



Integrating Large Language Model APIs into Enterprise Backend Services: Design Patterns and REST API Considerations

Prem Reddy Nomula*

Independent Researcher, India

* **Corresponding Author Email:** nomulapremreddy@gmail.com - **ORCID:** 0009-0001-6660-7724

Article Info:

DOI: 10.22399/ijcesen.5195

Received : 06 July 2023

Revised : 18 August 2023

Accepted : 09 September 2023

Keywords

Large Language Models,
REST APIs,
Enterprise Backend Architecture,
Design Patterns,
Observability

Abstract:

Large language models (LLMs) have rapidly evolved from experimental research artifacts into production-grade services that are widely accessible through RESTful APIs. Although these interfaces share structural similarities with conventional web services, their underlying characteristics—such as probabilistic outputs, token-based cost models, and high, variable latency—introduce fundamentally new challenges for enterprise system design. This paper presents a structured taxonomy of integration patterns for incorporating LLM APIs into enterprise backend architectures. Adopting a design science research methodology, the study derives and formalizes five core patterns—Gateway, Prompt Template, Retry and Fallback, Streaming Response, and Context Window Management—each addressing a distinct set of engineering concerns specific to LLM-based systems. In addition, the paper provides a comparative analysis of LLM APIs and traditional REST services, highlighting key architectural divergences. It further proposes observability and data governance strategies tailored to the operational realities of LLM integration. The applicability of the proposed patterns is demonstrated through a case-based validation, and practical implementation guidance is provided within the context of Java and Spring Boot environments. Together, these contributions offer a comprehensive framework for designing scalable, reliable, and compliant enterprise systems that leverage LLM capabilities.

1. Introduction

Transformer-based language models have fundamentally altered the landscape of software systems by enabling a transition from narrowly defined, task-specific machine learning solutions to flexible, general-purpose language interfaces capable of handling diverse cognitive tasks (Vaswani et al., 2017; Brown et al., 2020). The availability of these models through managed API services—most notably those provided by OpenAI, Anthropic, and Google—has significantly reduced the infrastructural and expertise barriers traditionally associated with deploying advanced AI systems. As a result, enterprises are increasingly embedding language model capabilities into backend services, integrating them into existing application ecosystems via REST-based interfaces. Despite this apparent continuity with established web service paradigms, large language model APIs introduce a set of architectural characteristics that diverge sharply from conventional RESTful

systems. Unlike deterministic APIs, LLM services produce probabilistic outputs that may vary across identical inputs, thereby challenging assumptions of repeatability and predictability (Bommasani et al., 2021; Bender et al., 2021). Furthermore, these systems rely heavily on context-sensitive inputs, where the structure and composition of prompts directly influence output quality. This stands in contrast to traditional schema-driven APIs, where input-output relationships are explicitly defined and tightly constrained (Fielding, 2000). In addition, LLM APIs are characterized by comparatively high and variable latency, as well as token-based cost models that directly impact architectural decisions related to request design and system scalability. These deviations necessitate a re-examination of backend integration strategies, as existing architectural patterns are insufficient to fully address the operational and design challenges introduced by LLM-based services. Accordingly, this paper investigates the following research question: How should enterprise backend systems

architecturally integrate LLM APIs while ensuring scalability, reliability, and compliance? By addressing this question, the study aims to bridge the gap between traditional service-oriented architectures and the emerging requirements of AI-driven backend systems.

2. Research Contributions

This paper makes three primary contributions to the emerging body of knowledge on enterprise integration of large language model (LLM) APIs. First, it proposes a formal taxonomy of LLM-specific backend integration patterns that extend and reinterpret classical enterprise integration patterns (Hohpe & Woolf, 2003) in the context of probabilistic, language-driven systems. By identifying recurring architectural challenges unique to LLM-based services—such as prompt management, context constraints, and model variability—the paper establishes a structured framework for reasoning about integration design.

Second, the paper presents a comparative analysis between LLM-based REST APIs and conventional RESTful services, highlighting key architectural divergences. This analysis focuses on dimensions such as determinism, latency, cost structure, and response representation, thereby demonstrating why existing backend design assumptions are insufficient when applied to LLM integrations. Through this comparison, the study provides a conceptual foundation for understanding the need for new architectural patterns.

Third, the paper introduces a design-oriented framework that is validated through a representative enterprise case scenario. This framework demonstrates how the proposed patterns can be systematically applied within a realistic backend architecture to address concerns of scalability, reliability, and compliance. By grounding theoretical constructs in a practical scenario, the paper bridges the gap between conceptual design and real-world implementation, offering actionable insights for practitioners and researchers alike.

3. Methodology

This study adopts a design science research methodology, which is widely recognized as an appropriate approach for developing and evaluating artifacts aimed at solving real-world engineering problems (Hevner et al., 2004). Within this framework, the proposed integration patterns are treated as design artifacts that address the practical challenges associated with incorporating large language model (LLM) APIs into enterprise

backend systems. The methodology emphasizes iterative construction, abstraction, and evaluation, ensuring that the resulting patterns are both theoretically grounded and practically applicable.

3.1 Pattern Derivation

The derivation of the proposed patterns is based on a multi-faceted analytical process. First, a comparative analysis was conducted on publicly documented API behaviors across major LLM providers, including those offered by OpenAI, Anthropic, and Google. This analysis focused on identifying common interaction models, limitations, and operational characteristics. Second, emerging LLM integration libraries and frameworks were examined to understand how early tooling attempts to abstract and manage these complexities within application architectures. Finally, the study synthesizes recurring challenges observed in industrial practice, such as handling rate limits, managing prompt construction, and dealing with context window constraints. Together, these sources informed the identification of recurring design problems and their corresponding pattern-based solutions.

3.2 Validation Approach

The validity of the proposed patterns is established through a combination of conceptual and practical evaluation strategies. A case-based architectural scenario is employed to demonstrate how the patterns can be applied cohesively within an enterprise backend system, illustrating their effectiveness in addressing real-world constraints. In addition, the patterns are assessed for consistency with established principles in distributed systems design, particularly those related to latency management, fault tolerance, and system resilience (Dean & Barroso, 2013; Nygard, 2007). This dual approach ensures that the patterns are not only contextually relevant to LLM integration but also aligned with foundational engineering practices.

4. The LLM API Integration Landscape

Large language model (LLM) APIs have largely converged toward RESTful interfaces, typically characterized by JSON-based payloads and standard HTTP semantics. This convergence provides a familiar integration surface for enterprise developers, enabling the reuse of existing web service tooling and architectural practices. However, beneath this apparent uniformity lies significant variation in deployment strategies and

operational characteristics, which introduce important architectural considerations for backend systems.

Two primary deployment models define the current integration landscape. The first consists of hosted APIs, where model inference is performed on external infrastructure managed by providers such as OpenAI, Anthropic, and Google. These services offer ease of integration and elastic scalability but are accompanied by challenges related to network latency, rate limiting, and data governance. The second model involves self-hosted or privately deployed systems, often based on open-weight models such as those derived from Meta's LLaMA family. While this approach provides greater control over data handling and system behavior, it requires substantial investment in infrastructure, optimization, and maintenance (Zhang et al., 2022). This duality between hosted and self-managed deployments necessitates the adoption of architectural strategies that emphasize abstraction and portability. Enterprise backend systems must be designed to accommodate variability in provider capabilities, performance characteristics, and compliance requirements, thereby reinforcing the need for flexible integration patterns and decoupled system design.

5. Core Design Patterns

This section presents a structured taxonomy of core design patterns for integrating large language model (LLM) APIs into enterprise backend systems. Each pattern is formalized using a consistent structure comprising context, problem, solution, and consequences. Collectively, these patterns address recurring architectural challenges arising from the probabilistic, high-latency, and context-sensitive nature of LLM services. Table 2 summarizes the proposed patterns, the specific problems they address, and the associated trade-offs, providing a consolidated view of their roles within system design.

5.1 Gateway Pattern

Context

Enterprise systems often interact with multiple LLM providers, each exposing heterogeneous APIs with evolving semantics and operational constraints.

Problem

Direct integration with individual providers results in tight coupling, leading to vendor lock-in, duplicated logic, and fragmented observability across services.

Solution

To address these challenges, a gateway abstraction layer is introduced to encapsulate provider-specific interactions. This layer defines a unified interface for LLM operations, enabling backend services to remain agnostic to underlying providers.

```
-----
#Pseudo Code
public interface LlmClient {
    String generate(String prompt);
}

public class GatewayService {
    private LlmClient client;

    public String process(String input) {
        return client.generate(input);
    }
}
-----
```

Consequences

The gateway pattern enhances portability and enables centralized management of concerns such as logging, authentication, and rate limiting. However, it introduces an additional abstraction layer, which may increase system complexity and necessitate normalization of differing provider capabilities. This pattern extends established API gateway principles in microservices architecture while incorporating LLM-specific considerations such as token accounting and model selection (Newman, 2015).

5.2 Prompt Template Pattern

Context

In LLM-based systems, prompt construction directly influences system behavior, effectively acting as a form of dynamic programming interface.

Problem

Embedding prompts directly within application logic limits flexibility, complicates testing, and hinders iterative improvement.

Solution

Prompts are externalized as version-controlled artifacts, allowing them to be managed independently of application code and dynamically injected at runtime.

```
-----
#Pseudo Code
String template =
loadTemplate("summary_v2.txt");
String prompt = template.replace("{input}",
userInput);
-----
```

Consequences

This pattern enables systematic experimentation, reproducibility, and evaluation of prompt variations. It also facilitates collaboration between developers and domain experts. However, it introduces the need for prompt lifecycle management, including versioning, validation, and deployment strategies. The approach aligns with emerging research on prompt engineering and optimization (Liu et al., 2023).

5.3 Retry and Fallback Pattern**Context**

LLM APIs frequently exhibit transient failures due to rate limits, timeouts, or capacity constraints.

Problem

Such failures can degrade system reliability and disrupt user-facing services if not properly managed.

Solution

A combination of exponential backoff and fallback strategies is employed, where failed requests are retried and, if necessary, redirected to alternative models.

#Pseudo Code

```
try {
    return primary.generate(prompt);
} catch (Exception e) {
    return fallback.generate(prompt);
}
```

Consequences

This pattern improves system resilience and enables cost-aware routing by leveraging alternative models when appropriate. However, fallback strategies may introduce variability in output quality and consistency. The pattern builds upon established fault-tolerance techniques in distributed systems, adapted to the probabilistic nature of LLM services (Nygard, 2007).

5.4 Streaming Response Pattern**Context**

Many LLM APIs support incremental token generation, allowing responses to be streamed rather than delivered as a single complete output.

Problem

Traditional blocking request–response models increase perceived latency and degrade user experience, particularly for long-form outputs.

Solution

Reactive streaming architectures are adopted to process and forward token streams in real time.

#Pseudo Code

```
Flux<String> stream = webClient.post()
    .retrieve()
    .bodyToFlux(String.class);
```

Consequences

Streaming significantly enhances responsiveness and user experience by delivering partial results as they are generated. However, it requires the adoption of reactive programming paradigms and introduces complexity in error handling and state management. This approach aligns with event-driven system design principles (Kreps et al., 2011).

5.5 Context Window Management Pattern**Context**

LLMs operate within fixed token limits, constraining the amount of input context that can be processed in a single request.

Problem

Enterprise applications often involve large inputs, such as documents or conversation histories, that exceed these limits.

Solution

Techniques such as summarization, chunking, and retrieval are applied to transform inputs into manageable representations.

#Pseudo Code

```
if (tokenCount(input) > limit) {
    input = summarize(input);
}
```

Consequences

This pattern enables the processing of large-scale inputs while optimizing cost and performance. However, it introduces potential information loss and requires additional preprocessing pipelines. The approach builds on established research in text summarization and information retrieval (Nallapati et al., 2016).

6. REST API Considerations for LLM Systems

Although large language model (LLM) APIs are typically exposed through RESTful interfaces, their operational characteristics diverge significantly from those of conventional REST services. These differences arise from the probabilistic nature of language models, their computational requirements, and their reliance on unstructured inputs and outputs. As a result, traditional assumptions

underlying REST-based system design—such as determinism, low latency, and schema-driven communication—are challenged, requiring backend architectures to adopt new strategies for handling variability, performance, and cost.

6.1 Non-Deterministic Semantics

A defining characteristic of LLM APIs is their non-deterministic behavior. Unlike traditional APIs, which produce consistent outputs for identical inputs, LLMs generate responses based on probabilistic inference mechanisms. Consequently, repeated requests with the same prompt may yield different outputs, depending on factors such as sampling parameters and internal model states (Goodfellow et al., 2016). This variability necessitates that downstream systems are designed to tolerate ambiguity and incorporate mechanisms for validation, ranking, or post-processing of responses.

6.2 Latency Characteristics

LLM inference introduces significantly higher and more variable latency compared to typical microservice calls. While conventional REST APIs often respond within milliseconds, LLM responses may take several seconds, depending on input size, model complexity, and system load. This latency profile aligns with observations in large-scale distributed systems, where tail latency can dominate performance (Dean & Barroso, 2013). As a result, backend systems must adopt asynchronous processing models, incorporate timeout strategies, and, where possible, leverage streaming to improve perceived responsiveness.

6.3 Token-Based Cost Model

Another critical distinction lies in the cost structure of LLM APIs, which is typically based on token usage rather than per-request billing. Both input and output tokens contribute to overall cost, making prompt design a key architectural concern. Efficient prompt construction, context management, and response truncation become essential strategies for controlling operational expenses (Strubell et al., 2019). This cost-awareness influences not only individual API calls but also broader system design decisions, such as caching, batching, and model selection.

6.4 Schema-less Outputs

Unlike conventional REST APIs that return structured and schema-validated JSON responses,

LLM APIs produce unstructured natural language outputs. This lack of strict schema guarantees introduces challenges in parsing, validation, and integration with downstream systems. Developers must implement heuristic or rule-based approaches to extract meaningful information, and in some cases, employ secondary validation models or structured prompting techniques (Bender et al., 2021). This shift from schema-driven to interpretation-driven processing represents a fundamental change in API consumption.

Taken together, these characteristics demonstrate that LLM APIs, while superficially aligned with REST principles, diverge in ways that significantly impact backend system design. In particular, assumptions of statelessness, idempotency, and predictable response structures are weakened, necessitating the adoption of new architectural patterns and safeguards.

7. Reliability and Observability

The integration of large language model (LLM) APIs into enterprise backend systems necessitates an expanded approach to reliability and observability that goes beyond traditional service monitoring practices. While conventional systems primarily focus on availability, latency, and error rates, LLM-based services introduce additional dimensions such as token consumption, cost variability, and output quality. These factors require more comprehensive instrumentation and monitoring strategies to ensure system robustness and operational efficiency.

7.1 Metrics

Effective observability in LLM-integrated systems begins with the identification and tracking of key performance and cost-related metrics. Among the most critical are token usage, which directly influences operational cost; latency distributions, which capture variability in response times across requests; cost per request, reflecting the cumulative impact of token-based pricing; and error rates, including rate limit violations, timeouts, and service disruptions. Unlike traditional APIs, where cost is often fixed per request, LLM systems require continuous monitoring of usage patterns to maintain cost efficiency and scalability.

7.2 Instrumentation

To capture these metrics, backend systems must incorporate structured logging and instrumentation at the point of API interaction. This typically involves recording relevant parameters such as

token counts, response times, and computed cost values for each request. For example, a logging statement such as:

```
-----
#Pseudo Code
log.info("tokens={}, latency={}, cost={}", tokens,
latency, cost);
-----
```

enables systematic tracking of resource consumption and performance characteristics. In practice, such logs are often integrated with centralized monitoring systems and observability platforms, allowing for real-time analysis, alerting, and long-term trend evaluation. This level of instrumentation supports both operational debugging and strategic optimization of LLM usage.

7.3 Quality Monitoring

A distinguishing challenge in LLM-based systems is that failures are not always explicit. Unlike traditional services, where errors are typically signaled through status codes or exceptions, LLM outputs may appear syntactically valid while being semantically incorrect, incomplete, or misleading. As a result, reliability must be assessed not only in terms of system availability but also in terms of output quality. This necessitates the development of evaluation pipelines that incorporate automated checks, heuristic validation, or human-in-the-loop review processes. Feedback loops can then be used to refine prompts, adjust model selection, or trigger fallback mechanisms. Such approaches align with established practices in machine learning system monitoring, where hidden performance degradation is a recognized challenge (Amershi et al., 2019; Sculley et al., 2015).

Overall, ensuring reliability in LLM-integrated systems requires a holistic observability framework that combines traditional system metrics with model-specific indicators, enabling organizations to manage not only system performance but also the quality and cost of generated outputs.

8. Data Privacy and Compliance

The integration of large language model (LLM) APIs into enterprise systems introduces significant risks related to data exposure and regulatory compliance. Because LLM interactions often involve transmitting user-generated or proprietary data to external services, organizations must carefully manage how information is processed, stored, and shared. These concerns are particularly critical in domains subject to strict regulatory requirements, where unauthorized data disclosure

can lead to legal and financial consequences. As a result, backend architectures must incorporate robust privacy-preserving mechanisms to ensure that sensitive information is adequately protected throughout the LLM interaction lifecycle.

8.1 Data Sanitization

A fundamental strategy for mitigating data exposure risks is the implementation of data sanitization processes prior to API invocation. Preprocessing pipelines are designed to detect and remove or mask sensitive information, such as personally identifiable information (PII), confidential business data, or regulated content. Techniques may include rule-based filtering, anonymization, or tokenization of sensitive fields. By ensuring that only necessary and non-sensitive data is transmitted to external LLM services, organizations can significantly reduce the risk of unintended data leakage while maintaining the functional integrity of the application.

8.2 Controlled Deployment

Controlled deployment strategies provide an additional layer of data protection by restricting where and how LLM inference is performed. Organizations may opt for private deployments, on-premise solutions, or region-specific cloud services to comply with data residency and sovereignty requirements. This approach allows enterprises to retain greater control over data flow and storage, ensuring that sensitive information does not leave designated jurisdictions. While such deployments may introduce additional infrastructure and operational complexity, they are often essential for meeting compliance obligations in regulated environments.

8.3 Audit Logging

Maintaining comprehensive audit logs of LLM interactions is critical for ensuring accountability and traceability. This includes recording prompts, responses, timestamps, and relevant metadata, while applying appropriate redaction techniques to protect sensitive information. Audit logs enable organizations to monitor data usage, investigate anomalies, and demonstrate compliance with regulatory standards during audits. Furthermore, they support governance practices by providing visibility into how LLM systems are utilized across the enterprise.

Collectively, these strategies—data sanitization, controlled deployment, and audit logging—form a foundational framework for managing privacy and

compliance risks in LLM-integrated systems. They align with established principles in privacy-preserving data processing and machine learning security, emphasizing the importance of minimizing data exposure while maintaining system functionality (Dwork, 2008; Shokri & Shmatikov, 2015).

9. Case-Based Validation

To illustrate the practical applicability of the proposed design patterns, consider a customer support backend system that integrates a large language model (LLM) API to automate query resolution and assist human agents. In such a system, incoming user requests are processed through a backend service that orchestrates interactions with one or more LLM providers. The application of the Gateway Pattern enables the system to dynamically switch between providers, ensuring flexibility in model selection and reducing dependency on a single vendor. This abstraction also facilitates centralized management of logging, authentication, and rate limiting.

The Prompt Template Pattern is employed to standardize the structure and tone of generated responses, ensuring consistency across different types of customer interactions. By externalizing prompts, the system allows for iterative refinement of response quality without requiring modifications to core application logic. In scenarios where API limits or transient failures occur, the Retry and Fallback Pattern ensures service continuity by reattempting requests or routing them to alternative models, thereby mitigating disruptions caused by rate limits or service unavailability.

To address latency concerns inherent in LLM inference, the Streaming Response Pattern is utilized to deliver partial responses incrementally, improving perceived responsiveness and enhancing the user experience. Additionally, the Context Window Management Pattern plays a crucial role in handling extended conversation histories, which are common in customer support workflows. By applying summarization or chunking techniques, the system ensures that relevant context is preserved while adhering to token constraints.

Collectively, this scenario demonstrates how the integration patterns operate in concert to address real-world constraints associated with LLM-based systems, including high latency, cost sensitivity, and output variability. The coordinated application of these patterns enables the construction of a robust, scalable, and user-centric backend architecture.

10. Java and Spring Boot Implementation

The Java and Spring Boot ecosystem provides a practical and mature foundation for integrating large language model (LLM) APIs into enterprise backend systems. Given the widespread use of Spring-based architectures in enterprise environments, the framework offers a familiar programming model for implementing the patterns discussed in this paper. In particular, Spring's support for dependency injection, REST client abstractions, and reactive programming makes it well suited to the demands of LLM integration, where systems must handle external API communication, provider flexibility, and streaming responses efficiently.

10.1 REST Integration

At the level of HTTP communication, Spring provides two principal mechanisms for interacting with external LLM services. The first is RestTemplate, which supports synchronous and blocking request-response interactions. This approach is suitable for simpler use cases where immediate, complete responses are acceptable and where the application architecture already follows a traditional blocking model. The second is WebClient, introduced as part of Spring's reactive stack, which supports asynchronous and non-blocking communication. WebClient is particularly advantageous for LLM-based systems because it enables efficient handling of long-running requests and token-streaming responses, making it more appropriate for modern AI-driven backend services.

10.2 Dependency Injection

A major advantage of using Spring Boot lies in its dependency injection mechanism, which allows developers to abstract provider-specific implementations behind common interfaces. This supports the Gateway Pattern by enabling flexible substitution of LLM providers without changing business logic. For example, an LLM client can be registered as a Spring bean:

```
-----
#Pseudo Code
@Bean
public LlmClient llmClient() {
    return new OpenAiClient();
}
-----
```

Through this approach, backend services can depend on the abstract LlmClient interface rather than a concrete implementation, thereby improving maintainability, testability, and portability across providers.

10.3 Reactive Streaming

For applications requiring incremental response delivery, Spring WebFlux provides strong support for reactive streaming pipelines. In particular, it enables server-sent events (SSE)-based communication, allowing partial outputs from an LLM to be relayed to frontend clients as they are generated. This capability is especially valuable for reducing perceived latency in interactive applications such as chat interfaces and customer support tools. By leveraging reactive streams, backend systems can process and transmit tokens efficiently while maintaining scalability under concurrent workloads.

10.4 Emerging Tooling

Alongside core Spring infrastructure, emerging libraries such as LangChain4j represent early attempts to standardize abstractions for LLM integration in the Java ecosystem. These tools provide higher-level support for prompt management, memory handling, retrieval-augmented generation, and provider abstraction, reducing the amount of boilerplate code required in enterprise applications. Although such libraries are still evolving, they signal the maturation of LLM integration practices and the growing recognition that traditional application frameworks require specialized extensions to accommodate language model workflows.

Overall, the Java and Spring Boot ecosystem offers both the foundational infrastructure and the extensibility needed to support enterprise-grade LLM integrations. By combining established framework capabilities with emerging LLM-specific tooling, developers can construct backend systems that are robust, adaptable, and aligned with contemporary enterprise requirements.

11. Discussion

The design patterns proposed in this study extend classical enterprise integration paradigms by

incorporating the unique characteristics of large language model (LLM) systems, particularly their probabilistic computation model and token-based cost structure. Unlike traditional service integrations, where behavior is deterministic and performance characteristics are relatively stable, LLM-based integrations require architectures that can accommodate variability in outputs, latency, and operational cost. The patterns presented—such as Gateway, Prompt Template, and Context Management—demonstrate how established software engineering principles can be adapted to address these emerging challenges while maintaining system modularity and scalability.

However, the adoption of these patterns is not without trade-offs. The introduction of additional abstraction layers, such as gateways and prompt management systems, increases architectural complexity and may impact system performance. Similarly, strategies designed to improve reliability, such as retry and fallback mechanisms, can introduce inconsistencies in output quality due to differences between models. Moreover, the inherent non-determinism of LLM outputs complicates testing, validation, and debugging processes, reducing the predictability that is typically expected in enterprise systems. These challenges highlight the need for careful system design and the incorporation of robust monitoring and evaluation mechanisms.

Looking forward, several avenues for future research emerge. One important direction is the development of formal methods for verifying the behavior of LLM-integrated systems, particularly in safety-critical or regulated domains. Another promising area is the automation of prompt optimization, where systematic techniques could be used to improve output quality while minimizing cost and latency. Additionally, further work is needed to establish standardized frameworks for evaluating LLM performance within backend systems, bridging the gap between machine learning evaluation metrics and software engineering quality attributes.

Table 1: Key architectural differences between conventional REST APIs and LLM-based APIs

Feature	Traditional REST APIs	LLM APIs
Determinism	Deterministic	Probabilistic
Latency	Milliseconds	Seconds (variable)
Response Structure	Structured JSON	Unstructured text
Cost Model	Per request	Token-based
Idempotency	Strong	Weak / non-deterministic
Output Validation	Schema-based	Heuristic / parsing required

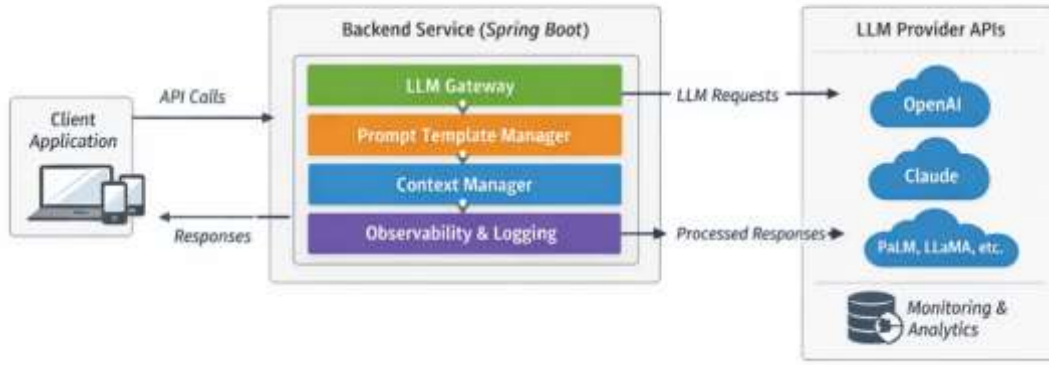


Figure 1: High Level Architecture of enterprise backend integration with LLM API's illustrating abstraction layers and supporting components

Table 2: Summary of proposed LLM integration patterns, their target problems, and associated trade-offs

Pattern	Problem Addressed	Solution	Trade-off
Gateway	Vendor lock-in	Abstraction layer	Added complexity
Prompt Template	Maintainability	Externalized prompts	Lifecycle management
Retry & Fallback	API failures	Backoff + alternative model	Output inconsistency
Streaming	High latency	Token streaming	Complex handling
Context Management	Token limits	Summarization/chunking	Information loss

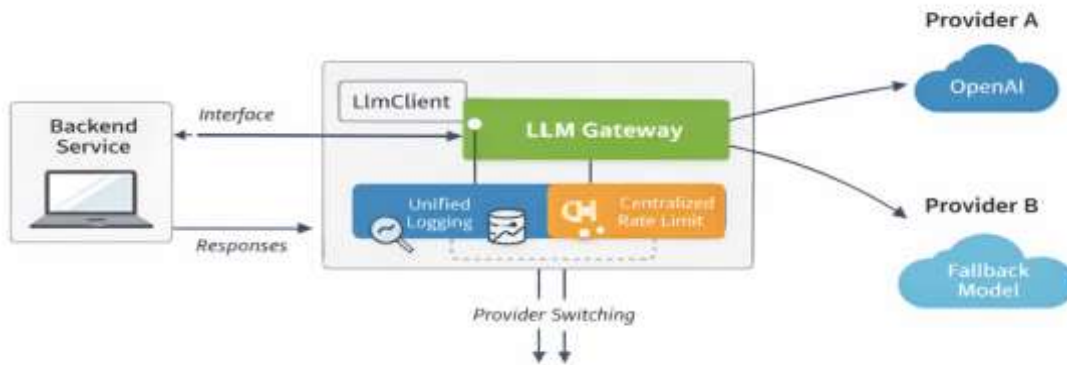


Figure 2: gateways pattern enabling abstraction over multiple LLM providers, supporting portability, centralized logging and rate management

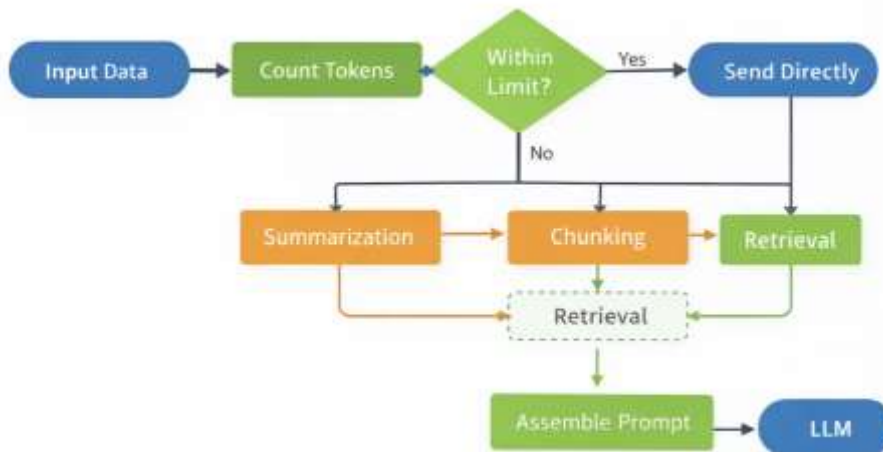


Figure 3: Context window management pipeline for handling token constraint using summarization and input transformation

Table 3: Illustrative trade-offs between latency, cost, and output quality across different LLM integration strategies

Strategy	Latency	Cost	Output Quality
Large model only	High	High	High
Fallback model	Medium	Medium	Medium
Summarized input	Low	Low	Slightly reduced

12. Conclusions

The integration of large language model (LLM) APIs into enterprise backend systems necessitates a fundamental rethinking of conventional REST-based architectures. While REST remains the primary interface for communication, the underlying characteristics of LLM services—such as non-deterministic outputs, high and variable latency, and token-based cost models—introduce challenges that are not adequately addressed by traditional design approaches. This paper has presented a set of structured design patterns that collectively provide a systematic framework for managing these challenges, enabling more effective integration of LLM capabilities into production systems.

The proposed patterns—spanning abstraction, prompt management, resilience, streaming, and context handling—demonstrate how established software engineering principles can be extended to accommodate the unique demands of LLM-driven applications. In addition, considerations related to observability, reliability, and data governance highlight the importance of adopting a holistic architectural perspective when deploying such systems in enterprise environments.

As LLM technologies continue to evolve, the patterns outlined in this work offer a foundational framework for designing backend systems that are scalable, reliable, and compliant. While future advancements may refine or extend these approaches, the core principles identified here are expected to remain central to the development of robust AI-enabled enterprise architectures.

Author Statements:

- **Ethical approval:** The conducted research is not related to either human or animal use.
- **Conflict of interest:** The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper
- **Acknowledgement:** The authors declare that they have nobody or no-company to acknowledge.
- **Author contributions:** The authors declare that they have equal right on this paper.
- **Funding information:** The authors declare that there is no funding to be acknowledged.
- **Data availability statement:** The data that support the findings of this study are available on request from the corresponding author. The data are not publicly available due to privacy or ethical restrictions.
- **Use of AI Tools:** The author(s) declare that no generative AI or AI-assisted technologies were used in the writing process of this manuscript.

References

- [1] Amershi, S., Begel, A., Bird, C., DeLine, R., Gall, H., Kamar, E., Nagappan, N., Nushi, B., & Zimmermann, T. (2019). Software engineering for machine learning: A case study. *Proceedings of the 41st International Conference on Software Engineering (ICSE)*.
- [2] Bender, E. M., Gebru, T., McMillan-Major, A., & Shmitchell, S. (2021). On the dangers of stochastic parrots: Can language models be too big? *Proceedings of the ACM Conference on Fairness, Accountability, and Transparency (FAccT)*.
- [3] Bommasani, R., Hudson, D. A., Adeli, E., Altman, R., Arora, S., von Arx, S., Bernstein, M. S., Bohg, J., Bosselut, A., & Brunskill, E. (2021). On the opportunities and risks of foundation models.
- [4] Brown, T. B., Mann, B., Ryder, N., Subbiah, M., Kaplan, J., Dhariwal, P., Neelakantan, A., Shyam, P., Sastry, G., & Askell, A. (2020). Language models are few-shot learners. *Advances in Neural Information Processing Systems (NeurIPS)*.
- [5] Dean, J., & Barroso, L. A. (2013). The tail at scale. *Communications of the ACM*, 56(2), 74–80.
- [6] Dwork, C. (2008). Differential privacy. *Proceedings of the International Colloquium on Automata, Languages, and Programming (ICALP)*.
- [7] Fielding, R. T. (2000). *Architectural styles and the design of network-based software architectures* (Doctoral dissertation, University of California, Irvine).
- [8] Goodfellow, I., Bengio, Y., & Courville, A. (2016). *Deep learning*. MIT Press.
- [9] Hevner, A. R., March, S. T., Park, J., & Ram, S. (2004). Design science in information systems research. *MIS Quarterly*, 28(1), 75–105.
- [10] Hohpe, G., & Woolf, B. (2003). *Enterprise integration patterns: Designing, building, and deploying messaging solutions*. Addison-Wesley.

- [11] Kreps, J., Narkhede, N., & Rao, J. (2011). Kafka: A distributed messaging system for log processing. *Proceedings of NetDB*.
- [12] Liu, P., Yuan, W., Fu, J., Jiang, Z., Hayashi, H., & Neubig, G. (2023). Pre-train, prompt, and predict: A systematic survey of prompting methods in natural language processing. *ACM Computing Surveys*.
- [13] Nallapati, R., Zhou, B., Gulcehre, C., & Xiang, B. (2016). Abstractive text summarization using sequence-to-sequence RNNs and beyond. *Proceedings of CoNLL*.
- [14] Newman, S. (2015). *Building microservices*. O'Reilly Media.
- [15] Nygard, M. (2007). *Release it!: Design and deploy production-ready software*. Pragmatic Bookshelf.
- [16] Rahman, M. M., et al. (2019). Configuration as code: Challenges and opportunities. *Proceedings of ICSE*.
- [17] Sculley, D., Holt, G., Golovin, D., Davydov, E., Phillips, T., Ebner, D., Chaudhary, V., Young, M., Crespo, J-F., & Dennison, D. (2015). Hidden technical debt in machine learning systems. *Advances in Neural Information Processing Systems (NeurIPS)*.
- [18] Shokri, R., & Shmatikov, V. (2015). Privacy risks in machine learning. *Proceedings of IEEE Symposium on Security and Privacy*.
- [19] Strubell, E., Ganesh, A., & McCallum, A. (2019). Energy and policy considerations for deep learning in NLP. *Proceedings of ACL*.
- [20] Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., Kaiser, Ł., & Polosukhin, I. (2017). Attention is all you need. *Advances in Neural Information Processing Systems (NeurIPS)*.
- [21] Zhang, S., Roller, S., Goyal, N., Artetxe, M., Chen, M., Chen, S., Dewan, C., Diab, M., Li, X. L., & Lin, X. V. (2022). OPT: Open pre-trained transformer language models. *arXiv preprint arXiv:2205.01068*.