizes two critical variables: the selection of the most capable expert agent for the task and the current transmission state of the network.⁵⁵ The central base station routes token embeddings exclusively to edge devices that exhibit both the requisite computational parameters and the strongest wireless signal strength.⁵⁵ Theoretical simulations combined with practical hardware testbed experiments validate this architecture, demonstrating an average reduction in end-to-end latency by 45.75% on standardized

reasoning datasets (such as PIQA) without any measurable degradation in LLM capability.⁵⁵ This proves that intelligent, decentralized parallel routing can entirely negate the inherent instability of wireless networks, paving the way for multi-agent swarms operating across decentralized urban IoT arrays.⁵⁶ Further hardware improvements, such as utilizing bi-directional optical transceivers and optical breakout cables in physical networking (e.g., MixNet), allow decentralized domains to scale throughput on par with traditional non-blocking Fat-tree topologies used in supercomputers.⁵⁸

Combined Solution Architecture and Use Cases

The final ReEnvision AI architectural solution merges hierarchical MoA orchestration, M1-Parallel execution, and DSP-guided speculative actions into a unified decentralized intelligence mesh.

High-Complexity Software Engineering: By employing this architecture in autonomous coding environments, swarms utilize Git Worktrees to execute codebase modifications entirely in parallel.³ Speculative actions empower the Product Manager orchestrator to draft code implementations and predict test suite outcomes locally, rapidly culling invalid logic branches without waiting for heavy external environment feedback. This transforms sequential coding loops into highly predictive, parallel synthesis, fundamentally accelerating modern engineering workflows.

Edge-Based and Wireless Deployments: In environments reliant on 5G or localized IoT edge devices, transferring massive context windows is heavily constrained by bandwidth.⁵⁶ By combining speculative prefetching with WDMoE and Exo's decentralized networking, the base station orchestrator can pre-cache expert parameters and speculatively execute agent responses over lightweight protocols. This successfully reduces end-to-end inference latency by over 45% while preserving logic fidelity.¹⁰

Overall Value Proposition:

The value derived from layering speculative decoding directly onto dynamic agent routing shatters the traditional AI scaling laws. It allows organizations to orchestrate an ecosystem of specialized, decentralized 7B parameter models that operate at high velocity, circumventing the exorbitant capital expenditure, power density, and centralized control of monolithic 100B+ parameter data center clusters.

Synthesizing the Future of Distributed Intelligence

The synthesis of these diverse research vectors—Mixture-of-Agents orchestration, statistical sampling and voting ensembles, hierarchical task delegation, and decentralized physical networking—projects a stark departure from the current, capital-intensive trajectory of artificial intelligence development.

The prevailing industry assumption that the path to Artificial General Intelligence (AGI) strictly requires the consolidation of tens of thousands of GPUs inside centralized, gigawatt data

centers to train singular, monolithic models is directly challenged by the structural efficacy of decentralized swarm architectures.⁵⁹ The mathematical and empirical evidence confirms that an organized ensemble of smaller, weaker agents, equipped with independent context windows and governed by a competent product-manager orchestration layer, can recursively synthesize, verify, debate, and vote their way to conclusions that consistently match or exceed the intelligence capabilities of heavily capitalized monolithic systems.²

By physically hosting these 7B and 25B models across disparate, localized hardware—ranging from globally decentralized consumer PC networks via Petals, local high-bandwidth unified memory clusters via Exo, to 5G edge devices managed by WDMoE algorithms—organizations can construct formidable AI capabilities that are highly resilient, intrinsically private, and fundamentally unbound by the token constraints of a single monolithic attention mechanism.⁹

The transition from single-model prompting to dynamic, multi-agent swarm orchestration represents the true maturation of generative AI. It evolves the technology from a conversational novelty into a robust, parallelized computational substrate. Through the disciplined application of strict agent isolation, consensus-free debate tracking, and aggressive parallel subtask execution via frameworks like M1-Parallel, decentralized multi-agent swarms will undoubtedly define the next generation of scalable, high-speed machine reasoning.

Works cited

1. Mixture of Agents (MoA) | LLM Knowledge Base - Promptmetheus, accessed May 17, 2026, <https://promptmetheus.com/resources/llm-knowledge-base/mixture-of-agents-moa>
2. Mixture-of-Agents Enhances Large Language Model Capabilities - arXiv, accessed May 17, 2026, <https://arxiv.org/html/2406.04692v1>
3. Claude Swarm Mode Complete Guide: 5 Steps to Master the New Paradigm of Multi-agent Collaborative Development, accessed May 17, 2026, <https://help.apiyi.com/en/claude-code-swarm-mode-multi-agent-guide-en.html>
4. Learning Unified Agent-Tool Orchestration with Parallel Subtask Decomposition - arXiv, accessed May 17, 2026, <https://arxiv.org/html/2604.17009v1>
5. SWARM ARENA - ETHGlobal, accessed May 17, 2026, <https://ethglobal.com/showcase/swarm-arena-qj1cn>
6. Agent Swarms: The Future of Decentralized AI | by manav ghosh | Medium, accessed May 17, 2026, <https://medium.com/@manavghosh/agent-swarms-the-future-of-decentralized-ai-c153feb33d4e>
7. NeurIPS Poster Multi-Agent Debate for LLM Judges with Adaptive Stability Detection, accessed May 17, 2026, <https://neurips.cc/virtual/2025/poster/117644>
8. Free-MAD: Consensus-Free Multi-Agent Debate | OpenReview, accessed May 17, 2026, <https://openreview.net/forum?id=46jbtZZWen>
9. Why is nobody talking about decentralized inference? (BitTorrent-style for LLMs) - Reddit, accessed May 17, 2026,

- https://www.reddit.com/r/LLM/comments/1rbc6l1/why_is_nobody_talking_about_decentralized/
10. Goodbye AI Cloud Bills... Exo Runs AI on Your Own Devices - YouTube, accessed May 17, 2026, <https://www.youtube.com/watch?v=L8468vMTX2s>
 11. Towards Crowdsourced Training of Large Neural Networks using Decentralized Mixture-of-Experts, accessed May 17, 2026, https://papers.neurips.cc/paper_files/paper/2020/file/25ddc0f8c9d3e22e03d3076f98d83cb2-Paper.pdf
 12. Mixture of Agents (MoA) on GPU Cloud: Deploy Multi-LLM Voting Architectures (2026 Guide) | Spheron Blog, accessed May 17, 2026, <https://www.spheron.network/blog/mixture-of-agents-gpu-cloud/>
 13. What is mixture of experts? | IBM, accessed May 17, 2026, <https://www.ibm.com/think/topics/mixture-of-experts>
 14. [2512.18452] Secret mixtures of experts inside your LLM - arXiv, accessed May 17, 2026, <https://arxiv.org/abs/2512.18452>
 15. Mixture of Experts Explained - Hugging Face, accessed May 17, 2026, <https://huggingface.co/blog/moe>
 16. What Is Mixture of Experts (MoE) and How It Works? | NVIDIA Glossary, accessed May 17, 2026, <https://www.nvidia.com/en-us/glossary/mixture-of-experts/>
 17. Distributed Mixture-of-Experts and Expert Parallelism - Bruno Magalhaes, accessed May 17, 2026, <https://brunomaga.github.io/Mixture-of-Experts>
 18. MIXTURE-OF-AGENTS ENHANCES LARGE LANGUAGE MODEL CAPABILITIES - ICLR Proceedings, accessed May 17, 2026, https://proceedings.iclr.cc/paper_files/paper/2025/file/5434be94e82c54327bb9dcaf7fca52b6-Paper-Conference.pdf
 19. Mixture-of-Agents (MoA): Improving LLM Quality through Multi-Agent Collaboration, accessed May 17, 2026, <https://a-nikishaev.medium.com/mixture-of-agents-moa-improving-llm-quality-through-multi-agent-collaboration-eb0bcbbdbe9f>
 20. Mixture Of Agents Using Groq. What is an Agent ? | by Plaban Nayak | The AI Forum, accessed May 17, 2026, <https://medium.com/the-ai-forum/mixture-of-agents-using-groq-09e2d46f3ce7>
 21. Together Mixture of Agents (MoA), accessed May 17, 2026, <https://docs.together.ai/docs/mixture-of-agents>
 22. togethercomputer/MoA: Together Mixture-Of-Agents (MoA) – 65.1% on AlpacaEval with OSS models - GitHub, accessed May 17, 2026, <https://github.com/togethercomputer/moa>
 23. Multi-LLM Verification for Question Answering under Conflicting Contexts - ACL Anthology, accessed May 17, 2026, <https://aclanthology.org/2025.ranlp-1.116.pdf>
 24. More Agents Is All You Need: A Deep Dive into Scaling Multi-Agent Systems | by Maheedhar, accessed May 17, 2026, <https://medium.com/@mahadheer/more-agents-is-all-you-need-a-deep-dive-into-scaling-multi-agent-systems-b0d3a83b47c1>
 25. [2402.05120] More Agents Is All You Need - arXiv, accessed May 17, 2026, <https://arxiv.org/abs/2402.05120>

26. More Agents Is All You Need - arXiv, accessed May 17, 2026, <https://arxiv.org/html/2402.05120v1>
27. More Agents Is All You Need | OpenReview, accessed May 17, 2026, <https://openreview.net/forum?id=bgzUSZ8aeg>
28. More Agents Is All You Need - YouTube, accessed May 17, 2026, <https://www.youtube.com/watch?v=8UvYm9HAQUQ>
29. [2604.17009] Small Model as Master Orchestrator: Learning Unified Agent-Tool Orchestration with Parallel Subtask Decomposition - arXiv, accessed May 17, 2026, <https://arxiv.org/abs/2604.17009>
30. How Kimi, Cursor, and Chroma Train Agentic Models with RL - Philschmid, accessed May 17, 2026, <https://www.philschmid.de/kimi-composer-context>
31. Voting or Consensus? Decision-Making in Multi-Agent Debate - ACL Anthology, accessed May 17, 2026, <https://aclanthology.org/2025.findings-acl.606.pdf>
32. Claude Code's Hidden Multi-Agent System - Emergent Minds | paddo.dev, accessed May 17, 2026, <https://paddo.dev/blog/claude-code-hidden-swarm/>
33. What is the Claude Code Swarm feature? Agent Teams, explained - Cyrus, accessed May 17, 2026, <https://www.atcyrus.com/stories/what-is-claude-code-swarm-feature>
34. Claude Code Swarms: Multi-Agent AI Coding Is Here - Zen van Riel, accessed May 17, 2026, <https://zenvanriel.com/ai-engineer-blog/claude-code-swarms-multi-agent-orchestration/>
35. Multi-Agent Systems: Design Patterns and Orchestration - Tetrade, accessed May 17, 2026, <https://tetrade.io/learn/ai/multi-agent-systems>
36. Optimizing Sequential Multi-Step Tasks with Parallel LLM Agents - arXiv, accessed May 17, 2026, <https://arxiv.org/html/2507.08944v1>
37. [2507.08944] Optimizing Sequential Multi-Step Tasks with Parallel LLM Agents - arXiv, accessed May 17, 2026, <https://arxiv.org/abs/2507.08944>
38. Aligned Delegation with Performance Guarantees for Multi-Agent LLM Reasoning - arXiv, accessed May 17, 2026, <https://arxiv.org/html/2602.00127v1>
39. [2502.19130] Voting or Consensus? Decision-Making in Multi-Agent Debate - arXiv, accessed May 17, 2026, <https://arxiv.org/abs/2502.19130>
40. Improving LLM Reasoning with Multi-Agent Tree-of-Thought Validator Agent - arXiv, accessed May 17, 2026, <https://arxiv.org/html/2409.11527v2>
41. Learning to Reason Across Parallel Samples for LLM Reasoning - OpenReview, accessed May 17, 2026, <https://openreview.net/forum?id=e95povsLp5>
42. Intelligent LLM Routing: How Multi-Model AI Cuts Costs by 85%, accessed May 17, 2026, <https://www.swfte.com/blog/intelligent-llm-routing-multi-model-ai>
43. Towards Efficient Multi-LLM Inference: Characterization and Analysis of LLM Routing and Hierarchical Techniques - arXiv, accessed May 17, 2026, <https://arxiv.org/html/2506.06579v1>
44. LLM Routing - Intuitively and Exhaustively Explained - Towards Data Science, accessed May 17, 2026, <https://towardsdatascience.com/llm-routing-intuitively-and-exhaustively-explained-5b0789fe27aa/>

45. LLM Routing: The \$10M Cost Trap Most Orgs Miss - DEV Community, accessed May 17, 2026, https://dev.to/dr_hernani_costa/llm-routing-the-10m-cost-trap-most-orgs-miss-5d0l
46. Distributed Local LLM Swarm using multiple computers instead of one powerful GPU : r/LocalLLaMA - Reddit, accessed May 17, 2026, https://www.reddit.com/r/LocalLLaMA/comments/1seli2p/distributed_local_llm_swarm_using_multiple/
47. Petals – Run LLMs at home, BitTorrent-style, accessed May 17, 2026, <https://petals.dev/>
48. Parallax: Efficient LLM Inference Service over Decentralized Environment - arXiv, accessed May 17, 2026, <https://arxiv.org/html/2509.26182v1>
49. GitHub - bigscience-workshop/petals: Run LLMs at home, BitTorrent-style. Fine-tuning and inference up to 10x faster than offloading, accessed May 17, 2026, <https://github.com/bigscience-workshop/petals>
50. 1 Introduction - arXiv, accessed May 17, 2026, <https://arxiv.org/html/2602.08019v1>
51. exo-explore/exo: Run frontier AI locally. - GitHub, accessed May 17, 2026, <https://github.com/exo-explore/exo>
52. Exo 1.0 means you can cluster mac studios for large models... can I cluster macbooks? : r/LocalLLaMA - Reddit, accessed May 17, 2026, https://www.reddit.com/r/LocalLLaMA/comments/1pqfs1l/exo_10_means_you_can_cluster_mac_studios_for/
53. Distributed approach similar to EXO/LocalAI · containers ramalama · Discussion #1092 - GitHub, accessed May 17, 2026, <https://github.com/containers/ramalama/discussions/1092>
54. [PDF] WDMoE: Wireless Distributed Mixture of Experts for Large, accessed May 17, 2026, <https://www.semanticscholar.org/paper/33177181c599e06535ca6714801af0da66585eb2>
55. WDMoE: Wireless Distributed Mixture of Experts for Large Language Models - arXiv, accessed May 17, 2026, <https://arxiv.org/html/2411.06681v1>
56. WDMoE: Wireless Distributed Large Language Models with Mixture of Experts - arXiv, accessed May 17, 2026, <https://arxiv.org/html/2405.03131v1>
57. WDMoE: Wireless Distributed Large Language Models with Mixture of Experts | Request PDF - ResearchGate, accessed May 17, 2026, https://www.researchgate.net/publication/389773526_WDMoE_Wireless_Distributed_Large_Language_Models_with_Mixture_of_Experts
58. MixNet: A Runtime Reconfigurable Optical-Electrical Fabric for Distributed Mixture-of-Experts Training - arXiv, accessed May 17, 2026, <https://arxiv.org/html/2501.03905v3>
59. Language Model Agents in 2025: Society Mind Revisited | by iSolutions - Medium, accessed May 17, 2026, <https://isolutions.medium.com/language-model-agents-in-2025-897ec15c9c42>
60. Help Us Build the Future of Decentralized AI - From 1 Node to a Global Network! · Issue #848 · exo-explore/exo - GitHub, accessed May 17, 2026,

<https://github.com/exo-explore/exo/issues/848>

61. GitHub - cabelo/multicortex-exo: Run your own AI cluster at home with everyday devices, accessed May 17, 2026, <https://github.com/cabelo/multicortex-exo>