

# It is Time to Virtualize Foundation Models with a Self-evolving Operating System Layer

Suparna Bhattacharya<sup>1</sup>, Tarun Kumar<sup>1\*</sup>, Cong Xu<sup>1\*</sup>, Satish Kumar Mopur<sup>1</sup>, Jiahao Li<sup>1</sup>, Ashish Mishra<sup>1</sup>, Aalap Tripathy<sup>1</sup>, Annmary Justine Koomthanam<sup>1</sup>, Martin Foltin<sup>1</sup>, Ian Foster<sup>2,3\*</sup>  
<sup>1</sup>Hewlett-Packard Enterprise      <sup>2</sup>Argonne National Laboratory      <sup>3</sup>University of Chicago

## Abstract

AI applications have shifted from single, mono-lithic foundation models (FM) to compound agentic systems. Yet today's stacks remain fragmented: even as protocols (e.g., MCP, A2A) ease tool/agent connectivity, each framework embeds an implicit runtime for state, memory, budgets, and trust, making behavior non-portable and governance brittle. It mirrors computing before operating systems, when every program re-implemented basic services. In this position paper, we argue that foundation models should be treated as virtualized system resources, analogous to CPUs, memory, and accelerators, and managed through an operating-system-like layer. We propose a Foundation Model Operating System (FMOS) that introduces a virtualization boundary between applications and heterogeneous physical foundation models. FMOS exposes Virtual Foundation Models (VFM) as stable abstractions that compose model processors, contextual knowledge and memory, and reasoning and trust mechanisms behind a unified interface. Like the human brain switching between fast intuition and slow deliberation, the FMOS learns when to intervene and when to let inference proceed directly and continuously adapting its policies based on operational experience. By reframing foundation model deployment as a systems problem rather than a purely learning-centric one, we outline how OS-inspired abstractions can enable safe, efficient, and portable FM-based applications across heterogeneous models and infrastructures.

**Keywords:** foundation models, abstractions, virtualization

## ACM Reference Format:

Suparna Bhattacharya<sup>1</sup>, Tarun Kumar<sup>1\*</sup>, Cong Xu<sup>1\*</sup>, Satish Kumar Mopur<sup>1</sup>, Jiahao Li<sup>1</sup>, Ashish Mishra<sup>1</sup>, Aalap Tripathy<sup>1</sup>, Annmary

\*Both authors contributed equally to this research.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from [permissions@acm.org](mailto:permissions@acm.org).  
*AgenticOS ASPLOS 2026, Pittsburgh, PA, USA*

© 2026 Copyright held by the owner/author(s). Publication rights licensed to ACM.

ACM ISBN 978-x-xxxx-xxxx-x/YYYY/MM  
<https://doi.org/10.1145/nnnnnnn.nnnnnnn>

Justine Koomthanam<sup>1</sup>, Martin Foltin<sup>1</sup>, Ian Foster<sup>2,3</sup>. 2026. It is Time to Virtualize Foundation Models with a Self-evolving Operating System Layer. In *Proceedings of AgenticOS Workshop @ ASPLOS 2026 (AgenticOS ASPLOS 2026)*. ACM, New York, NY, USA, 17 pages. <https://doi.org/10.1145/nnnnnnn.nnnnnnn>

## 1 Introduction

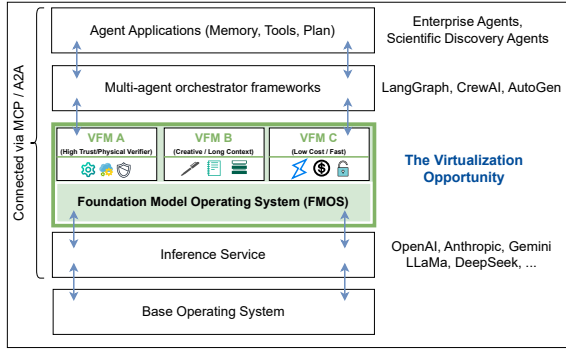
Modern computing systems increasingly rely on Foundation Models (FMs) as shared computational resources across diverse applications, ranging from interactive agents to data analytics pipelines. However, today's FM deployments lack a principled systems abstraction: applications bind directly to specific services (agent runtimes, memory, guardrails), configurations, and inference stacks. As a result FM-based systems conflate application logic with model behavior, leading to brittle designs that are difficult to secure, optimize, or evolve. This tight coupling mirrors pre-virtualization computing, where applications were bound to physical hardware, inhibiting scalability, safety, and evolution.

Virtualization succeeded in systems by introducing a stable abstraction boundary between applications and physical resources, enabling isolation, multiplexing, and controlled sharing. We argue that FMs require a similar abstraction layer. In this analogy, models act as virtual compute devices, prompts and contexts resemble virtual address spaces, and agent capabilities correspond to privileged operations.

### 1.1 Missing System Layer for Agentic AI

Compound agentic systems are quickly becoming the dominant deployment pattern for foundation models, yet the stack lacks a *system layer* analogous to what operating systems provided for traditional software: stable execution semantics, shared services, and enforceable governance.

**The Current Fragmentation: Proliferating Frameworks Without Common Foundations** Despite rapid progress toward compound FM systems, the supporting software stack remains fractured. In practice, adopting an orchestration framework (e.g., LangChain, Claude Code) means inheriting a framework-specific *runtime* that fixes execution semantics: how prompts are assembled, how state is represented, how tool outputs are retained, how failures are retried, how budgets are tracked, and how safety checks are applied. Because most semantics are not exposed as portable interfaces, two logically similar agents can exhibit materially different behavior, reliability, and governance properties across harnesses. Protocol efforts (MCP, A2A) [13,



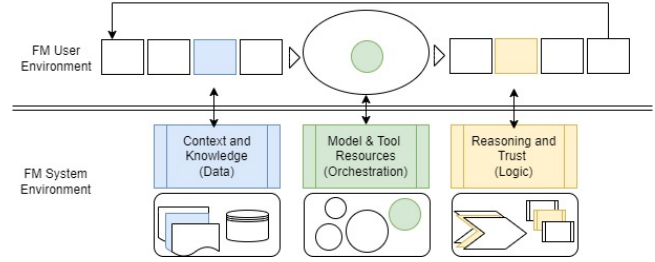
**Figure 1.** FMOS decouples agent logic from model resources, enabling scalable, secure agents. VFMs provide an illusion of infinite dedicated resources to applications.

26] ease **integration**, but they don’t specify the **system layer**—portable contracts for reliable, governable execution. The key gaps are **state/memory semantics** (identity, persistence, sharing, replay/checkpoints), **observability/auditability** (end-to-end traces: requests, tool calls, provenance, decision paths), **resource governance** (budgets/quotas, multi-tenancy), **trust enforcement** (policy, escalation, safe-by-default mediation), and **model mediation** (routing/caching/materialization with controlled upgrades/rollbacks).

Absent these contracts, teams rebuild a bespoke “mini-platform” inside each framework. System-level optimization and governance remain brittle—a pre-OS pattern where libraries existed, but shared execution semantics did not. Recent LLMOS systems gesture in this direction—AIOS [24], MemGPT [32], and Llumnix [40] explore scheduling, memory virtualization, or serving-level orchestration—but the field still lacks a unifying virtualization boundary that *jointly* governs knowledge, model reasoning, verification, and trust.

**Unifying Workflows and Agentic Loops: The Missing Execution Substrate** Enterprise deployments increasingly mix two execution regimes: (i) *workflow-driven* pipelines with explicit structure and auditability [6], and (ii) *agentic loops* that plan-act-reflect over long horizons [5, 21, 37, 43]. These regimes rarely share a common substrate: workflows assume typed steps, stable boundaries, and predictable logging; loops use open-ended control flow, opportunistic tool use, backtracking, and adaptive context growth.

As shown in Figure 1, the core gap is the absence of portable system-layer contracts for *state, memory, budgets, and mediation* that apply uniformly across execution forms. A system layer with a single virtualization boundary can treat workflows and loops as two schedulable *execution forms* over shared primitives, enforcing the same context/memory management, tracing, resource governance, and trust controls regardless of whether the next step is “run a node” or “plan the next move.” This shared substrate enables hybrid systems that combine workflow stability with loop flexibility.



**Figure 2.** Agentic flows with system-environment services

## 2 Position: Virtual Foundation Models Enabled by an FMOS

We formalize our position: agent deployments now require a distinct system layer—a Foundation Model Operating System (FMOS)—whose primary abstraction is the *Virtual Foundation Model* (VFM). FMOS is a virtualization layer between FM-based applications and heterogeneous model execution backends (refer to Figure 1). Analogous to a hypervisor, FMOS intercepts application requests, enforces capability constraints, manages shared state, and dynamically maps virtual model interfaces to physical model instances. This separation allows applications to be written against a stable virtual FM interface while enabling the system to optimize placement, enforce policies, and evolve underlying models independently. A VFM presents applications with the illusion of a dedicated, trustworthy FM instance with effectively unbounded capabilities, while the FMOS mediates how knowledge is retrieved and updated, how models and reasoning are tailored to tasks, and how resources are allocated.

**Self-evolution as a core principle.** Unlike traditional OS virtualization, which preserves fidelity to underlying hardware, FMOS virtualization is designed for *progressive quality gain through self-evolution*. From longitudinal interaction traces, the FMOS learns to update prompts and memories—enabling the system to improve without retraining underlying models. This evolution is managed through versioning, canary deployments, and rollback mechanisms, ensuring that improvements propagate safely across applications while maintaining reproducibility when required.

**Why Now?** Three developments make FMOS timely. First, enterprise agentic deployments have outpaced infrastructure: organizations are scaling beyond pilots but lack adequate governance and integration frameworks. Second, protocol standardization has reached critical mass: MCP and A2A enable tool/agent interoperability, but system-layer abstractions for knowledge management, reasoning verification, and policy enforcement remain missing. Third, compound systems have proven superior to monolithic scaling: state-of-the-art FMs are themselves compound architectures, validating that the future lies in flexible orchestration rather than ever-larger individual models.

This system-layer framing accomplishes three objectives:

1. **Unified services with co-evolution:** Context management, tool orchestration, and verification are handled at the FMOS layer; learned improvements propagate across dependent workloads.
2. **Joint optimization:** The FMOS coordinates knowledge retrieval, model selection, inference quotas, and verification depth as a combined optimization problem—achieving efficiencies that fragmented stacks cannot.
3. **Cross-enterprise reusability:** Domain-specific skills, knowledge augmentations, and reasoning policies can be defined once and reused across applications and teams. Realizing FMOS requires collaborative research to define abstractions, learning mechanisms, and governance frameworks that make FM virtualization practical and principled.

### 3 Virtual FM System Environment Services

We now outline key mechanisms and FMOS system-level services required to realize a virtual FM environment.

#### 3.1 FM Virtualization Conditions

Just as virtual memory allows applications to behave as if they have access to unbounded memory, a virtual FM provides AI applications with the ability to provision and request foundation models with (potentially) unbounded capabilities.

We draw parallels from virtualization requirements in computer architecture [34] (details in Appendix B). A virtual environment (virtual machine) provided by a virtual machine monitor (VMM), is characterized by three key properties: efficiency, resource control (safety) and fidelity (equivalence). An architecture is virtualizable if the set of sensitive operations is a subset of privileged operations, where non-privileged operations execute natively while privileged operations trap if invoked from user environment, thus passing control to the VMM. Analogously, VFMs virtualize higher-level FM capabilities such as knowledge, models, and reasoning. The key properties of efficiency and resource control apply to these resources (e.g., context windows or reasoning operations). The third property is not fidelity, but rather progressive quality gain through self-evolution.

LLMs/FMs have been informally likened to processors that interpret human language. Thus, defining sensitive and privileged operations is more complex and must be *learned* (offline or in-context) based on VFM-controlled capabilities, potentially customized via a model virtualization protocol for FM traps. Unlike hardware VMMs where trap conditions are deterministically defined by the instruction set, FM trap classification must be learned from operational traces and contextual signals such as domain, risk level, budget state, and inferred intent, since the same natural-language request may be innocuous in one context and sensitive in another. Once a trap fires, the activated FMOS capability (e.g., knowledge augmentation) is invoked via a structured interface

carrying a trap type drawn from {knowledge, model, trust, budget}, a context snapshot, and a capability invocation spec.

#### 3.2 VFM System Environment Requirements

A VFM must support a prototypical FM user workflow (Figure 2) comprising input context preparation, model execution passes, and output processing—in a manner that allows for transparent system-level interception and control.

Such an interception enables services in FMOS system environment to support the core intent of simplifying and optimizing FM-based agent applications, while enabling the system’s capacity for self-evolution.

Self-evolution often involves closed-loop, trajectory-driven adaptation and may use both parametric and non-parametric adaptation. Rather than necessarily fine-tuning weights or just adding data, FMOS learns from interaction traces and updates prompts, policies, and structured memories for the VFM. The system improves behavior through curated, FMOS-managed learning—enabling standardized interception, diagnosis, and safe rollout across models and applications.

#### 3.3 Context Management & Knowledge Augmentation

A FM combines internal (parametric) knowledge acquired during training with (non-parametric) knowledge that it receives as input context (prompts). System environment services control this context both to elicit (with selective focus) what the model knows and to expand (augment) it with external knowledge. Appendix D describes a few capabilities that fall under this category, such as (1) context memory management, (2) knowledge compression and retrieval, and (3) handling knowledge-oriented abstractions for different data modalities. A key challenge in realizing these services is learning to adapt to what is most relevant for the FM application and current context.

**We propound context-management policies as first-class objects.** Today’s agent stacks lack a portable intelligent way to bind an application’s *intent* for context (what must stay in-window, what can be summarized, what must be recoverable) to the *mechanisms* that actually construct prompts and manage tool outputs. This missing interface matters because application developers often *know* which pieces of context are valuable (and when), but cannot express that knowledge to the serving layer. Developers may want to specify rules such as: (1) post-tool-call offload, (2) adaptive agent skills unloading, (3) retain thoughts, prune observations, etc. More details in Appendix D.

These examples share a common structure: each is an application-level *policy* over a system-level *mechanism* (buffering, summarization, pruning, offloading, retrieval). The absence of a policy-to-mechanism mapping forces developers to either (i) accept brittle defaults embedded in a particular harness, or (ii) reimplement context plumbing in application code, undermining composability and reuse.

A declarative context-policy interface enables VFMs to expose stable semantics while allowing the FMOS to learn & optimize the concrete realization of those policies over time.

### 3.4 Reasoning and Trust Augmentation

Reasoning is essential for discovering, evolving or assimilating new knowledge or tools and for ensuring trustworthy FM outputs. Trust in context processing and model behavior cannot be driven by static assessors. The VFM environment services can regulate FM output selection and processing (e.g., through constrained decoding, sampling, representation engineering [50, 51] and relevant tool invocations) with an implicit trust assessment. As described in Appendix D, such capabilities include (1) expanding and managing reasoning resources, (2) switching between reasoning modes, and (3) low-overhead trust and verification, protection. FMOS steering mechanisms help regulate this flow dynamically.

### 3.5 Model Resource Sharing and Orchestration

FM inference demands significant GPU resources, which escalate with inference-time scaling and multi-agent setups. Using smaller FMs or distilled models can alleviate this resource pressure. Underlying model serving platforms typically perform optimizations for requests to a given FM. We propose that VFM system environment services can intercept them [1] and use its higher-level intent to enable deeper co-optimizations and to manage tradeoffs involved in model selection and orchestration. Appendix D describes three capabilities: Scheduling and mapping, model composition and instantiation, and profiling, measurement, and tracing.

### 3.6 Broader Considerations

Beyond the three elements and their interactions, FMOS must also address long-term needs:

**Continual self-evolution:** as new scenarios emerge, with knowledge continuously generated, validated, and refined.

**Continual adoption of latest techniques:** Continuously incorporate improved models, frameworks, and methods so workflows automatically benefit from rapid advances.

**External control:** Support admin-applied controls for global policies (e.g., safety, compliance, and resource bounds), such as privileged instruction hierarchies in FMs [41]

## 4 FMOS Architecture

Drawing on OS virtualization, we present the *VFM* in Figure 3 as the interface exposed to applications, while the *FMOS* mediates access to underlying *physical FMs* (*pFMs*). By default, requests execute on a lightweight *fast path*; when task conditions warrant (e.g., insufficient context, elevated risk, or tight budgets), FMOS triggers learned *traps* to activate a more deliberative *slow path* that performs targeted knowledge augmentation, model routing, and verification. Because this virtualization boundary sits beneath diverse

agent frameworks and APIs, it preserves programming flexibility while enabling shared system policies and learned artifacts to evolve under standard controls (e.g., versioning, canaries, rollback). This design emphasizes three coupled challenges a virtualized VFM layer must address:

**Data and Knowledge.** Represent, expand, and continually update multimodal knowledge while managing context/memory under finite budgets.

**Trust and Reasoning.** Adapt reasoning depth, domain-specific verification and policy enforcement, user-specified risk levels, and time horizons.

**Efficiency and Adaptability.** Share and switch among a growing pool of models and tools to optimize quality-latency-cost trade-offs under multi-tenant constraints.

FMOS operationalizes these goals through three cooperating subsystems behind the VFM interface: a *Data Agent* that manages context and non-parametric memory; a *Composition Optimizer* that selects, routes, and optionally composes *pFMs*; and a *Trust & Reasoning Agent* that governs escalation, verification, and guardrails. These subsystems are invoked via interception points (e.g., virtual model endpoints, MCP services, framework hooks) and remain optional: a VFM can be realized as a near-direct call to a base model, or as a progressively richer orchestration that optimizes multiple objectives. FMOS is grounded in three design principles:

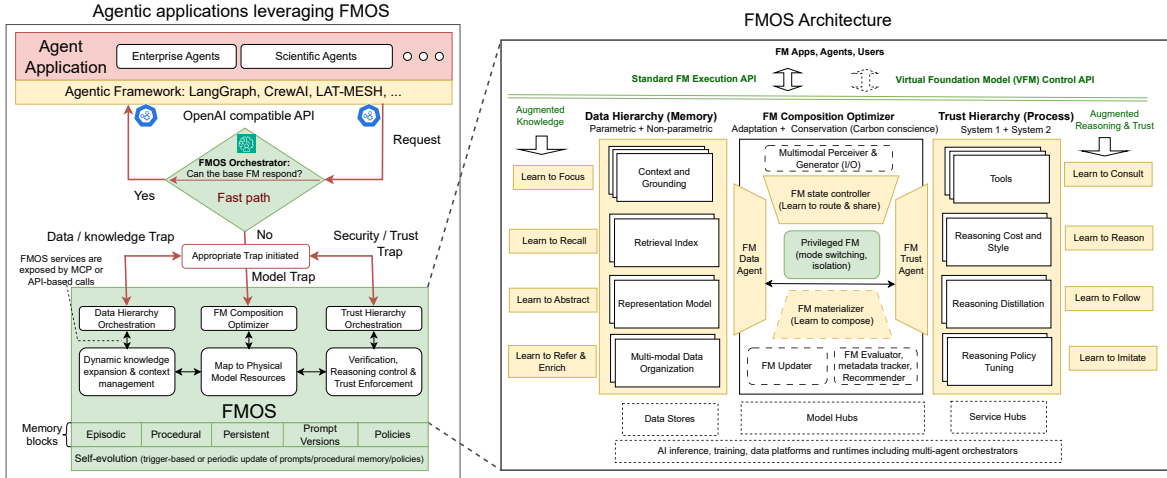
**1. Virtualize capabilities behind stable semantics:** Expose a stable behavioral contract (state/memory persistence, budget, trust) where applications rely on invariants than implementation. FMOS evolves mechanisms (routing, context tiering, verification) as long as a contract holds.

**2. Demand-driven orchestration:** keep the common case fast, and escalate only via explicit traps when quality, budget, or trust requirements require it.

**3. Ecosystem compatibility:** integrate with existing APIs, agent frameworks, and protocols so improvements propagate without forcing application rewrites.

### 4.1 Data Agent: Dynamic Knowledge Expansion and Context Management

The Data Agent functions as a virtual memory subsystem, responsible for abstraction, compaction, and controlled exposure of contextual state. It implements the VFM’s *knowledge plane*: it mediates what enters the model’s active context and what is persisted externally, enabling reusable knowledge augmentation under explicit context and cost budgets (Section 3). Concretely, it selects, transforms, and retrieves multimodal artifacts (e.g., documents, tables, code, images, and time series) so that model execution remains grounded while state and evidence remain recoverable across steps and sessions. Like virtual memory, it maintains a *memory hierarchy* organized by functional role rather than raw latency: a small, fast “working set” (prompt context) backed by lower tier (files, vector or graph databases) that store episodic, procedural, and semantic artifacts. Paging and caching are



**Figure 3.** (left) FMOS intercepts the execution flow where needed and orchestrates an optimized execution path. (right) An expanded view of common fundamental capabilities offered in FMOS.

*content-aware*: retrieved items may be summarized, schema-fied, embedded, or replaced with retrieval handles, preserving access semantics while keeping the active window within budget. This hierarchy is driven by learnable policies:

**Learn to focus:** produce task-conditioned context slices and grounding (e.g., self-RAG [7] or adaptive zooming [49]) to steer generation, filter irrelevant material; when needed, iterate between generation, augmentation, and refinement.

**Learn to recall:** index and retrieve previously encountered artifacts via suitable representations (embeddings, graphs, files, or representation-engineered features [9, 50]) to support reuse without repeatedly re-deriving the same evidence.

**Learn to abstract:** when retrieval alone is insufficient, train or finetune specific components (e.g., domain embedding models, parameter-efficient adapters, distilled auxiliaries) to improve representations and extend coverage to additional modalities such as time series [23].

**Learn to refer and enrich:** pre-process, synthesize, and align heterogeneous sources so they become retrieval- and adaptation-ready; operate efficiently and, where necessary, approximately and on-demand at multiple granularities.

By virtualizing the movement of knowledge between short-term context and long-term stores, the Data Agent gives VFM the practical illusion of boundless, continually improving access to essential information.

#### 4.2 FM Composition Optimizer

The FM Composition Optimizer maps VFM calls to *physical* model resources under explicit quality–latency–cost–policy constraints. It supports both (i) *static provisioning* of a fit-for-purpose model bundle for a class of workloads, and (ii) *dynamic scheduling* at runtime (routing, caching, sharing) as task demands and resource conditions change. Decisions are informed by continual measurement and task-specific performance traces rather than fixed, framework-level heuristics.

**FM materializer (Learn to Compose):** provisions the execution substrate for a VFM by operations like selecting

and distilling models from a candidate pool to meet capability and deployment constraints.

**FM controller (Learn to Route):** performs runtime model routing [20], caching, and sharing across instantiated models, while enforcing instruction-privileges, security boundaries and respecting SLOs (latency, throughput, and cost).

**FM updater:** manages controlled evolution of materialized models via continual (un)learning and editing, with versioning and rollback consistent with FMOS governance.

Operationally, this subsystem plays the role of a scheduler-plus-hypervisor for model capability: it decides *what* physical models a VFM is backed by and *when* to switch, reuse, or refresh them as workloads and the model ecosystem evolve.

#### 4.3 Trust & Reasoning Hierarchy

The FMOS Trust & Reasoning Agent mediates how a VFM allocates deliberation and verification under explicit safety, cost, and latency budgets. In routine cases it stays on a lightweight “fast” path; when tasks become ambiguous, high-stakes, or policy-sensitive, it escalates to slower reasoning, tool-assisted checks, and stricter guardrails, and logs the outcomes to improve future decisions.

**Learn to consult:** Decide when to invoke tools (e.g., search, checkers, human-in-loop) and how to integrate their outputs as evidence avoiding uncontrolled context expansion with associated verification.

**Learn to reason:** Control the reasoning *cost* and *style* (planning, decomposition, reflection), including switching between fast heuristics and deliberate search. When slow-path reasoning is validated, distill reusable reasoning templates or policies to reduce future compute for similar cases.

**Learn to follow:** When constraints and domain logic are stable and recurring, then update prompts, policies [19], or lightweight adapters so that common rules and responsible behaviors are enforced.

**Learn to imitate:** When gaps reflect missing coverage in the underlying FMs, trigger broader upgrades (e.g., continual pretraining, or specialized reasoning), gated by evaluation to manage catastrophic forgetting.

This hierarchy mirrors OS protection mechanisms: as actions become more privileged or risky, they trigger progressively stronger validation, scaling computational commitment with task criticality and trust requirements.

#### 4.4 Minimal Overhead with Maximum Capability

FMOS is designed around a fast path. By default, requests pass through the VFM interface with minimal interception beyond lightweight accounting and tracing. When learned “traps” fire (e.g., uncertainty, policy sensitivity, budget pressure, or anomaly signals), FMOS activates only the required subset of capabilities—escalating context retrieval, switching models, or deepening verification—and then returns execution to the fast path. This preserves low latency and cost while allowing the same VFM endpoint to provide stronger guarantees when needed, and keeps context management, routing, and trust policies transparent to applications. Importantly, FMOS does not subsume inference-layer optimizations such as KV-cache management or paged attention. Its fast path is transparent to the serving stack, intercepting only at the semantic VFM boundary and passing context policies as hints to underlying serving infrastructure. In Appendix A, we illustrate how FMOS mechanisms are applicable in practice in two typical application scenarios. In Appendix B we provide an analytical perspective of how FMOS virtualization enables co-evolution, co-optimization, and reuse.

## 5 Alternative Views

*FMs will become so good at everything that we will no longer need to augment them.* FMs are expanding across modalities, context length, and large reasoning models (LRMs) [10, 44], solving hard problems through inference-time scaling. It introduces new challenges including reasoning cost, “overthinking” (Appendix C), and trustworthiness [16]. Even as FMs/LRMs improve, longer inference-time reasoning traces alone cannot gather new evidence or adapt behavior in dynamic tasks, so “thinking more” is insufficient without interaction [38]. Recent agentic-reasoning work instead treats capability as a plan-act-learn loop with tools, feedback, and memory—so augmentation remains fundamental rather than optional [42]. (more details are given in the Appendix E)

*Agent frameworks like LangChain, AutoGen will encompass everything, when combined with query and pipeline optimization techniques for compound AI systems.* Agent frameworks help with wiring, but they do not provide system-layer guarantees. Empirical evidence shows multi-agent workflows still break on validation, context loss, rollback, and coordination, yielding inconsistent state and poor recovery [11]; developer data likewise highlights orchestration and reliability as persistent bottlenecks [8]. Optimizers inherit these gaps, and

even single agents require OS-like memory virtualization to escape fixed context limits [33]. Hence an FMOS-like layer is needed for portable semantics over state, memory, and trust.

*Model Context Protocol (MCP) and Agent-to-Agent communication protocol advancements will address most challenges.* They standardize how agents *connect*—to tools, resources, and other agents—and this interoperability has catalyzed rapid ecosystem growth. However, these protocols intentionally stop short of specifying *execution semantics* and *governance guarantees*. As deployments scale, teams still must define (and today, reimplement) the system-layer contracts identified in §1.1. Without a shared substrate, those capabilities get bolted onto frameworks or MCP servers in incompatible ways, yielding protocol-compliant but brittle “bloat” and fragmented control.

*The OS and virtualization analogy is misleading as it is a higher level layer and does not directly manage hardware resources.* Traditional OSs virtualize hardware resources while remaining largely unaware of workload intent due to separation-of-concerns principles. For agentic systems, this semantic gap has widened: workloads are expressed in terms of FM instructions, knowledge, and reasoning, where conventional OS abstractions offer limited control for efficiency, safety, and trust [24, 48]. FMOS addresses this by virtualizing higher-level FM operations above the base OS, while still leveraging OS signals and mechanisms to manage environments using OS-inspired principles [24, 32].

## 6 Conclusion and Call to Action

The “LLM as OS” metaphor [18] has gained traction, but its implications for **virtualization**, **isolation**, and **governance** remain underexplored as compound agentic systems become dominant. We argued that this complexity calls for an explicit virtualization layer—a FMOS—that mediates access to pFMs and exposes VFMs as stable system abstractions. By treating FMs as virtualized system resources rather than monolithic learning artifacts, FMOS can improve performance reliability in deployment and provide clearer hooks for safety, trust, and long-term evolvability, echoing the role of operating systems and hypervisors in modern computing.

Turning this position into practice requires a focused systems agenda. The community should converge on a minimal set of FMOS abstractions—VFMs, traps, and system-layer mediation—and build open reference implementations that interpose on existing APIs and agent frameworks without requiring application rewrites. These should be evaluated with benchmarks targeting system-level properties (e.g., context efficiency, cost-quality trade-offs, robustness under evolution, and auditability). Protocols like MCP and A2A support interoperability, but should be paired with explicit governance contracts for privilege, safety, and external control. We invite the systems and machine learning communities to explore these challenges together and chart the future of foundation-model operating systems.

## References

- [1] Reyna Abhyankar, Zijian He, Vikranth Srivatsa, Hao Zhang, and Yiyang Zhang. 2024. InferCept: Efficient Intercept Support for Augmented Large Language Model Inference. *Preprint arXiv:2402.01869* (2024). arXiv:2402.01869 [cs.LG]
- [2] AgentFS. 2026. *AgentFS: Filesystem Isolation for AI Agents*. <https://www.agentfs.ai/>
- [3] AgentFS. 2026. *tursodatabase/agentfs: The filesystem for agents*. <https://github.com/tursodatabase/agentfs>
- [4] Pranjal Aggarwal and Sean Welleck. 2025. L1: Controlling How Long A Reasoning Model Thinks With Reinforcement Learning. *Preprint arXiv:2503.04697* (2025). arXiv:2503.04697 [cs.CL]
- [5] Duncan Anderson. 2025. The 4 Levels of AI Agents: When to Use Workflows vs Autonomous Systems. Blog post, Barnacle Labs. <https://www.barnacle.ai/blog/2025-09-25-agents-intro> Accessed 27 Jan 2026.
- [6] Anthropic. 2024. Building effective agents. Engineering blog. <https://www.anthropic.com/engineering/building-effective-agents> Accessed 27 Jan 2026.
- [7] Akari Asai, Zeqiu Wu, Yizhong Wang, Avirup Sil, and Hannaneh Hajishirzi. 2023. Self-RAG: Learning to retrieve, generate, and critique through self-reflection. *Preprint arXiv:2310.11511* (2023).
- [8] Ali Asgari, Annibale Panichella, Pouria Derakhshanfar, and Mitchell Olsthorn. 2026. What Challenges Do Developers Face in AI Agent Systems? An Empirical Study on Stack Overflow. *arXiv preprint arXiv:2510.25423* (2026).
- [9] Lukasz Bartoszcze, Sarthak Munshi, Bryan Sukidi, Jennifer Yen, Zejia Yang, David Williams-King, Linh Le, Kosi Asuzu, and Carsten Maple. 2025. Representation Engineering for Large-Language Models: Survey and Research Challenges. arXiv:2502.17601 [cs.AI] <https://arxiv.org/abs/2502.17601>
- [10] Maciej Besta, Julia Barth, Eric Schreiber, Ales Kubicek, Afonso Catarino, Robert Gerstenberger, Piotr Nyczyk, Patrick Iff, Yueling Li, Sam Houlston, Tomasz Sternal, Marcin Copik, Grzegorz Kwaśniewski, Jürgen Müller, Lukasz Flis, Hannes Eberhard, Hubert Niewiadomski, and Torsten Hoefler. 2025. Reasoning Language Models: A Blueprint. *Preprint arXiv:2501.11223* (2025). arXiv:2501.11223 [cs.AI]
- [11] Edward Y. Chang and Longling Geng. 2025. SagaLLM: Context Management, Validation, and Transaction Guarantees for Multi-Agent LLM Planning. *Proceedings of the VLDB Endowment* 18, 12 (2025), 4874–4886. doi:10.14778/3750601.3750611
- [12] Pekka Enberg and Glauber Costa. 2025. *The Missing Abstraction for AI Agents: The Agent Filesystem*. <https://turso.tech/blog/agentfs> Turso.
- [13] Google Developers. 2025. A2A: A New Era of Agent Interoperability. <https://developers.googleblog.com/en/a2a-a-new-era-of-agent-interoperability/> Last accessed May 21, 2025.
- [14] Shibo Hao, Sainbayar Sukhbaatar, DiJia Su, Xian Li, Zhiting Hu, Jason Weston, and Yuandong Tian. 2024. Training Large Language Models to Reason in a Continuous Latent Space. *Preprint arXiv:2412.06769* (2024).
- [15] Nick Huang. 2025. *How agents can use filesystems for context engineering*. <https://www.blog.langchain.com/how-agents-can-use-filesystems-for-context-engineering/> LangChain.
- [16] Ben Hylak and Latent Space. 2025. o1 isn't a chat model: and that's the point. <https://www.latent.space/p/o1-skill-issue>.
- [17] Daniel Kahneman. 2011. *Thinking, Fast and Slow*. Farrar, Straus and Giroux, New York. [https://www.amazon.de/Thinking-Fast-Slow-Daniel-Kahneman/dp/0374275637/ref=wl\\_it\\_dp\\_o\\_pdT1\\_nS\\_nC?ie=UTF8&colid=1511935NGKJT9&coliid=I3OCESLZCVDL7](https://www.amazon.de/Thinking-Fast-Slow-Daniel-Kahneman/dp/0374275637/ref=wl_it_dp_o_pdT1_nS_nC?ie=UTF8&colid=1511935NGKJT9&coliid=I3OCESLZCVDL7)
- [18] Andrej Karpathy. 2023. LLM as an OS. <https://x.com/karpathy/status/1723140519554105733>.
- [19] Tarun Kumar, Aalap Tripathy, Gayathri Saranathan, Martin Foltin, Suparna Bhattacharya, Scott Hinchley, Donald M Bahls, David Brookshire, Larry Kaplan, and Robert W. Wisniewski. 2026. Infrastructure-Sentinel: Policy Enforced Guardrails for Secure MCP-driven Infrastructure Agents. *Proceedings of the AAAI Conference on Artificial Intelligence* 40, 47 (Mar. 2026), 40295–40301. doi:10.1609/aaai.v40i47.41468
- [20] Tarun Kumar, Cong Xu, Arpit Shah, Baradji Diallo, Martin Foltin, and Suparna Bhattacharya. 2025. Co-optimizing Recommendation and Evaluation for LLM Selection. In *ICLR 2025 Workshop on Foundation Models in the Wild*.
- [21] LangGraph. 2025. Deep Agents. Engineering blog. <https://www.blog.langchain.com/deep-agents/>
- [22] Xiaomin Li, Zhou Yu, Zhiwei Zhang, Xupeng Chen, Ziji Zhang, Yingying Zhuang, Narayanan Sadagopan, and Anurag Beniwal. 2025. When Thinking Fails: The Pitfalls of Reasoning for Instruction-Following in LLMs. *Preprint arXiv:2505.11423* (2025). arXiv:2505.11423 [cs.CL]
- [23] Yuxuan Liang, Haomin Wen, Yuqi Nie, Yushan Jiang, Ming Jin, Dongjin Song, Shirui Pan, and Qingsong Wen. 2024. Foundation Models for Time Series Analysis: A Tutorial and Survey. In *Proceedings of the 30th ACM SIGKDD Conference on Knowledge Discovery and Data Mining (KDD '24)*. ACM, 6555–6565. doi:10.1145/3637528.3671451
- [24] Kai Mei, Zelong Li, Shuyuan Xu, Ruosong Ye, Yingqiang Ge, and Yongfeng Zhang. 2024. AIOS: LLM agent operating system. *Preprint arXiv:2403.16971v3* (2024).
- [25] MiroMind Team. 2025. MiroThinker: Pushing the Performance Boundaries of Open-Source Research Agents via Model, Context, and Interactive Scaling. *arXiv preprint arXiv:2511.11793* (2025). <https://arxiv.org/abs/2511.11793>
- [26] Model Context Protocol. 2025. Supported Clients – Model Context Protocol. <https://modelcontextprotocol.io/clients>. Accessed: 2025-05-22.
- [27] Model Context Protocol Team. 2025. Model Context Protocol Specifications. <https://github.com/modelcontextprotocol/modelcontextprotocol/blob/main/docs/specification/2025-03-26/basic/authorization.mdx>.
- [28] Niklas Muennighoff, Zitong Yang, Weijia Shi, Xiang Lisa Li, Li Fei-Fei, Hannaneh Hajishirzi, Luke Zettlemoyer, Percy Liang, Emmanuel Candès, and Tatsunori Hashimoto. 2025. s1: Simple test-time scaling. *Preprint arXiv:2501.19393* (2025). arXiv:2501.19393 [cs.CL]
- [29] Nexi. 2026. *nexi-lab/nexus: Nexus AI-native filesystem*. <https://github.com/nexi-lab/nexus>
- [30] Nexi. 2026. *Nexus: AI-Native Filesystem for Building Intelligent Agents*. <https://nexi-lab.github.io/nexus/>
- [31] Isaac Ong, Amjad Almahairi, Vincent Wu, Wei-Lin Chiang, Tianhao Wu, Joseph E. Gonzalez, M Waleed Kadous, and Ion Stoica. 2025. RouteLLM: Learning to Route LLMs from Preference Data. In *13th International Conference on Learning Representations*. <https://openreview.net/forum?id=8sSqNntaMr>
- [32] Charles Packer, Sarah Wooders, Kevin Lin, Vivian Fang, Shishir G Patil, Ion Stoica, and Joseph E Gonzalez. 2023. MemGPT: Towards LLMs as operating systems. *Preprint arXiv:2310.08560* (2023).
- [33] Charles Packer, Sarah Wooders, Kevin Lin, Vivian Fang, Shishir G. Patil, Ion Stoica, and Joseph E. Gonzalez. 2024. MemGPT: Towards LLMs as Operating Systems. *arXiv preprint arXiv:2310.08560* (2024).
- [34] Gerald J Popek and Robert P Goldberg. 1974. Formal requirements for virtualizable third generation architectures. *Commun. ACM* 17, 7 (1974), 412–421.
- [35] Ramya Prabhu, Ajay Nayak, Jayashree Mohan, Ramachandran Ramjee, and Ashish Panwar. 2025. vAttention: Dynamic Memory Management for Serving LLMs without PagedAttention. *Preprint arXiv:2405.04437* (2025). arXiv:2405.04437 [cs.LG]
- [36] Alireza Rezazadeh, Zichao Li, Wei Wei, and Yujia Bao. 2024. From Isolated Conversations to Hierarchical Schemas: Dynamic Tree Memory Representation for LLMs. *Preprint arXiv:2410.14052* (2024).

- arXiv:2410.14052 [cs.CL]
- [37] Phil Schmid. 2025. Zero to One: Learning Agentic Patterns. Blog post. <https://www.philschmid.de/agent-ic-pattern> Accessed 27 Jan 2026.
- [38] Junhong Shen, Hao Bai, Lunjun Zhang, Yifei Zhou, Amrith Setlur, Shengbang Tong, Diego Caples, Nan Jiang, Tong Zhang, Ameet Talwalkar, and Aviral Kumar. 2025. Thinking vs. Doing: Agents that Reason by Scaling Test-Time Interaction. arXiv:2506.07976 [cs.LG] doi:10.48550/arXiv.2506.07976
- [39] Tal Shnitzer, Anthony Ou, Mirian Silva, Kate Soule, Yuekai Sun, Justin Solomon, Neil Thompson, and Mikhail Yurochkin. 2023. Large Language Model Routing with Benchmark Datasets. *Preprint arXiv:2309.15789* (2023). arXiv:2309.15789 [cs.CL]
- [40] Biao Sun, Ziming Huang, Hanyu Zhao, Wencong Xiao, Xinyi Zhang, Yong Li, and Wei Lin. 2024. Llumnix: Dynamic scheduling for large language model serving. In *18th USENIX Symposium on Operating Systems Design and Implementation*. 173–191.
- [41] Eric Wallace, Kai Xiao, Reimar Leike, Lilian Weng, Johannes Heidecke, and Alex Beutel. 2024. The Instruction Hierarchy: Training LLMs to Prioritize Privileged Instructions. *Preprint arXiv:2404.13208* (2024). arXiv:2404.13208 [cs.CR]
- [42] Tianxin Wei, Ting-Wei Li, Zhining Liu, Xuying Ning, Ze Yang, Jiuru Zou, Zhichen Zeng, Ruizhong Qiu, Xiao Lin, Dongqi Fu, Zihao Li, Mengting Ai, Duo Zhou, Wenxuan Bao, Yunzhe Li, Gaotang Li, Cheng Qian, Yu Wang, Xiangru Tang, Yin Xiao, Liri Fang, Hui Liu, Xianfeng Tang, Yuji Zhang, Chi Wang, Jiakuan You, Heng Ji, Hanghang Tong, and Jingrui He. 2026. Agentic Reasoning for Large Language Models. arXiv:2601.12538 [cs.AI] doi:10.48550/arXiv.2601.12538
- [43] Simon Willison. 2025. Designing agentic loops. Blog post. <https://simonwillison.net/2025/Sep/30/designing-agent-ic-loops/> Accessed 27 Jan 2026.
- [44] Fengli Xu, Qianyue Hao, Zefang Zong, Jingwei Wang, Yunke Zhang, Jingyi Wang, Xiaochong Lan, Jiahui Gong, Tianjian Ouyang, Fanjin Meng, Chenyang Shao, Yuwei Yan, Qinglong Yang, Yiwen Song, Sijian Ren, Xinyuan Hu, Yu Li, Jie Feng, Chen Gao, and Yong Li. 2025. Towards Large Reasoning Models: A Survey on Scaling LLM Reasoning Capabilities. *Preprint arXiv:2501.09686* (2025). arXiv:2501.09686 [cs.AI]
- [45] Silei Xu, Wenhao Xie, Lingxiao Zhao, and Pengcheng He. 2025. Chain of Draft: Thinking Faster by Writing Less. *Preprint arXiv:2502.18600* (2025).
- [46] Yuan Yang, Siheng Xiong, Ehsan Shareghi, and Faramarz Fekri. 2024. The Compressor-Retriever Architecture for Language Model OS. *Preprint arXiv:2409.01495* (2024).
- [47] Shunyu Yao, Noah Shinn, Pedram Razavi, and Karthik Narasimhan. 2024.  $\tau$ -bench: A Benchmark for Tool-Agent-User Interaction in Real-World Domains. *Preprint arXiv:2406.12045* (2024). arXiv:2406.12045 [cs.AI]
- [48] Yifan Zhang, Xinkui Zhao, Jianwei Yin, Lufei Zhang, and Zuoning Chen. 2024. Operating System And Artificial Intelligence: A Systematic Review. *Preprint arXiv:2407.14567* (2024).
- [49] Jiyang Zheng, Jialiang Shen, Yu Yao, Min Wang, Yang Yang, Dadong Wang, and Tongliang Liu. 2025. Chain-of-Focus Prompting: Leveraging Sequential Visual Cues to Prompt Large Autoregressive Vision Models. In *13th International Conference on Learning Representations*.
- [50] Andy Zou, Long Phan, Sarah Chen, James Campbell, Phillip Guo, Richard Ren, Alexander Pan, Xuwang Yin, Mantas Mazeika, Ann-Kathrin Dombrowski, Shashwat Goel, Nathaniel Li, Michael J. Byun, Zifan Wang, Alex Mallen, Steven Basart, Sanmi Koyejo, Dawn Song, Matt Fredrikson, J. Zico Kolter, and Dan Hendrycks. 2023. Representation Engineering: A Top-Down Approach to AI Transparency. *Preprint arXiv:2310.01405* (2023). arXiv:2310.01405 [cs.LG]
- [51] Andy Zou, Long Phan, Justin Wang, Derek Duenas, Maxwell Lin, Maksym Andriushchenko, J Zico Kolter, Matt Fredrikson, and Dan Hendrycks. 2024. Improving Alignment and Robustness with Circuit Breakers. In *38th Annual Conference on Neural Information Processing*

## A Case Studies

We highlight two agentic applications where FMOS helps manage complexity and support system evolution.

### A.1 Acceleration of scientific discovery

Multi-agent systems increasingly support hypothesis generation, experiment planning and control, simulation, and analysis. These workloads stress three FMOS capabilities:

**Knowledge augmentation:** Scientific evidence is distributed across heterogeneous, multimodal sources (tables, figures, time series). FMOS’s Data Agent provides a policy-driven substrate that unifies visual and textual evidence via guided retrieval. For example, a query such as “Assess catalytic activity for hydrogen evolution of this MoS<sub>2</sub> microscopy tile” can trigger a knowledge trap that retrieves relevant image regions and supporting literature, retaining only task-critical context in the active window.

**Domain-aware model selection:** Scientific tasks often require both domain knowledge and strong visual reasoning. The FM Composition Optimizer routes the request to an appropriate physical FM (or a composite) based on prior task performance under compute and latency budgets.

**Verification under physical constraints:** Outputs must respect scientific priors and physical laws. FMOS escalates to a verification path when needed, invoking multi-step checks and constraint-aware reasoning. Over time, it can reuse validated associations (e.g., between defect signatures, free-energy diagrams, and polarization curves) to guide subsequent retrieval and reduce unnecessary re-computation.

### A.2 Technical Support: Adaptive Context Management

Enterprise technical support agents use complex decision trees to troubleshoot specific customer technical issues. The challenge is to enable agents to flexibly manage both a detailed, “zoomed-in” view of the current decision-tree node and a broader, “zoomed-out” context of the overall troubleshooting path. This ensures continuity across complex support flows. Without infrastructure to fluidly switch and validate these contexts, agents risk losing track of node states, leading to guidance errors. FMOS enables this through dynamic context handling. Rather than using a fixed window, it employs a hierarchical decision tree of past interactions that evolves over time. The Data Agent routes context based on the query’s position in this tree, maintaining state across sessions. This demand-paging-like approach draws from OS memory principles and avoids manual memory engineering.

## B Why FMOS Enables Co-evolution, Co-optimization, and Reuse: A Virtualization Lens

### B.1 A minimal FM-virtualizability condition (and what it buys us)

An application interacts with the foundation-model stack through an operation set  $\mathcal{A}$  (e.g., generate, retrieve, cite-check, tool-call, write/read memory, route to another FM). Following virtualization tradition, we partition operations into: (i) *innocuous* operations  $\mathcal{A}_{\text{ino}}$  that can run on the fast path, and (ii) *sensitive* operations  $\mathcal{A}_{\text{sen}}$  whose effects depend on (or can change) shared *resources* (budgets/quotas), *knowledge* (memory tiers/indexes), or *trust state* (policy gates, provenance).

FMOS designates a set of *privileged* operations  $\mathcal{A}_{\text{priv}} \subseteq \mathcal{A}$  that must “trap” to the FMOS control plane (dispatcher/allocator/interpreters). We assume a minimal virtualizability condition:

$$\boxed{\mathcal{A}_{\text{sen}} \subseteq \mathcal{A}_{\text{priv}}} \tag{1}$$

i.e., every operation that can impact shared budgets/knowledge/trust is mediated by FMOS.

**Conservative guarantees (analogous to VM goals).** Under (1), FMOS can *target* three properties (we phrase them conservatively to avoid overclaim):

- **Efficiency (fast path):** operations in  $\mathcal{A}_{\text{ino}}$  do not require orchestration and can execute with minimal FMOS involvement.
- **Resource control:** effects on shared resources/trust/knowledge occur only via trapped operations, enabling enforceable budgeting and policy checks *within the VFM interface*.
- **Interface-level equivalence:** applications program against a stable VFM interface (ABI); FMOS may change internal realizations while preserving agreed semantics (up to latency/noise/stochasticity).

Crucially, cross-cutting improvements (retrieval, verification, routing, memory updates) live in a small trapped surface, while most application logic remains unchanged on the fast path.

## B.2 Co-evolution of capabilities via FMOS (shared learning of privileged mechanisms)

FMOS centralizes a small set of *privileged* operations (e.g., retrieval, routing, verification, memory updates) whose execution must be mediated to enforce budgets, knowledge consistency, and trust. These operations are heterogeneous and may exhibit different dynamics: some components can stabilize quickly (e.g., routing heuristics in a fixed model pool), while others require continual refinement (e.g., verification under domain shift). We therefore model FMOS as a modular control program that selects a *structured* privileged action  $a_t = (a_t^{\text{ret}}, a_t^{\text{route}}, a_t^{\text{ver}}, a_t^{\text{mem}})$  given an observed request state  $s_t$  (domain, uncertainty, risk, budget, intent, etc.). A convenient factorized form is

$$\pi_\theta(a | s) = \pi_{\theta_{\text{ret}}}(a^{\text{ret}} | s) \pi_{\theta_{\text{route}}}(a^{\text{route}} | s) \pi_{\theta_{\text{ver}}}(a^{\text{ver}} | s) \pi_{\theta_{\text{mem}}}(a^{\text{mem}} | s), \quad (2)$$

where each module has its own parameters but shares the same state representation  $s$  and contributes to the end-to-end quality–trust–cost objective.

Consider  $K$  applications/tenants producing traces  $\tau \sim \mathcal{D}_k$  with loss  $\ell_k(\tau; \theta)$  capturing task error, trust penalties, and resource cost. FMOS learns shared parameters  $\theta = (\theta_{\text{ret}}, \theta_{\text{route}}, \theta_{\text{ver}}, \theta_{\text{mem}})$  by minimizing

$$\theta^* \in \arg \min_{\theta} J(\theta) \triangleq \sum_{k=1}^K w_k \mathbb{E}_{\tau \sim \mathcal{D}_k} [\ell_k(\tau; \theta)]. \quad (3)$$

The co-evolution effect arises because privileged decisions generate repeated, comparable feedback across applications, allowing FMOS to refine shared modules from pooled trapped events. Let  $N_m = \sum_k N_{m,k}$  denote the total number of events that provide learning signal for module  $m \in \{\text{ret}, \text{route}, \text{ver}, \text{mem}\}$  (e.g., invocations or audited outcomes). For bounded module complexity  $\text{Comp}(\Pi_m)$ , standard generalization arguments yield a module-wise estimation error of the form

$$\sup_{\theta_m \in \Theta_m} |J_m(\theta_m) - \widehat{J}_m(\theta_m)| \leq O\left(\sqrt{\frac{\text{Comp}(\Pi_m)}{N_m}}\right), \quad (4)$$

where  $J_m$  denotes the population objective attributable to module  $m$  and  $\widehat{J}_m$  its empirical estimate. Since  $N_m$  aggregates feedback across all dependent applications, learning signal for a given module can accumulate substantially faster than in isolated per-application tuning (where only  $N_{m,k}$  samples are available). This yields a precise and conservative statement: *for any privileged module that receives pooled feedback (not necessarily all modules), FMOS can improve its decisions more sample-efficiently and propagate those improvements to all applications that share the VFM interface.*

Heterogeneous convergence is naturally handled by module-specific update schedules, e.g.,

$$\theta_m^{t+1} = \theta_m^t - \eta_m(t) \nabla_{\theta_m} \ell_t(\theta), \quad (5)$$

allowing stable components to change slowly while continually adapting components update more aggressively. The same logic applies in the single-application setting: as long as trapped events recur over time (across sessions, users, or subtasks), the effective  $N_m$  grows and FMOS can refine the corresponding privileged modules accordingly.

## B.3 2) Co-optimization of resources (global allocator + trust-aware budgets)

**Coupled budgets are the point.** Let there be  $R$  shared resources: tokens, GPU time, tool-call quota, latency budget, memory writes, verifier invocations. At time  $t$ , application  $k$  chooses privileged action  $a_{k,t}$  with utility  $u_k(a_{k,t})$  and consumption  $c_r(a_{k,t})$ . FMOS solves a global constrained optimization:

$$\max_{\{a_{k,t}\}} \sum_{k,t} u_k(a_{k,t}) \quad \text{s.t.} \quad \sum_{k,t} c_r(a_{k,t}) \leq B_r, \quad \forall r \in \{1, \dots, R\}. \quad (6)$$

**Shadow prices yield coordinated decisions (under standard assumptions).** Introduce multipliers  $\lambda_r \geq 0$  (“shadow prices”) and consider the Lagrangian

$$\mathcal{L}(\{a_{k,t}\}, \lambda) = \sum_{k,t} \left( u_k(a_{k,t}) - \sum_{r=1}^R \lambda_r c_r(a_{k,t}) \right) + \sum_{r=1}^R \lambda_r B_r. \quad (7)$$

Given  $\lambda$ , each application selects actions locally:

$$a_{k,t}^*(\lambda) \in \arg \max_{a \in \mathcal{A}_{\text{priv}}} \left( u_k(a) - \sum_{r=1}^R \lambda_r c_r(a) \right). \quad (8)$$

FMOS updates  $\lambda$  to satisfy budgets (dual ascent), yielding a globally optimal allocation for (6) under convexity/regularity (and a principled heuristic otherwise). This is the formal meaning of *co-optimization*: a shared allocator sets system-wide prices/policies so that many local decisions collectively respect shared budgets and maximize total value.

**Trust/safety as a first-class coupled constraint (not an afterthought).** Let  $S(\{a_{k,t}\})$  denote an aggregate risk measure (e.g., expected policy violation / hallucination / unsafe tool side-effect rate). FMOS can enforce a risk budget:

$$\max_{\{a_{k,t}\}} \sum_{k,t} u_k(a_{k,t}) \quad \text{s.t.} \quad \sum_{k,t} c_r(a_{k,t}) \leq B_r (\forall r), \quad S(\{a_{k,t}\}) \leq \varepsilon. \quad (9)$$

A Lagrangian form adds a risk multiplier  $\mu \geq 0$ :

$$u_k(a) \mapsto u_k(a) - \mu s(a), \quad (10)$$

where  $s(a)$  is the per-action risk contribution (e.g., skipping verification, calling external tools, writing memory). This makes “trust” compatible with the same allocator logic: verification/trust checks become privileged actions whose use is optimized subject to explicit risk budgets.

**Single-application case.** With  $K = 1$ , co-optimization still applies because the decision is *intra-application*: FMOS allocates resources across the application’s own components (retrieve vs. verify vs. generate vs. tool-use), and across concurrent sessions/agents, under shared budgets and risk constraints.

### B.4 3) Reusability across enterprises (interface contract + approximate equivalence)

**VFM ABI: program to the interface, not the implementation.** A core virtualization promise is that applications target a stable interface while the substrate may change. For FMOS, applications program to a VFM “ABI” (API + semantics) independent of the underlying physical FMs, vector stores, tools, or verifiers.

Let  $S_P$  be the physical FMOS state (models, caches, indexes, policies, tool handles) and  $S_V$  the virtual state exposed to applications (virtual memory/context, virtual budgets, virtual trust guarantees). FMOS implements a mapping  $f : S_P \rightarrow S_V$  such that for any application-visible operation sequence  $e$  there exists an FMOS-internal realization  $e'$  satisfying an interface-commutation condition:

$$\boxed{f(e(S)) \approx_C e'(f(S))} \quad (11)$$

where  $\approx_C$  denotes *approximate equivalence under a contract C* (e.g., budgets respected, provenance attached, safety policy enforced, memory consistency semantics, and task-level acceptance metrics). We use  $\approx$  (not strict equality) to acknowledge stochastic generation and changing model backends.

**Why this yields enterprise reuse.** Eq. (11) formalizes that applications depend on VFM semantics, not physical realization. Therefore, *to the extent that FMOS maintains the contract C*: (i) agent code ports across organizations with different model stacks, (ii) domain capabilities can be packaged as FMOS “drivers” (retrievers, memory schemas, verifiers, policy modules), and (iii) upgrades to physical models/tools can occur with limited application changes—provided the VFM contract remains stable and sensitive operations continue to trap via (1).

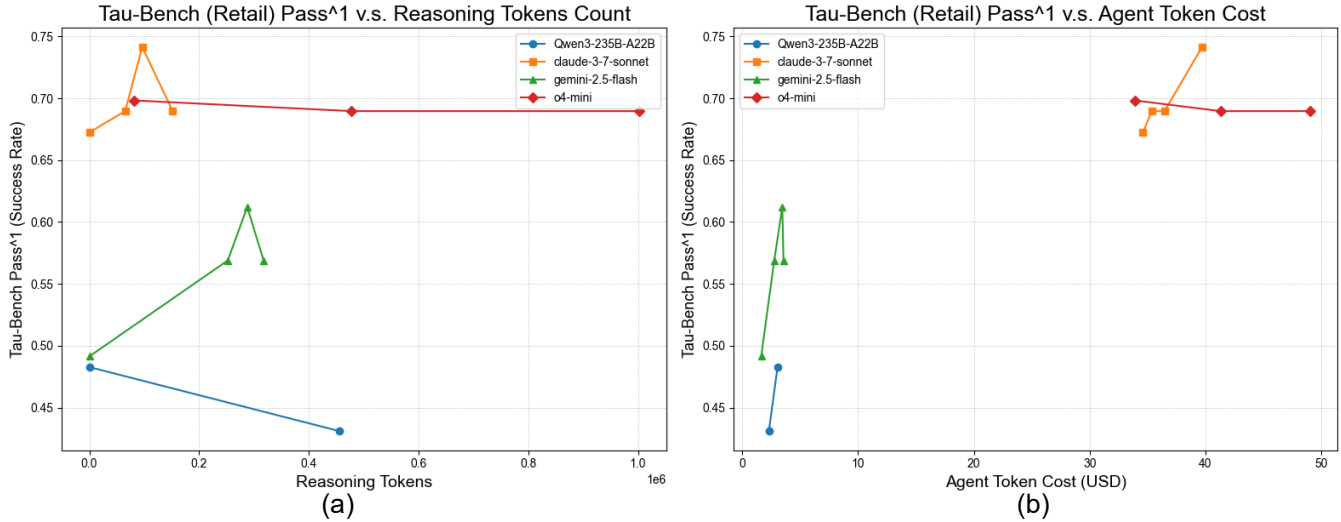
### B.5 Summary

- **Co-evolution:** FMOS centralizes privileged mechanisms and learns them from pooled traces; generalization improves with total trapped samples  $N$  (Eq. ??). This holds across many apps ( $K > 1$ ) and over time within one app ( $K = 1$ ).
- **Co-optimization:** FMOS acts as a global allocator for coupled budgets and risk constraints; shadow prices coordinate local choices into system-level policies under standard assumptions (Eqs. 6–9).
- **Reusability:** FMOS provides a stable VFM contract; approximate interface-level equivalence enables portability and upgradability without claiming identical outputs (Eq. 11).

## C A Control Reasoning Knob to Optimize Agent Execution

Figure 4 contrasts the *success rate* (Pass<sup>1</sup>) of  $\tau$ -bench [47] with (a) the *reasoning tokens* produced by an agent LLM and (b) its total *input/output token cost*. The study spans four recent reasoning LLMs—Qwen3-235B-A22B, gemini-2.5-flash-preview-05-20, claude-3.7-sonnet, o4-mini—each exercised under multiple “reasoning-budget” settings. Two clear trends emerge:

1. **Moderate deliberation improves reliability.** Most LLMs except Qwen3-235B-A22B exhibit a lift when moving from *no-thinking* to a modest level of reasoning tokens: e.g., gemini-2.5-flash rises from 0.49 (no reasoning budget) to 0.61 Pass<sup>1</sup> with a 4096-token budget.



**Figure 4.** Tau-bench (Retail) success rate (Pass<sup>1</sup>) vs.: (a) reasoning tokens in agent execution and (b) total cost of all tokens consumed by agent LLM

2. **Excessive reasoning degrades performance and cost efficiency.** Beyond a task-dependent “sweet spot” the accuracy curve turns downward while cost grows super-linearly. For example, increasing the gemini-2.5-flash reasoning budget to 16,384 tokens actually lowers Pass<sup>1</sup> to 0.57, even though token usage grows by more than four times. Similarly, o4-mini reaches its limit at the low effort setting; switching to high raises the cost per run by 50% without improving performance. Qwen3-235B-A22B does not offer reasoning-level control, and enabling reasoning mode causes the LLM to overthink, dropping Pass<sup>1</sup> from 48.3% to 43.1%.

The challenges of setting up these experiments highlighted the heterogeneity of control knobs across models: A practitioner cannot simply “dial in” the same budget across different models. For example, gemini-2.5-flash and claude-3.7-sonnet allow the setting of an explicit reasoning token budget, while o4-mini offers only three opaque effort levels (low|medium|high) without a direct token cap. Some open-source LLMs, such as Qwen3-235B-A22B, provide only a binary reasoning on/off switch. Fine-grained control might be achieved with techniques like Chain-of-Draft [45], latent-space reasoning [14], appending tokens such as “Wait” to extend reasoning or using end-of-thinking delimiters like “Final Answer:” to shorten it [4, 28]. However, integrating these techniques into a specific LLM serving framework can be non-trivial for application developers. This heterogeneity leaves today’s agent developers with brittle, model-specific heuristics that must be hand-tuned for every workflow and re-tuned as models evolve. Worse, mis-configuration can induce over-thinking [22], wasting compute while reducing correctness.

These findings motivate the *Learn-to-Reason* in FMOS’s Trust & Reasoning hierarchy (Section 3). By abstracting the notion of reasoning effort behind a `set_reasoning_level()` API, FMOS can: (a) **normalize heterogeneous control** that maps the user’s high-level budget request to the appropriate switch (on/off, effort level, or token cap) for each concrete LLM; (b) **self-evolve budgets on-line** by monitoring success signals, execution trace, and cost to iteratively converge to the task-specific sweet spot. In short, *reasoning is a double-edged sword*: indispensable for difficult tasks yet detrimental when uncontrolled. An OS-like abstraction layer that virtualizes reasoning budgets, just as classical OSs virtualizes memory, enables cost-aware and model-agnostic optimization of agent performance.

## D Virtual System Environment Services Enabling Agent Workflows

A virtual FM system must support a prototypical user workflow (see Figure 2) comprising input context preparation, model execution passes, and output processing in a manner that allows for system-level interception and control. Such interception enables a suite of environment services that abstract and manage these stages, for simplifying the development of FM-based agent applications and for enabling the system’s capacity for self-evolution.

## D.1 Agent Filesystems

Recent work on *agent filesystems* proposes such an abstraction by building an OS-like filesystem substrate tailored for AI agents [3, 12]. Instead of scattering agent state across ad hoc databases, logs, and local files, these systems encapsulate an agent’s runtime artifacts, key–value state, and tool-call audit trails behind a familiar filesystem interface, enabling post-hoc inspection, debugging, and reproducibility [12].

AgentFS, for example, implements an agent-oriented filesystem on top of a single SQLite file, making an agent session portable and snapshot-friendly while supporting queryable audit logs for observability and compliance [2, 12]. Isolation mechanisms such as copy-on-write overlays allow agents to safely use real command-line tools without mutating the underlying host project until changes are reviewed and applied [2]. This “single durable artifact” design also makes it practical to fork state for subagents and to time-travel (rollback) during development and evaluation [2, 12].

Nexus generalizes this direction into a programmable, backend-agnostic filesystem for AI agents that combines file storage, memory across sessions, fine-grained (relationship-based) permissions, and semantic search under a unified API [29, 30]. This consolidates several system concerns—persistent memory, access control, and multi-agent sharing—into one substrate that can be deployed locally or in multi-tenant settings [30].

From the perspective of context engineering, filesystems provide agents with an interface to store, retrieve, and update an effectively unbounded amount of context without bloating the prompt window [15]. Deep agents can offload large tool outputs (e.g., web-search dumps) to files and selectively pull back only the needed spans using filesystem search primitives such as `ls`, `glob`, and `grep` [15]. Files also serve as a natural mechanism for long-horizon plans, subagent handoffs, and skill/instruction libraries that can be loaded on demand rather than permanently occupying the system prompt [15]. Finally, because agents can write to their own filesystem, user feedback and operational lessons can be persisted as editable artifacts, providing a concrete substrate for longitudinal self-improvement via versioning and rollback [12, 15].

Within an FMOS architecture, agent filesystems complement the Data Agent by providing a durable memory tier and audit substrate shared across workflows and agents. They operationalize system-level policies (e.g., size-based offloading, skill paging, and trace capture) while supplying OS-like primitives—permissions, versioning, and event triggers—that are essential for safe, governable, self-evolving VFMs [2, 30].

## D.2 Context Management and Knowledge Augmentation

A FM combines internal (parametric) knowledge acquired during training with (non-parametric) knowledge that it receives as input context (prompts). System environment services control this context both to elicit (selectively focus on) what the model knows and to expand (augment) it with external knowledge.

**Context memory management** Context memory management [7, 24, 32] methods “virtualize” limited LLM context window space by retrieving or swapping appropriate content in/out from a persistent store, enabling the creation of stateful agents. Letta/MemGPT [32] achieves this using an elaborate system prompt that “teaches” an LLM to summarize, recall, and edit information in its context memory using a toolbox of functions. AIOS [24] splits conversation context into blocks and use a k-LRU policy. Self-RAG [7] fine-tunes the target LLM with retrieval and critic tokens to trigger retrieval and assess value of retrieved information. These techniques remain useful even for FMs supporting long contexts by filtering irrelevant data and saving costs. However a key challenge in realizing context memory management services lies in automatically learning to identify and focus on what is relevant, which may vary across different domains and use cases.

**Missing interface for mapping context heuristics into the context management policies** This missing interface matters because application developers often *know* which pieces of context are valuable (and when), but cannot express that knowledge to the serving layer that owns the prompt budget. Concretely, developers may want to specify rules such as:

1. **Web research agent (post-answer offload).** After a scraped page has been fully used to answer the local question, proactively offload the full page content to a file and keep only a pointer plus a brief summary in the active window. Today this is typically implemented manually at the application level; however, since the harness already performs compaction/offloading, the same rule could be enforced at the serving layer and reuse any underlying memory tier (files, vector stores, databases) uniformly.
2. **Enterprise infrastructure agent (size-based tool-output offload).** If an MCP server returns a JSON result larger than 20KB (developer-specified threshold), store the payload in a file and insert only a pointer + schema/summary into the prompt before the next action. Current frameworks do not offer a portable way to bind such application-specific thresholds to the underlying context manager.
3. **Adaptive Agent skills unloading.** After a skill has been used, unload the corresponding SKILL .md from active context (even if it might be needed later), retaining only a compact “capability header” and a retrieval handle. This becomes important

as skills evolve and become lengthy; yet serving frameworks have ad-hoc and non-programmable defaults for progressive skill loading.

4. **Deep research agents (retain thoughts, prune observations).** Empirically, aggressively pruning accumulated web-search/tool results while preserving the agent’s internal reasoning trace can improve final outputs; MiroThinker operationalizes a related principle via recency-based retention of tool responses while preserving the full thought/action trajectory [25]. Today, such policies are mostly hand-engineered in application code instead of being declaratively enforced at the serving layer where they could be reused across tasks.

**Knowledge compression and retrieval** Typically, indexing and retrieval of external knowledge, prior to context augmentation, is achieved using multiple components such as embedding models, retriever models and ranking models, sometimes jointly tuned along with the target FM. Yang et al. [46] fine tune a base LLM to compress and retrieve augmented knowledge in terms of a hierarchical state representation, for lifelong context management across sessions, while MemTree [36] maintains a dynamic tree structured memory representation. System environment services can optimize encoding and retrieval performance and resource efficiency, based on workload patterns and context (e.g. context aware pre-fetching, caching, or switching encoding and retrieval algorithms).

**Knowledge oriented abstractions for different modalities** External knowledge/data may available in a variety of structures and modalities. Further, many scientific explorations involve large scale data and a continual influx of new observations and data. System environment services can enable the efficient curation of such data into knowledge oriented abstractions suitable for augmenting FMs. This curation may in turn be aided by using an FM, and could require optimizations to support high throughput and scale. Different scientific domains and even different scenarios for the same scientific domain may require slightly different abstractions.

### D.3 Reasoning and Trust Augmentation

Reasoning capabilities are essential both in the discovery or evolution (and integration) of new knowledge and when leveraging existing knowledge, tools, simulators etc and for ensuring that FMs generate safe, responsible and trustworthy results. System environment services can control FM output selection and processing (e.g., through constrained decoding, sampling, invoking verification and planning tools, representation engineering [50, 51]), either at the final or intermediate layers.

**Expand and manage reasoning resources (system 1, system 2)** FM augmented reasoning and planning may be characterized into different modes, analogous to human cognition, e.g., a fast thinking *System 1* (e.g., a direct inference) and a slow deliberative (multi-step) thinking *System 2* [17]. These modes have different resource (and reliability) profiles: System 2 places a heavier load on inference time resources, while System 1 improvements require training and fine tuning resource. System services which control and activate suitable modes of reasoning as needed, would help enable effective tuning, management and sharing of reasoning resources across applications.

**Reasoning at multiple levels (tiers): Abstract and specialized reasoning** Answering complex questions and formulating hypotheses in science requires multi-step, hierarchical reasoning that draws from different domain specific concepts. Effective reasoning process templates can be stored and enhanced over time by the system in procedural memory. Analysis of token probability distribution (and associated properties such as entropy and variance) can be used to guide reasoning decisions, for instance by providing insight into how certain tokens dominate or diversify the reasoning space [10].

#### **Low-overhead verification, protection and steering mechanisms**

When FMs produce unreliable, biased or unsafe outputs this could have cascading implications in autonomous FM agent workflows. External verifiers, evaluators, critics and guardrail models/agents are often used to moderate and control FM outputs, along with alignment techniques built into FMs. Controls close to the source ensure broad protection but have limited context and higher resource costs. Workflow-specific checks reduce false positives and negatives but are harder to generalize and leave other workloads vulnerable. System level services can be designed to tackle these trade-offs, ensuring consistency, relevance, flexibility and efficiency in verification, protection and steering mechanisms, including measuring and tracking uncertainty of reasoning paths.

### D.4 Model Resource Sharing and Orchestration

FM inference and adaptation demand significant GPU resources, which escalate with reasoning stages, RLMs, and multi-agent setups, causing resource contention. Using smaller FMs, distilled models, and optimized orchestration can mitigate these issues. The underlying model inference/serving platforms typically perform several optimizations for all requests to a given FM, but system environment services can intercept them [1] and use its awareness of higher level intent enabling much deeper co-optimizations and management of tradeoffs involved in both model selection and orchestration.

**Scheduling and mapping:** Given a pool for underlying FMs and tools, and a set of application agents/workflows, there are two primary aspects of scheduling to be considered: mapping requests from application agents to one or more suitable

FMs from underlying pool of FMs (analogous to mapping process threads to CPU and accelerators) [31, 39] and allocation and scheduling of these FM resources across user agents/FM applications (analogous to scheduling / context switching between user processes and threads)[24]. System environment services that have more direct context can learn to guide both of these optimizations better and pass them as hints to underlying subsystems such as inference platforms, which in turn can leverage underlying base OS [35].

**Model composition and instantiation:** In domain specific research where certain models may possess deep domain capabilities that are valuable for a given application but lack important capabilities which are present in other models. System environment services can enable the creation, configuration and provisioning of suitable composite FMs or distilled FMs to meet both capability augmentation and resource/latency constraints (e.g. for agents that perform in the loop steering of experiments or high throughput IT system events monitoring).

**Profiling, measurement, and tracing:** In order to optimize, evolve and control FM systems effectively, new profiling and tracing services would be needed that provide visibility into key FM states and operational conditions.

## D.5 Broader Considerations

Beyond application and system considerations related to the three elements and their interactions at any given point, the design of an operating system for FM workloads in open ended domains also involves some broad long term considerations.

**Continual self-evolution:** The world around us is always changing. The ability to evolve and adapt with these changes is especially important for system environments that support FMs in scientific discovery and other open ended domains, where new observations, new scenarios to learn from keep emerging and new knowledge is being constantly being generated, verified and refined.

**Continual adoption of latest techniques:** The world of AI also continues to advance at a phenomenal pace (calling into question the shelf life of a position paper like this). Systems environments that support emerging FM workloads, therefore need to be designed with this reality in mind, and must be able to continually adopt better models, frameworks and methods at the same pace, so that every workflow automatically benefits from those advancements.

**External control:** System environments should have mechanisms for external controls to be applied automatically across all workflows by administrators to allow incorporation of global policies, particularly with respect to safety guidance, compliance and resource bounds, e.g. using techniques for incorporation of privileged instruction hierarchy in FMs [41]

## E Additional Alternative Views

*AI workloads are not that different; few if any changes are needed to conventional OS concepts and methods.*

The FMOS is built as a layer over a conventional OS [24]; hence it can use insights from observing these resources to provide hints to the underlying OS, as well as control the agent application environments using OS-inspired principles [24, 32]. This opens up fresh approaches to address a classical cross-layer dichotomy in OS design: how to provide workload intent to the OS, and system resource awareness to workloads, without breaking abstraction boundaries. Bridging this gap enables better coordination between workflows and the system in order to make the best decisions at both the application and system level.

*Model Context Protocol (MCP) and Agent to Agent communication Protocol advancements will address most of the challenges*

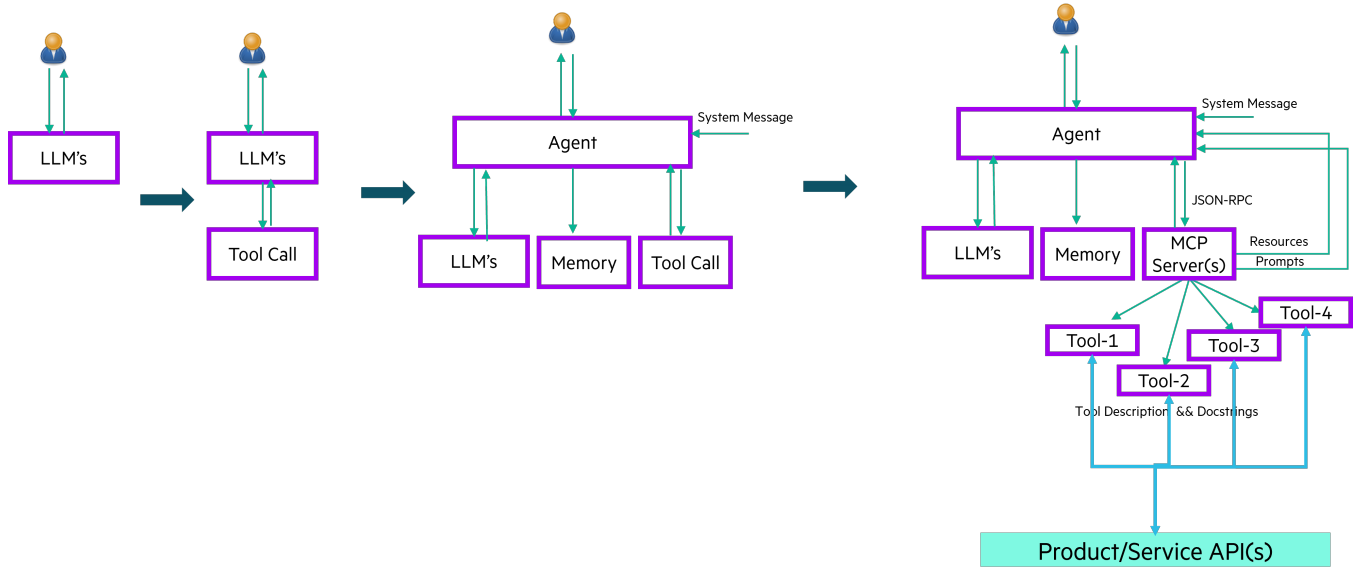
The community has progressed from directly interacting with Foundation Models and devising ways to establish context (such as RAG, GraphRAG etc.) to perform tool calling (introduced by OpenAI in 2023). This led to a cacophony of orchestration frameworks (Langchain, Langgraph, Langflow, n8n, etc.) each on a journey to create a viable ecosystem of libraries/components/modules to enable users and developers to select between different LLM providers (cloud or on-prem), libraries to interact with product or service/providers. This led to tools, resources and prompts being created and a flourishing ecosystem evolved in except that these implementations were captive to their orchestration framework. The introduction of Model Context Protocol as an unifying standard to connect various Agents to business product/service tools, exposing prompts and resources has been a recent game-changer.

An AI agent (typically acting as an MCP client) can now interrogate MCP Server(s) over JSON-RPC to list and execute tools (actionable functions), list and express prompts (interactive templates) and list and provide resources (data). This has enabled an explosive growth of MCP servers of various persuasions across the industry. A parallel contribution has been the introduction of the Agent-to-Agent (A2A) protocol which enables Agents to advertise their capabilities through a template endpoint (./well-known/agent.json). The combination of these two enables an agent to subscribe to one or more MCP servers; and one/more such agents being able to interoperate in a standardized fashion with each other.

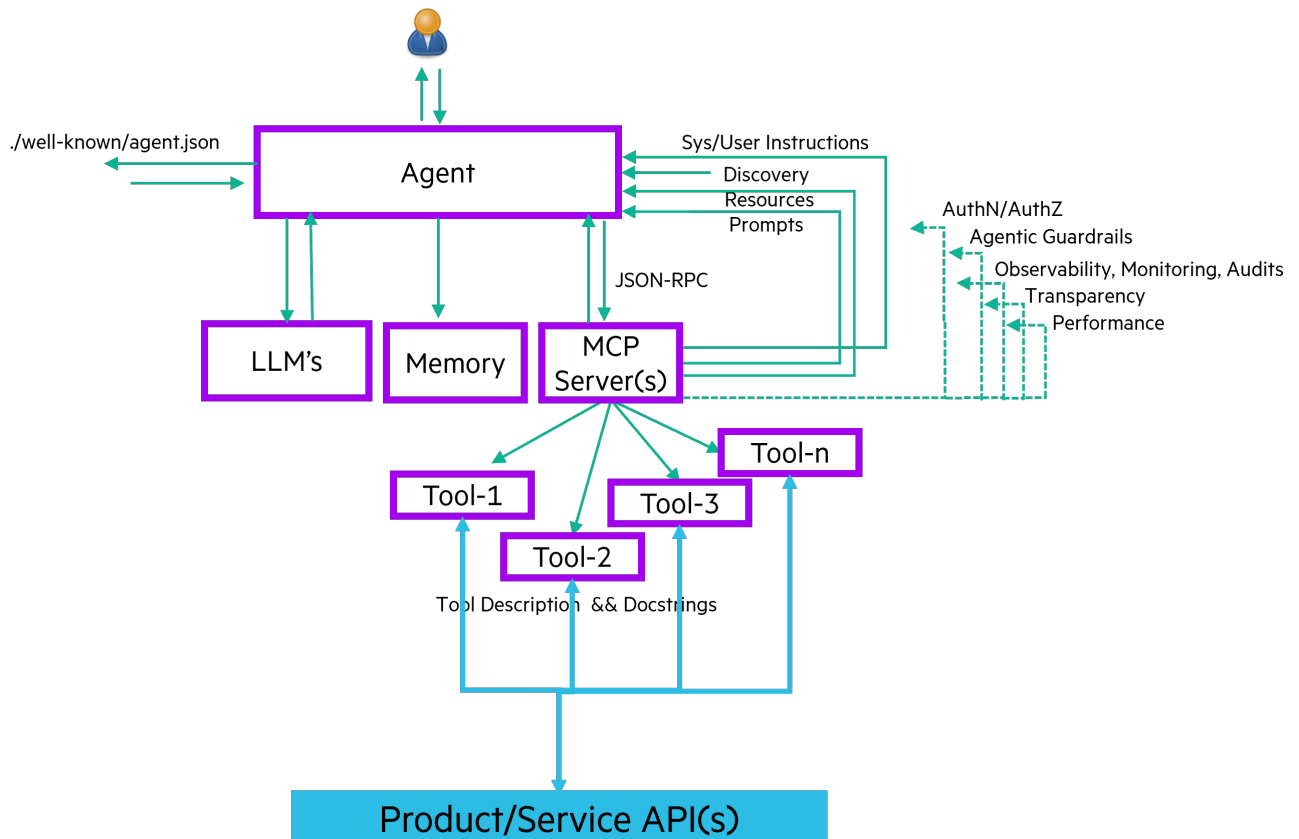
These developments have enabled agents and multi-agent systems in the personal assistant space to become wildly successful. However, these still does not address all requirements required for reliability and scalability in Enterprise computing

environments. We outline a few such requirements and therefore justify why the FMOS approach outlined in this paper might be a more suitable option.

1. **AuthZ/AuthN:** The first generation of agents have evolved from developer centric workflows. Several have automated tasks which otherwise needed multiple steps in an IDE and/or access to a Cloud resource. In those circumstances directly passing API Keys (with developer equivalent credentials) or OAuth2.1 credentials with callbacks (good enough for a human to access) was sufficient and passed on to an assistant. However this causes problems for more sophisticated enterprise products or services where tighter Authentication and Authorization separation is necessary for a downstream tool. Recently Authorization support with OAuth 2.0 has been adapted into the MCP spec for HTTP transport (SSE Server Side Events [27]).
2. **Agentic Guardrails:** Enterprise applications typically require Service Level Agreements (SLA's) and are tightly coupled iwth one another. Exposing an AI agent to autonomously make changes into these systems exposes the environment to excessive risk. There is a need to intercept the agent-tool, agent-data, and agent-LLM interactions and add appropriate guardrails in a similar fashion as with LLM's. This might differ based on the class, degree of sophistication and need of the downstream Enterprise application or service. This is necessary to develop since it is easy to cause irreversible changes to Enterprise platforms and databases with a malformed prompt, hallucinating LLM or a badly implemented tool function exposed by an MCP server. Finally, since there is no established marketplace for MCP servers, there are often tens of MCP servers for the same application or cloud-service in open-source (and progressively getting worse). This requires a rethink to FMOS structured approach as advocated earlier in this paper.
3. **Observability, Monitoring and Audits:** It is necessary to debug, trace and monitor, agents and agentic workflows during their entire lifecycle (design and operation). An entirely new set of observability tools such as LangSmith and Langfuse are now being created in the application layer (where modern agent development is taking place). While they are necessary, they will lead to "software bloat" at a fundamental level. Bringing agentic AI development into the FMOS realm as advocated in this paper would therefore be a better approach.
4. **Transparency, Reliability and Trust:** It is well known that the best performing Foundation Models can still hallucinate. Since FM's are essentially the brains of a modern AI agent, it is therefore inherently untrustworthy. Special attention has to be made at various layers of the hierarchy (such as tool description/Docstrings, system/user prompts, resources exposed via an MCP server) to bring an element of dependability and reliability in agents. While observability, monitoring and auditability are necessary features, it is unlikely to be added into the MCP or A2A standards since they will likely make the spec bloated and unsustainable. This also calls for a rethink in the design of future agentic AI systems along the FMOS approach and interface advocated in this paper. In this way a dedicated Trust agent can be built and augmented independent of the MCP server or Agent and the appropriate fine-tuning or policy or guardrails implemented.
5. **Performance and Multi-tenancy:** First generation agentic AI workflows have focused on functionality over performance. Much of the performance optimizations have been limited to the AI inference layer. An end-to-end evaluation needs to be performed from the workflow being triggered to the chat prompt being generated. Modern containerized stacks enable inference, MCP servers, underlying service and the agent (MCP client) to be physically and logically separated across network segments. This will be unsustainable to optimize without a standardized FMOS construct. Finally, most first generation agents are built to run single instances per client. The next generation of agentic workflows will likely need to support multiple tenants in a single instance for which no easy enhancements can be done to prevailing MCP and A2A specifications without bloating them unsustainably.
6. **Client experience:** The overall experience with an MCP enabled agentic workflow unfortunately depends on which parts of the MCP spec are implemented by it. The MCP Client Feature support matrix [26] shows a wide variety of clients with varying levels of integration with MCP servers. This is unsustainable for the ecosystem in the long-term.



**Figure 5. Evolution of Tool calling with LLMs:** Architectural progression from LLM’s performing tool-calling to orchestration platforms to the construct exposed by MCP servers



**Figure 6. Looking into the future with MCP:** Constructs exposed by MCP protocol today, limitations and future possibilities