



Kelvin Open Science Publishers  
Connect with Research Community

Research Article

Volume 1 / Issue 1

KOS Journal of AIML, Data Science, and Robotics

<https://kelvinpublishers.com/journals/aiml-data-science-robotics.php>

# A Governance Oriented Study of Fine-Tuning Domain Specific Large Language Models with Transaction and Operations Data

Ramani Teegala

Lead Engineer-Java Full Stack, USA

\*Corresponding author: Ramani Teegala, Lead Engineer-Java Full Stack, USA

Received: March 06, 2024; Accepted: March 17, 2024; Published: March 19, 2024

**Citation:** Ramani T. (2024) A Governance Oriented Study of Fine-Tuning Domain Specific Large Language Models with Transaction and Operations Data. *KOS J AIML, Data Sci, Robot.* 1(1): 1-10.

**Copyright:** © 2024 Ramani T., This is an open-access article published in *KOS J AIML, Data Sci, Robot* and distributed under the terms of the Creative Commons Attribution License, which permits unrestricted use, distribution, and reproduction in any medium, provided the original author and source are credited.

## 1. Abstract

By early 2024, large language models were increasingly explored beyond general purpose language tasks and into enterprise environments that manage transaction and operations data. Organizations across finance, retail, logistics, healthcare administration, and platform operations evaluated fine tuning as a mechanism to adapt generic models to domain specific terminology, workflows, and decision contexts. Unlike prompt-based adaptation, fine tuning directly modifies model behavior through exposure to curated datasets, making it a materially different control surface with distinct architectural, operational, and governance implications. When transaction records and operational telemetry are introduced into the training process, concerns related to data sensitivity, correctness, traceability, and long-term model stability become central rather than peripheral. This paper examines fine tuning of domain specific large language models using transaction and operations data as practiced and understood up to February 2024. It situates fine tuning within the historical evolution of enterprise AI systems, contrasting it with earlier approaches such as rule engines, feature based machine learning, and retrieval augmented inference. While fine tuning offers potential benefits in terms of domain fluency and reduced prompt complexity, it also introduces risks that are less visible at inference time, including irreversible data leakage into model weights, reduced behavioral transparency, and challenges in validating model outputs against deterministic business rules. These risks are amplified in environments where transactional correctness, auditability, and regulatory compliance are non negotiable. The paper adopts a governance-oriented perspective that treats fine tuning not as a model optimization exercise but as a system level decision that reshapes accountability boundaries. It analyzes how transaction and operations data differ from general text corpora in terms of structure, sensitivity, and lifecycle, and how these differences affect data selection, training validation, and post deployment monitoring. Particular attention is given to failure modes that emerge when probabilistic models are trained on historically contingent operational data, including the propagation of outdated policies, reinforcement of process anomalies, and erosion of explainability. Rather than advocating for or against fine tuning, this study provides an architectural framework for evaluating when fine tuning is appropriate, how it should be bounded, and what governance mechanisms are required to make it defensible in enterprise contexts. The analysis emphasizes augmentation over automation and positions fine-tuned models as assistive components embedded within controlled system architectures. The goal is to provide practitioners and architects with a structured basis for decision making that aligns technical capability with operational realism and institutional responsibility.

**2. Keywords:** Domain specific large language models, Fine tuning governance, Enterprise AI systems, Transaction

data modeling, Operations data analytics, Probabilistic AI behavior, AI auditability, Model lifecycle management,

regulated enterprise systems, AI risk management, Human in the loop AI, Enterprise data governance

### 3. Introduction: Fine Tuning Large Language Models Within Enterprise Transaction and Operations Contexts

By early 2024, enterprise interest in large language models had shifted from exploratory experimentation toward deeper forms of adaptation intended to align model behavior with domain specific realities. While early deployments relied heavily on prompt engineering and retrieval based techniques to contextualize generic models, many organizations began to evaluate fine tuning as a more structural approach to embedding domain knowledge. This shift was driven by practical limitations observed in prompt based systems, particularly in environments dominated by transactional and operational data. In such settings, prompts alone often proved insufficient to capture implicit business rules, procedural nuance, and domain vocabulary that had evolved over years of system operation. Fine tuning appeared to offer a path toward models that internalized these patterns rather than repeatedly reconstructing them at inference time. Transaction and operations data occupy a distinct position within enterprise information landscapes. Unlike unstructured text such as documentation or knowledge articles, transactional records encode state changes, obligations, and outcomes that are often legally and operationally binding. Operations data similarly reflects system behavior, incident histories, performance characteristics, and human intervention patterns under real world constraints. When these data types are introduced into model training processes, they do more than improve linguistic familiarity. They influence how models generalize behavior, infer causality, and frame responses in contexts that may directly affect business decisions. This makes fine tuning on such data materially different from adapting models using publicly available corpora or curated domain literature. The introduction of historically contingent and context sensitive data into model weights raises questions about correctness, accountability, and long term stability that are central to enterprise system design.

Historically, enterprise software systems evolved under assumptions of determinism, traceability, and explicit control. Rule engines, workflow systems, and later machine learning models were designed so that decision logic could be inspected, tested, and constrained within well defined boundaries. Even supervised learning systems trained on historical data were typically scoped to narrow tasks with measurable outputs and retraining strategies that aligned with governance processes. Large language models, by contrast, operate as probabilistic systems whose internal reasoning is not directly inspectable. Fine tuning deepens this opacity by embedding domain specific patterns directly into model parameters, making it difficult to distinguish between learned behavior derived from transaction data and generalized reasoning inherited from pretraining. This challenges established enterprise practices for validation, audit, and change control. Another defining characteristic of transaction and operations data is their lifecycle volatility. Business rules change, regulatory interpretations evolve, and operational processes adapt in response to incidents and market conditions. Data captured at one point in time often reflects assumptions that later become invalid. When such data is used for fine tuning, models may internalize outdated logic or exceptional behaviors that were contextually correct but no longer appropriate. Unlike prompts or retrieval sources,

which can be updated or withdrawn with immediate effect, fine-tuned behavior persists until retraining occurs. This persistence introduces a form of architectural inertia that must be carefully managed. Enterprises therefore face a trade off between reducing inference time complexity and increasing the cost and risk associated with model updates.

This paper positions fine tuning of domain specific large language models on transaction and operations data as a system level architectural decision rather than a purely data science optimization. It argues that the appropriateness of fine tuning cannot be evaluated solely on output quality or task performance. Instead, it must be assessed in terms of governance impact, failure modes, and alignment with enterprise accountability structures. By grounding the analysis in practices and constraints observed prior to March 2024, the paper aims to provide a realistic and defensible framework for architects and engineers evaluating fine tuning within complex operational environments. The sections that follow trace the evolution of enterprise AI adaptation approaches, examine relevant literature, and develop conceptual and architectural models that clarify when and how fine tuning can be responsibly applied.

### 4. Evolution of Enterprise AI Adaptation Approaches for Transaction and Operations Data

Enterprise approaches to adapting intelligent systems to domain specific needs evolved long before the emergence of large language models. Early enterprise systems relied almost exclusively on deterministic logic encoded through rules engines, workflow orchestration platforms, and hard coded business processes. Transaction systems such as order management, billing, payments, and inventory control were designed around strict schemas and state transitions, ensuring that every operation could be validated, replayed, and audited. Operations data generated by these systems was primarily used for monitoring, reporting, and incident analysis rather than for driving autonomous decision making. In this period, adaptability was achieved through configuration and rule updates, not through learning from data, reflecting a strong preference for predictability and control. As enterprises began adopting machine learning techniques in the late 2000s and early 2010s, adaptation shifted toward supervised models trained on historical datasets. These models were typically applied to narrow tasks such as fraud detection, demand forecasting, anomaly detection, or classification of operational events. Transaction and operations data served as structured feature inputs rather than narrative context. Importantly, these models were embedded within architectures that preserved clear boundaries between prediction and action. Model outputs informed human decisions or downstream systems, but final authority remained explicit and traceable. Retraining cycles were formalized, datasets were curated with known labels, and performance degradation could be measured quantitatively. This approach aligned well with enterprise governance expectations, even as it introduced statistical uncertainty into decision support.

The introduction of large language models marked a departure from task specific learning toward general purpose reasoning over heterogeneous data. Initially, enterprises experimented with these models through prompt based adaptation, using carefully constructed instructions and contextual augmentation to guide behavior without

modifying model parameters. For transaction and operations domains, this often involved summarizing logs, explaining transaction flows, or answering questions over documentation and historical records. Prompt engineering and retrieval based techniques offered flexibility and reversibility, allowing organizations to adjust behavior without retraining models. However, as usage expanded, limitations became apparent. Prompts grew increasingly complex, context windows became saturated, and repeated exposure to the same domain patterns introduced inefficiencies and inconsistencies in model responses. By 2022 and into 2023, fine tuning began to attract attention as a means of embedding domain familiarity directly into model behavior. Enterprises observed that models fine tuned on domain specific text could reduce prompt complexity, improve terminology alignment, and exhibit more consistent reasoning patterns. In transactional and operational contexts, this raised the possibility of models that better understood process flows, exception handling, and system specific vocabulary. At the same time, fine tuning blurred established boundaries between data, logic, and behavior. Unlike supervised models with explicit targets, fine tuned language models absorbed patterns implicitly, making it difficult to isolate which aspects of transaction or operations data influenced specific outputs. This represented a significant shift in how enterprises thought about adaptation and control.

The evolution toward fine tuning also reflected broader operational pressures. Enterprises sought to reduce latency, simplify integration architectures, and lower the cognitive burden on users interacting with AI systems. Embedding domain knowledge into the model appeared to promise more natural interactions and reduced reliance on brittle prompt templates. However, this evolution introduced new risks that were not fully addressed by existing governance frameworks. Transaction and operations data often encode historical contingencies, workarounds, and exceptions that were contextually valid but not normative. When such data informs fine tuning, models may internalize behaviors that conflict with current policy or best practice. This evolution therefore created a tension between adaptability and institutional memory, forcing enterprises to reconsider how learning systems should evolve alongside changing operational realities. This section establishes that fine tuning did not emerge in isolation but as part of a continuum of enterprise adaptation strategies. Each stage traded off flexibility, control, and transparency in different ways. Understanding this progression is essential for evaluating fine tuning not as an inevitable next step, but as one option among several, each with distinct architectural and governance implications. The next section examines how academic and industry literature addressed these shifts, with particular attention to the risks and constraints associated with learning from transaction and operations data.

## 5. Literature Review on Fine Tuning Large Language Models with Enterprise Transaction and Operations Data

The body of literature available prior to March 2024 on fine tuning large language models reflects a gradual shift from model centric optimization toward system aware and governance oriented analysis. Early academic work on language model fine tuning primarily focused on improving task performance in areas such as text classification, summarization, and question answering. These studies typically relied on publicly available datasets or curated

domain corpora such as biomedical abstracts or legal documents. Transaction and operations data were largely absent from this early discourse, as such data posed challenges related to privacy, accessibility, and labeling. As a result, much of the foundational literature framed fine tuning as a technical mechanism for domain adaptation rather than as an architectural decision with enterprise wide implications. As large language models began to be evaluated in enterprise contexts, industry research and practitioner reports expanded the scope of analysis to include operational considerations. Studies published between 2021 and 2023 increasingly acknowledged that enterprise data differed fundamentally from general text corpora. Transaction data was characterized by high structural regularity, implicit semantics, and sensitivity, while operations data reflected system behavior under real world constraints rather than idealized processes. Literature in this period highlighted that fine tuning on such data risked encoding historical biases, procedural shortcuts, and exception handling patterns that were not always aligned with formal policy. However, many of these discussions remained descriptive, documenting observed behavior rather than proposing systematic governance frameworks.

A recurring theme in both academic and industry literature involved the trade off between fine tuning and alternative adaptation strategies. Retrieval augmented generation emerged as a prominent comparator, with researchers noting its advantages in terms of reversibility, transparency, and data isolation. Several studies argued that retrieval based approaches were better suited to environments with frequently changing rules, as updates to source documents immediately influenced model behavior without retraining. In contrast, fine tuning was associated with higher upfront effort and greater difficulty in isolating the influence of specific data sources on outputs. Despite this, literature also acknowledged scenarios where fine tuning produced more stable and coherent responses, particularly when domain terminology and process descriptions were repeatedly encountered. This tension underscored the need for context sensitive decision making rather than blanket recommendations. Another strand of literature examined risks associated with fine tuning from a safety and compliance perspective. Researchers noted that once transaction or operations data influenced model weights, traditional data access controls no longer applied in the same way. Even when original datasets were protected, traces of sensitive information could persist implicitly in model behavior. This raised concerns about inadvertent memorization, inferential leakage, and difficulty proving data deletion or correction. Studies on model inversion and membership inference attacks, while not always specific to enterprise data, were frequently cited as cautionary signals. In regulated industries, these risks were framed not merely as technical vulnerabilities but as governance challenges that affected auditability and regulatory defensibility.

The literature also reflected growing attention to lifecycle management and validation. Unlike supervised learning models, where performance metrics and ground truth labels guided evaluation, fine tuned language models lacked clear validation baselines for many enterprise tasks. Researchers emphasized the difficulty of establishing acceptance criteria when outputs were probabilistic and context dependent. Some proposed hybrid evaluation strategies combining qualitative review, rule based checks, and scenario testing, though consensus on best practices remained limited.

Importantly, several authors argued that evaluation must extend beyond output correctness to include stability over time, sensitivity to input variation, and alignment with evolving policies. This perspective resonated strongly with enterprise architects concerned about long term operational risk. Overall, the literature prior to March 2024 converged on the view that fine tuning large language models with enterprise data could not be evaluated solely through the lens of machine learning performance. While empirical evidence demonstrated potential gains in domain fluency and interaction efficiency, the associated risks to transparency, control, and governance were equally prominent. The absence of standardized frameworks for managing these trade offs left practitioners reliant on architectural judgment and institutional risk tolerance. This gap in the literature motivates the conceptual and architectural models developed in the following sections, which aim to translate fragmented research insights into a coherent system level perspective.

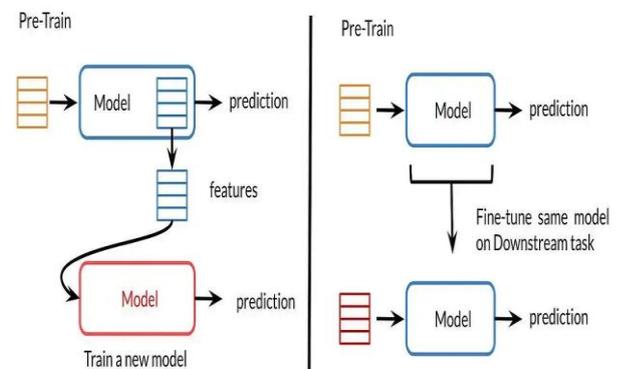
## 6. Conceptual Model for Fine Tuning Domain Specific Large Language Models Using Transaction and Operations Data

Fine tuning domain specific large language models using transaction and operations data can be understood as a control loop that transforms governed enterprise records into probabilistic behavioral change. This conceptualization differs from conventional views that treat fine tuning primarily as a data science activity conducted within isolated experimentation environments. In enterprise settings, fine tuning changes the behavior of an inference component that will be embedded inside workflows, user interfaces, and decision support systems that already operate under strict accountability constraints. A useful conceptual model therefore begins by treating fine tuning as a multi stage system process with explicit trust boundaries, lifecycle checkpoints, and governance obligations that persist beyond model training. The first stage in the model is domain intent definition, which establishes why fine tuning is being pursued and what domain capability is being targeted. For transaction and operations environments, intent is often framed as improving model fluency over transaction schemas, operational incident narratives, reconciliation procedures, or runbook guided responses. The model assumes that intent definition is a governance activity rather than a technical preference, since it determines which data classes may be used and what kinds of outputs will be considered acceptable. Without explicit intent, fine tuning risks becoming an uncontrolled absorption of domain history that produces inconsistent or unreviewable behavioral changes. Intent definition also clarifies whether fine tuning is intended to reduce prompt complexity, improve retrieval grounding, increase consistency of structured outputs, or encode procedural reasoning patterns that are difficult to restate reliably at inference time.

The second stage is data eligibility and boundary enforcement. Transaction and operations data are not interchangeable with general domain text because they often embed identifiers, behavioral fingerprints, and sensitive correlations even when obvious fields are removed. The conceptual model treats data selection as a classification exercise that determines which records can be used, under what transformations, and with what retention guarantees. This stage includes decisions about anonymization, aggregation, tokenization of structured fields, and exclusion of high risk attributes. Importantly, the model assumes that

data transformation is not a one time preprocessing step but a repeatable and auditable pipeline whose outputs must be versioned and reconstructible. In regulated enterprises, the ability to explain what data entered a training run and why it was permitted is as important as the final model performance. The third stage is adaptation execution, which encompasses the training procedure, the fine tuning objective, and the controls applied to ensure bounded behavioral change. In this conceptual model, adaptation execution is framed as controlled modification of an enterprise component rather than as open ended model improvement. This requires explicit parameter choices that reflect risk tolerance, including training duration, learning rates, and selection of fine tuning methods that constrain the magnitude of change. The model also assumes that training must be instrumented for traceability, capturing the provenance of datasets, the configuration of the run, and the resulting model artifact identifiers. Since language model behavior can shift in subtle ways, this stage must treat training as a change event that demands the same rigor as a major software release.

### Feature-based vs. Fine-Tuning

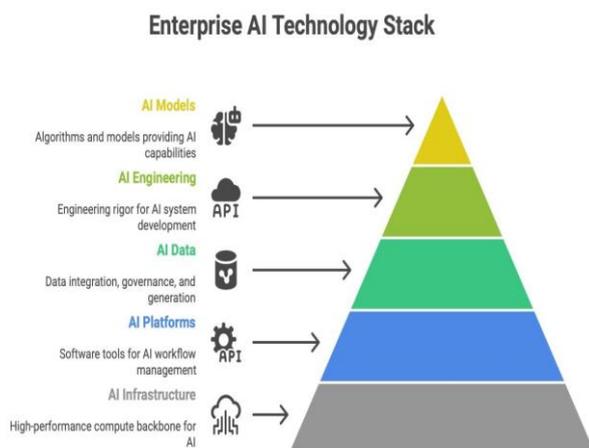


The fourth stage is behavioral validation and risk testing. Traditional model evaluation approaches such as accuracy or loss reduction are insufficient for most transaction and operations use cases, where the most significant risks are semantic errors, policy drift, and unsafe generalization. The conceptual model therefore requires validation to be scenario based and governance aligned. Validation includes test suites derived from transaction workflows, operational incidents, and boundary cases that reflect how the model will be used. This stage also includes stability testing, where small input variations are applied to measure output volatility, and policy adherence testing, where the model is challenged with prompts that should trigger refusals, redactions, or constrained answers. The goal is not to prove correctness exhaustively, which is unrealistic for probabilistic systems, but to establish defensible evidence that behavior is bounded within acceptable enterprise tolerances. The final stage is deployment governance and lifecycle control. Once a fine tuned model is deployed, it becomes part of a living operational system. The conceptual model requires mechanisms for version control, rollback, audit logging of inference interactions, and continuous monitoring for drift and misuse. It also requires a process for retiring models when policies change or when the underlying transaction and operations environment evolves materially. Unlike retrieval sources that can be updated instantly, fine tuned behavior persists until retraining, making lifecycle discipline essential. This stage closes the loop by feeding operational observations back into intent definition and data governance

decisions, ensuring that fine tuning remains a controlled and accountable practice rather than an accumulation of irreversible behavioral changes.

## 7. Layered Architecture for Fine Tuning Domain Specific Large Language Models in Transaction and Operations Environments

Translating the conceptual model of fine tuning into enterprise practice requires a layered architecture that explicitly separates concerns related to data governance, model adaptation, inference control, and operational oversight. In transaction and operations environments, this separation is not optional. These systems already operate under layered architectures designed to isolate data access, business logic, and execution authority. Fine tuned large language models must therefore be integrated in a way that preserves existing control boundaries rather than collapsing them into opaque model behavior. A layered architectural approach ensures that risks introduced during fine tuning do not propagate unchecked into downstream operational workflows. The outermost layer in this architecture is the domain data governance and preparation layer. This layer is responsible for curating, classifying, and transforming transaction and operations data before it is ever considered for model adaptation. It enforces policies related to data eligibility, anonymization, aggregation, and retention. In mature enterprise environments, this layer aligns with existing data governance programs, including data catalogs, classification frameworks, and access control mechanisms. Crucially, this layer produces versioned training datasets with documented provenance, ensuring that every fine tuning activity can be traced back to specific data sources and transformation rules. By isolating data governance from model training, the architecture prevents ad hoc or experimental data usage from silently influencing production model behavior.



The second layer is the model adaptation and fine tuning control layer. This layer encapsulates the training infrastructure, fine tuning configuration, and change management processes associated with modifying model behavior. It treats fine tuning as a controlled system change rather than as an iterative experiment. Configuration parameters, training objectives, and stopping criteria are explicitly defined and reviewed prior to execution. This layer also captures metadata about each fine tuning run, including dataset versions, base model identifiers, and resulting artifacts. In transaction and operations contexts, this documentation is essential for demonstrating that model changes were intentional, reviewed, and aligned with

approved use cases. Without this layer, fine tuning risks becoming an ungoverned accumulation of behavioral changes that cannot be explained or reversed. The third layer is the inference integration and behavior bounding layer. Once a model has been fine tuned, it must be embedded into enterprise workflows that often include user interfaces, automation pipelines, and decision support tools. This layer controls how the model is invoked, what inputs it receives, and how its outputs are consumed. Even with fine tuning, prompts, context assembly, and inference parameters remain critical control points. The architecture assumes that fine tuning does not eliminate the need for prompt discipline, but rather shifts some domain understanding into the model while preserving external constraints. This layer may enforce role based access, task specific invocation patterns, and output formatting rules that ensure the model operates within its intended scope. In high risk transaction workflows, this layer often routes model outputs through human review or downstream validation systems rather than allowing direct action.

The fourth layer is the validation, monitoring, and audit layer. This layer continuously evaluates model behavior in production against expected patterns, policy constraints, and operational tolerances. It logs inference interactions, captures deviations, and supports forensic analysis when anomalies occur. In transaction and operations environments, this layer is particularly important because errors may not manifest immediately. A model may produce plausible but subtly incorrect guidance that accumulates risk over time. Continuous monitoring enables organizations to detect drift, misuse, or unintended generalization early, before material impact occurs. Audit capabilities within this layer support regulatory inquiries and internal reviews by providing evidence of how and why model behavior evolved following fine tuning. The final layer is the lifecycle and change governance layer, which oversees the long term management of fine tuned models. This layer defines processes for periodic review, retraining, rollback, and retirement. It ensures that fine tuned models remain aligned with current transaction logic, operational procedures, and regulatory interpretations. Importantly, this layer recognizes that transaction and operations environments are not static. Business rules change, systems are modernized, and exceptional conditions that once dominated historical data may become irrelevant. Lifecycle governance prevents models from becoming repositories of outdated operational memory. By embedding fine tuning within a layered architecture that mirrors established enterprise system design principles, organizations can integrate domain specific large language models without sacrificing control, accountability, or operational resilience.

## 8. Comparative Analysis of Adaptation Strategies for Transaction and Operations Domains

Enterprises evaluating large language models for transaction and operations use cases prior to March 2024 were rarely deciding whether to adapt models, but rather how to do so within acceptable risk and governance boundaries. Fine tuning represented only one option within a broader spectrum of adaptation strategies that included prompt engineering, retrieval augmented inference, and hybrid approaches combining multiple techniques. Each strategy embodied a different balance between behavioral flexibility, control, transparency, and operational cost. A comparative analysis is

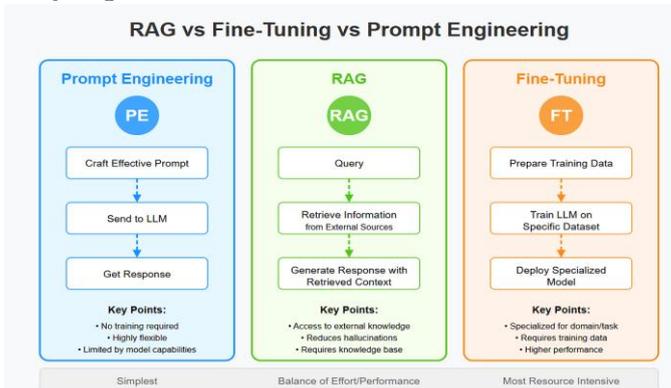
therefore essential to clarify when fine tuning is architecturally justified and when alternative approaches provide superior risk alignment. Prompt based adaptation was the earliest and most widely adopted strategy in enterprise environments. It allowed organizations to experiment with large language models without modifying model parameters or exposing sensitive data during training. For transaction and operations domains, prompts were often used to instruct models on how to interpret transaction records, summarize incidents, or explain process flows. This approach preserved strong reversibility and made governance comparatively straightforward, since prompts could be versioned, reviewed, and updated with immediate effect. However, as domains grew more complex, prompt based systems suffered from brittleness and inconsistency. Large prompts became difficult to maintain, context windows were strained by repetitive domain information, and subtle variations in wording could lead to materially different outputs. These limitations motivated interest in more durable forms of adaptation.

Retrieval augmented generation emerged as a response to some of these constraints by externalizing domain knowledge into controlled repositories. Instead of embedding transaction or operations data into model weights, relevant information was retrieved dynamically at inference time and injected as context. This strategy offered strong data isolation, since source records remained outside the model and could be updated or removed independently. For operational environments with frequently changing rules or documentation, retrieval based approaches provided a compelling balance between adaptability and control. Nevertheless, retrieval introduced its own architectural complexity. Systems had to manage relevance ranking, context filtering, and provenance tracking, and output quality depended heavily on retrieval accuracy. In transaction heavy environments, where meaning is often implicit rather than textual, retrieval alone sometimes failed to convey the procedural nuance needed for consistent reasoning. Fine tuning addressed some of these limitations by internalizing domain patterns directly into model behavior. When applied carefully, it reduced prompt complexity, improved consistency of terminology usage, and enabled models to reason more fluently about domain specific workflows. However, this came at the cost of reduced transparency and reversibility. Once transaction or operations data influenced model parameters, isolating or correcting specific learned behaviors became difficult. Fine tuning also complicated compliance narratives, as organizations had to demonstrate that sensitive data was handled appropriately during training and that its influence could be bounded. These trade offs made fine tuning most suitable for domains with relatively stable processes, well understood data semantics, and strong lifecycle governance.

Hybrid strategies attempted to combine the strengths of these approaches by fine tuning models on abstracted or synthetic domain representations while relying on retrieval for current or sensitive details. In such architectures, fine tuning encoded general domain fluency and reasoning patterns, while retrieval provided up to date transaction context and policy constraints. This reduced the amount of sensitive data embedded in model weights and preserved some degree of reversibility. However, hybrid approaches increased system complexity and required careful coordination between training, retrieval, and inference layers. Without disciplined governance, hybrids risked inheriting the weaknesses of both strategies rather than mitigating them.

Adaptation Strategy	Primary Control Mechanism	Governance Strength	Reversibility	Suitability for Transaction and Operations Data
Prompt Based Adaptation	Instruction design at inference	High	Immediate	Limited by complexity and consistency
Retrieval Augmented Generation	External knowledge injection	Very High	Immediate	Strong for evolving rules and reference data
Fine Tuning	Model parameter modification	Moderate	Low	Suitable for stable domain semantics
Hybrid Fine Tuning and Retrieval	Combined internal and external control	High but complex	Partial	Balanced but operationally demanding

This comparison illustrates that fine tuning is neither inherently superior nor inherently unsafe. Its value depends on the stability of transaction logic, the sensitivity of operational data, and the organization’s ability to govern model lifecycle changes. Enterprises that approached fine tuning as a substitute for governance encountered increased risk, while those that treated it as one component within a layered control strategy were better positioned to extract value without undermining accountability. The next section examines how these architectural choices translate into concrete methodological practices for implementing fine tuning in enterprise environments.



## 9. Methodology for Evaluating and Implementing Fine Tuning in Transaction and Operations Systems

The methodology for evaluating and implementing fine tuning in transaction and operations systems must be grounded in architectural realism rather than experimental optimization. In enterprise environments, the primary question is not whether fine tuning improves linguistic fluency, but whether it can be introduced without weakening existing guarantees around correctness, auditability, and operational control. As such, the methodology presented here

treats fine tuning as a governed system change that spans data management, model adaptation, validation, and ongoing oversight. This approach reflects practices observed in mature enterprises prior to March 2024, where AI initiatives were required to align with established change management and risk assessment frameworks. The first methodological step is use case qualification and risk classification. Not all transaction and operations scenarios are appropriate candidates for fine tuning. Enterprises must first determine whether the intended use case involves interpretive assistance, explanatory support, or decision influencing behavior. Fine tuning is methodologically inappropriate for scenarios where outputs directly trigger financial transactions or irreversible operational actions without human validation. Risk classification considers factors such as data sensitivity, regulatory exposure, tolerance for probabilistic variance, and the potential impact of subtle semantic errors. This step ensures that fine tuning is applied only where its benefits outweigh the structural risks introduced by embedding domain patterns into model behavior.

The second step involves controlled data selection and transformation. Transaction and operations data are rarely suitable for direct inclusion in training pipelines without modification. The methodology requires explicit criteria for data eligibility, including temporal relevance, policy alignment, and representativeness of current operational norms. Historical data reflecting deprecated processes, emergency workarounds, or anomalous events must be identified and either excluded or abstracted. Transformation techniques such as schema generalization, identifier removal, and aggregation are applied to reduce the risk of memorization and inferential leakage. Importantly, all data preparation steps are treated as reproducible processes, with versioned outputs and documented rationale, enabling later audit and review. The third step is bounded fine tuning execution. Rather than pursuing maximal behavioral change, enterprises apply conservative fine tuning strategies that limit deviation from the base model. This includes restricting training duration, constraining learning rates, and favoring methods that emphasize adaptation to domain language patterns over procedural inference. The methodology assumes that fine tuning should enhance familiarity, not encode authoritative decision logic. Training runs are executed within controlled environments, with detailed logging of configurations, inputs, and outputs. Each resulting model artifact is uniquely identified and associated with an approval record, reinforcing the notion that fine tuning represents a formal system modification rather than an experimental iteration.

The fourth step focuses on multi dimensional validation. Traditional machine learning evaluation metrics are insufficient for transaction and operations contexts, where the most significant risks involve misinterpretation and policy misalignment rather than syntactic errors. Validation therefore combines scenario based testing, policy adherence checks, and stability analysis. Scenario testing uses representative transaction workflows and operational incidents to evaluate whether the model produces consistent and contextually appropriate responses. Policy checks assess whether outputs respect domain constraints, such as avoiding prescriptive financial advice or unauthorized procedural guidance. Stability analysis examines how small input variations affect outputs, providing insight into behavioral volatility introduced by fine tuning. This validation phase

produces qualitative and quantitative evidence that supports a defensible deployment decision. The final step is integration with operational governance. Once deployed, fine tuned models are continuously monitored within their intended usage boundaries. Inference interactions are logged with sufficient context to support retrospective analysis, and deviations from expected behavior trigger review workflows. The methodology requires periodic reassessment of model alignment as transaction logic, operational processes, or regulatory interpretations evolve. Retraining or retirement decisions are made through the same governance channels that approved the original fine tuning. This closed loop methodology ensures that fine tuning remains an accountable practice embedded within enterprise system governance, rather than a one time optimization detached from operational reality.

## 10. Findings from Enterprise Experiments with Fine Tuned Models Using Transaction and Operations Data

Observations from enterprise experiments conducted prior to March 2024 indicate that the impact of fine tuning on transaction and operations use cases was uneven and highly dependent on governance maturity. Organizations that approached fine tuning as a controlled architectural intervention reported measurable improvements in domain fluency, terminology consistency, and reduction in prompt complexity. In these environments, models demonstrated improved ability to interpret transaction narratives, explain operational states, and summarize process outcomes without repeated contextual reinforcement. These gains were most visible in assistive scenarios, such as internal operations support, reconciliation analysis, and post incident reporting, where outputs were reviewed by humans before action was taken. In contrast, enterprises that treated fine tuning as a performance optimization step frequently encountered unexpected behavioral drift. Models fine tuned on large volumes of historical transaction data exhibited tendencies to over generalize exceptional cases, infer intent where none existed, or reproduce outdated procedural logic. These behaviors were not immediately apparent during initial validation but surfaced over time as the model encountered edge cases. This finding reinforced the insight that fine tuning amplifies historical signal, including anomalies and workarounds, unless data curation and validation are rigorously enforced. The risk was particularly pronounced in operations data, which often reflects system behavior under stress rather than normative conditions.

Another consistent finding related to transparency and explainability. While fine tuned models often produced more confident and fluent responses, stakeholders reported increased difficulty in explaining why a model responded in a particular way. This was especially challenging during audit reviews, where compliance teams sought to trace outputs back to specific data sources or policy references. Enterprises that relied solely on fine tuning without complementary retrieval or rule based validation struggled to provide defensible explanations. Conversely, organizations that paired fine tuned models with external validation layers and explicit usage boundaries were better able to contextualize model behavior within existing governance frameworks. Operational stability emerged as a differentiating factor. Fine tuned models deployed in environments with frequent policy changes required more frequent retraining to remain aligned. Enterprises underestimated the operational overhead

associated with retraining cycles, regression testing, and approval workflows. In some cases, the cost and delay associated with retraining eroded the initial benefits of fine tuning. This led several organizations to restrict fine tuning to domains with relatively stable semantics, while relying on retrieval based approaches for areas subject to frequent change. This selective application proved more sustainable than broad fine tuning across heterogeneous operational domains.

Observation Area	Low Governance Maturity Outcome	High Governance Maturity Outcome	Enterprise Impact
Domain Fluency	Inconsistent improvement	Stable terminology and reasoning	Productivity gains
Behavioral Drift	Undetected over time	Early detection and mitigation	Risk containment
Explainability	Limited traceability	Augmented with validation layers	Audit readiness
Operational Overhead	Underestimated retraining cost	Planned lifecycle management	Predictable operations
Policy Alignment	Implicit and fragile	Explicit and enforced	Compliance confidence

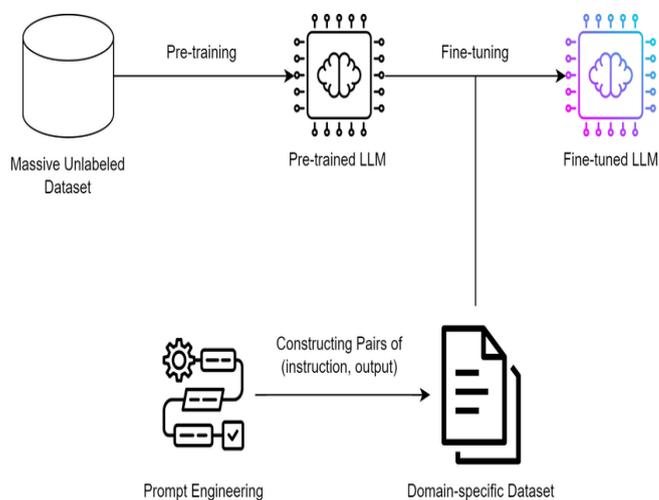
### 11. Challenges and Limitations of Fine-Tuning Large Language Models on Transaction and Operations Data

Fine tuning large language models on transaction and operations data introduces a set of challenges that extend beyond conventional concerns about model accuracy or performance. The most fundamental limitation arises from the probabilistic nature of large language models themselves. Even when fine tuning is performed on carefully curated enterprise datasets, the resulting behavior cannot be made fully deterministic. In transaction and operations environments, where correctness is often binary and errors can propagate downstream, this probabilistic behavior creates inherent tension. Fine tuning may improve fluency and apparent confidence, but it does not guarantee alignment with authoritative business logic. In some cases, increased confidence in responses can obscure uncertainty, making it harder for users to recognize when model outputs should be questioned or escalated. A second major challenge relates to data representativeness and historical bias. Transaction and operations data are produced within specific organizational, regulatory, and technological contexts. They often reflect temporary workarounds, legacy system constraints, or crisis driven behavior rather than idealized process design. When such data is used for fine tuning, models may internalize patterns that were contextually necessary at the time but are no longer valid. This risk is amplified when datasets span long time horizons without sufficient segmentation or temporal weighting. Unlike retrieval based approaches, which can prioritize current documentation, fine-tuned models blur historical boundaries, making it difficult to ensure that learned behavior reflects present day policy and practice.

Data sensitivity and irreversibility represent another structural limitation. Once transaction or operations data influences model weights, it becomes difficult to demonstrate that specific information has been fully removed or corrected. This poses challenges for regulatory requirements related to data deletion, correction, and subject rights. Even when explicit identifiers are removed prior to training, models may still encode correlations or behavioral signatures that raise concerns about inferential leakage. From a governance perspective, this irreversibility complicates audit narratives and increases the burden on organizations to justify why fine tuning was necessary and how risks were mitigated. In highly regulated environments, this alone may render fine tuning unsuitable for certain data classes. Validation and testing also present persistent challenges. Unlike traditional software systems or supervised learning models, fine tuned language models lack clear acceptance criteria that can be exhaustively verified. Transaction and operations use cases often involve nuanced interpretation rather than discrete outcomes, making it difficult to define comprehensive test suites. Scenario based testing can reduce risk, but it cannot cover the full space of possible interactions. As a result, organizations must accept a degree of residual uncertainty even after extensive validation. This limitation reinforces the need for human oversight and conservative deployment boundaries, particularly in workflows with financial or regulatory impact.

Challenge Area	Description	Impact on Enterprise Systems	Governance Implication
Probabilistic Behavior	Non deterministic outputs despite fine tuning	Risk of subtle transactional errors	Requires human oversight
Historical Bias	Learning outdated or exceptional patterns	Policy drift and misalignment	Demands data curation discipline
Data Irreversibility	Embedded influence of sensitive data	Compliance and audit difficulty	Limits eligible data classes
Validation Gaps	Incomplete test coverage	Residual uncertainty	Necessitates conservative scope
Operational Scalability	Growth of fine tuned variants	Increased maintenance burden	Requires centralized lifecycle control

Finally, organizational and operational constraints limit the scalability of fine tuning. Fine tuned models require ongoing lifecycle management, including retraining, regression testing, and approval workflows. As enterprises scale AI usage across multiple domains, the number of fine tuned variants can grow rapidly, increasing maintenance complexity and risk of inconsistency. Without centralized governance and reuse strategies, teams may diverge in how they apply fine tuning, leading to fragmented control and duplicated risk. These challenges suggest that fine tuning should be applied selectively and sparingly, guided by clear architectural principles rather than opportunistic experimentation.



## 12. Conclusion and Architectural Implications for Enterprise Adoption

The analysis presented in this paper demonstrates that fine tuning domain specific large language models using transaction and operations data is not a purely technical enhancement, but a structural architectural choice with long lasting consequences. When enterprises move from prompt based or retrieval augmented approaches toward fine tuning, they alter where domain knowledge resides, how behavior is controlled, and how accountability is enforced. Transaction and operations environments amplify these consequences because the data involved reflects binding business outcomes, regulatory obligations, and historically contingent processes. As observed prior to March 2024, organizations that underestimated this shift often encountered governance challenges that outweighed the immediate gains in fluency or usability. A central conclusion is that fine tuning should be treated as an act of institutional memory encoding rather than as model personalization. Transaction records and operational histories do not merely describe facts. They encode decisions made under constraints, exceptions tolerated under pressure, and interpretations shaped by temporal context. When these patterns are internalized by a probabilistic model, they persist beyond their original relevance unless actively governed. This persistence challenges traditional enterprise assumptions about change control, where logic can be updated deterministically and effects observed immediately. Fine tuning therefore demands stronger lifecycle discipline than many organizations initially anticipated, particularly in domains where policy and process evolve continuously.

The paper also highlights that architectural layering remains essential even when models are fine tuned. Fine tuning does not eliminate the need for prompt constraints, inference boundaries, validation layers, or human oversight. Instead, it shifts the balance of responsibility across these layers. Enterprises that preserved external control mechanisms were better able to contain risk and explain behavior, while those that relied on fine tuning as a substitute for governance experienced reduced transparency and increased audit complexity. This reinforces a broader architectural principle that probabilistic systems must be constrained by deterministic controls when embedded in high impact workflows. From an adoption perspective, fine tuning is most defensible when applied selectively to stable domains with well understood semantics and limited regulatory volatility. Use cases centered on explanation, summarization, and

analytical assistance align more naturally with fine tuned behavior than those involving direct execution or decision authority. In contrast, domains characterized by frequent rule changes, jurisdictional variation, or strict correctness requirements are often better served by retrieval based or hybrid approaches that preserve reversibility. This selective positioning allows enterprises to benefit from fine tuning without overextending it into areas where its risks cannot be adequately mitigated. Ultimately, the decision to fine tune domain specific large language models on transaction and operations data must be grounded in architectural realism and institutional responsibility. Fine tuning can enhance capability when supported by disciplined data governance, conservative validation, and explicit lifecycle control. It cannot replace foundational principles of enterprise system design, nor can it eliminate the need for human judgment in complex operational contexts. By framing fine tuning as a governed architectural choice rather than a technical optimization, enterprises can integrate large language models in ways that respect accountability, preserve trust, and align innovation with long term operational stability.

## 13. References

1. Dominik Kreuzberger, Niklas Kühn, Sebastian Hirschl. (2023) Machine Learning Operations (MLOps): Overview, Definition, and Architecture. *IEEE Access* 11: 31866-31879. <https://doi.org/10.1109/ACCESS.2023.3262138>
2. Nanchari N. (2022) Data Privacy And Security Challenges In Iot Healthcare. In *International Journal of Scientific Research & Engineering Trends*. 8(6). <https://doi.org/10.5281/zenodo.15796381>
3. Kranthi Kumar Routhu. (2019) Conversational AI in Human Capital Management: Transforming Self-Service Experiences with Oracle Digital Assistant. In *International Journal of Scientific Research & Engineering Trends*. 5(6). <https://doi.org/10.5281/zenodo.17678011>
4. Shraavan Kumar Reddy Padur. (2020) From Centralized Control to Democratized Insights: Migrating Enterprise Reporting from IBM Cognos to Microsoft Power BI" *International Journal of Scientific Research in Computer Science, Engineering and Information Technology (IJSRCSEIT)*. 6(1): 218-225. <https://doi.org/10.32628/CSEIT2390625>
5. Dana Alsagheer, Lei Xu, Weidong Shi. (2023) Decentralized Machine Learning Governance: Overview, Opportunities, and Challenges. *IEEE Access*. 11: 1-1. <https://doi.org/10.1109/ACCESS.2023.3311713>
6. Vinay Chamola, Vikas Hassija, A. Razia Sulthana, et al. (2023) A Review of Trustworthy and Explainable Artificial Intelligence (XAI). *IEEE Access*. 11: 1-1. <https://doi.org/10.1109/ACCESS.2023.3294569>
7. Kranthi Kumar Routhu. (2020) Strategic Compensation Equity and Rewards Optimization: A Multi-cloud Analytics Blueprint with Oracle Analytics Cloud. *KOS Journal of AIML, Data Science, and Robotics*. 1(1): 1-5. <https://doi.org/10.5281/zenodo.17531207>
8. Sudhir Vishnubhatla. (2020) Adaptive Real-Time Decision Systems: Bridging Complex Event Processing And Artificial Intelligence. In *International Journal of*

- Science, Engineering and Technology. 8(2).  
<https://doi.org/10.5281/zenodo.17471901>
9. Sudhir Vishnubhatla. (2021) Intelligent Loan Processing: Streaming, Explainability, and Customer 360 Platforms in Modern Banking. *Journal of Scientific and Engineering Research*. 8(2): 309-316.  
<https://doi.org/10.5281/zenodo.17639093>
  10. Shravan Kumar Reddy Padur. (2021) From Control to Code: Governance Models for Multi-Cloud ERP Modernization. In *International Journal of Scientific Research & Engineering Trends*. 7(3).  
<https://doi.org/10.5281/zenodo.17679693>
  11. Amina Adadi, Mohammed Berrada. (2018) Peeking Inside the Black-Box: A Survey on Explainable Artificial Intelligence (XAI). *IEEE Access*. 6: 52138-52160. <https://doi.org/10.1109/ACCESS.2018.2870052>
  12. Miloud Bagaa, Tarik Taleb, Jorge Bernal Bernabe, et al. (2020). A Machine Learning Security Framework for IoT Systems. *IEEE Access*. 8: 114066-114077.  
<https://doi.org/10.1109/ACCESS.2020.2996214>
  13. Ximeng Liu, Lehui Xie, Yaopeng Wang, et al. (2021). Privacy and Security Issues in Deep Learning: A Survey. *IEEE Access*. 9: 4566-4593.  
<https://doi.org/10.1109/ACCESS.2020.3045078>
  14. Nithin Nanchari. (2023) IoT for Mental Health Monitoring. *European Journal of Advances in Engineering and Technology*. 10(2): 75-77.  
<https://doi.org/10.5281/zenodo.15969008>
  15. Elif Ustundag Soykan, Leyli Karaçay, Ferhat Karakoç, et al. (2022) A Survey and Guideline on Privacy Enhancing Technologies for Collaborative Machine Learning. *IEEE Access*. 10: 97495-97519.  
<https://doi.org/10.1109/ACCESS.2022.3204037>
  16. Ahmed El Ouadrhiri, Ahmed Abdelhadi. (2022) Differential Privacy for Deep and Federated Learning: A Survey. *IEEE Access*. 10: 22359-22380.  
<https://doi.org/10.1109/ACCESS.2022.3151670>
  17. KM Jawadur Rahman, Faisal Ahmed, Nazma Akther, et al. (2021) Challenges, Applications and Design Aspects of Federated Learning: A Survey. *IEEE Access*. 9: 124682-124700.  
<https://doi.org/10.1109/ACCESS.2021.3111118>