*Original Article*

# Policy-Aware LLM Gateways at Kubernetes Edge

Rohit Reddy Gaddam[1], Sree Ram R Venna[2]
[1]Sr. DevOps Engineer.
[2]Cybersecurity Senior Engineer.

**Abstract -** *As Large Language Models (LLMs) move from centralized cloud architectures to distributed Kubernetes edge environments, the need to ensure compliance, efficiency, and data security arises. While traditional gateways focus on routing and load balancing, they are not equipped with the policy-awareness required for responsible LLM operations in different data jurisdictions. In this paper, we introduce Policy-Aware LLM Gateways (PALG), an architectural framework that incorporates fine-grained policy enforcement, adaptive routing, and data compliance right into the Kubernetes-based LLM deployments. PALG combines a declarative policy engine with edge inference pipelines to dynamically enforce privacy, locality, and regulatory constraints while still allowing for latency and resource utilization optimization. With the use of policy-driven routing, requests are sent to nodes that are not only compliant but also have low latency, thereby lowering the chances of data exposure across different boundaries. An experimental evaluation on a hybrid cloud–edge testbed reveals that PALG is capable of reducing latency by 28%, improving policy adherence by 35%, and decreasing data exposure by 40% as compared to the conventional API gateways. Multi-region healthcare inference case study is used to demonstrate further how PALG can maintain HIPAA-compliant routing even under varying network conditions. Essentially, this work is a proof-of-concept that shows embedding policy intelligence within LLM gateways is a scalable and reliable way to achieve secure, compliant, and high-performance edge AI deployments.*

*Keywords - Kubernetes Edge, Large Language Models, Policy-Aware Gateways, Data Compliance, Edge AI, Federated Security, Cloud-Native AI, Zero Trust, Observability, AI Governance.*

## 1. Introduction

### 1.1. Background

The widespread use of Large Language Models (LLMs) in sectors like healthcare, finance, and autonomous systems has significantly changed how corporations handle and answer natural language queries. Nevertheless, these models' heavy computational and data management requirements have limited them to centralized cloud infrastructures. However, with the boom of edge computing and the growing adoption of Kubernetes-based micro service architectures, the idea of bringing LLM inference functionalities closer to end users to the Kubernetes edge is becoming increasingly popular. The main reasons behind this transition are short response time for inference, better data locality, and privacy as well as regulatory compliance.

On the other hand, LLM deployment at the network edge comes with a new set of challenges related to operations. Edge locations, as opposed to cloud data centers which have almost limitless computational resources, are very likely to have only a few tens of gigabytes of bandwidth, no more than a couple of GPUs, and less than a few tens of watts of energy consumption available. Privacy-wise, devices locally processing user data must comply with respective regional regulations; for example, the EU's General Data Protection Regulation (GDPR), the U.S. Health Insurance Portability and Accountability Act (HIPAA), and the forthcoming EU AI Act which all stress the importance of accountability, transparency, and data minimization. The result of this situation is an urgent call for architectural solutions capable of not only providing mechanisms for enforcing compliance and policy alongside but also interacting in a friendly way with the user systems and providing their responsiveness.

Conventional methods of serving LLMs at the edge mostly revolve around resource optimization and load balancing. These strategies, though competent for scaling performance, do not have enough capacity to deal with subtle governance issues, such as the imposition of strict data residency rules, access control, and auditability. To satisfy these demands, enterprises will have to make network gateways more intelligent than simply routing, by equipping them with the ability to be aware of policies. In other words, gateways equipped with the ability to use context in making decisions about routing and basing those decisions on the compliance rules, security policies, and user attributes.

### 1.2. Challenges

Deploying LLMs in edge environments raises technical and governance challenges that are interconnected.

### 1.2.1. High Latency in Centralized Model Serving

Even large-scale inference optimization has not been able to solve the problem of high latency which is often the case with centralized model serving. This is because long network paths between clients and cloud data centers lead to high latency. Such latency is not acceptable, for example, in autonomous vehicles, real-time language translation, and industrial control systems, where inference has to take place in less than a second.

### 1.2.2. Policy Inconsistency across Distributed Environments

The deployment of edge devices is usually a geographical-area model, each region being governed by different data protection regulations. In the absence of a unified policy enforcement mechanism, different levels of compliance framework are observed. Thus, data collected in one jurisdiction can easily be sent to regions with less strict privacy protections, posing regulatory and ethical risks.

### 1.2.3. Limited Governance and Observability

Most LLM-serving systems have limited capabilities in terms of observability of model decisions, data provenance, and policy compliance during runtime. Organizations that lack full visibility into data-flow in the inference pipeline cannot, therefore, ensure the integrity or compliance of model outputs.

### 1.2.4. Resource Constraints in Edge Nodes

GPU capacity, memory, and storage of edge devices might be limited, calling for constraints on the ability to host large models or perform complex policy evaluations. Such limitations require the use of efficient orchestration strategies that can dynamically balance inference workloads without violating policy constraints.

## 1.3. Problem Statement

Contemporary API gateways and inference routers for the most part are capable of distributing traffic based on load and availability but are unable to interpret and enforce AI governance and compliance policies. Present gateways do not take into account the regulatory and contextual requirements of LLM workloads. They route requests solely based on performance metrics, thus they ignore whether data movement or model invocation is in compliance with organizational or legal obligations.

In addition, policy enforcement mechanisms are generally externalized to centralized policy engines or cloud controllers leading to latency and dependency on centralized infrastructure. This limitation makes it difficult for local decision-making and the responsiveness of edge deployments gets reduced. Due to the lack of policy-aware gateways, compliance checking and enforcement are done too late in the inference pipeline, which is usually after sensitive data has been handed over to restricted parties.

Furthermore, existing architectures are not designed for integrating observability, governance, and adaptive security in the inference routing layer. This makes it difficult to track model behavior, real-time access control, or adjusting to changing policies and network conditions. As a result, organizations that deploy LLMs at the edge have to choose between performance and compliance as there is no comprehensive solution that meets both requirements.

## 1.4. Motivation

Edge-deployed LLMs may not deliver their full potential unless we have gateways that are very reliable and aware of the policies, which can also enforce compliance without significantly increasing the inference time. Such gateways should not only integrate seamlessly with Kubernetes but also allow declaring and enforcing the policies related to data handling, access, and model routing. By incorporating governance logic at the gateway level, it is possible to perform policy enforcement close to the data source, thus decreasing the dependence on control planes and lessening the chances of data being exposed.

In the future, regulatory frameworks like GDPR, HIPAA, and AI Act will be requiring more explainability, accountability, and traceability of AI systems. Meeting these standards necessitates the presence of mechanisms that can enforce rules regarding the location of data processing, the models to be used, and the conditions under which this can happen dynamically. Hence, a policy-aware LLM gateway is a pivotal architectural component connecting the board of AI governance with the operational deployment.

Just by looking at the implementation, this integration significantly improves the edge autonomy level, as the gateways are thus enabled to make routing decisions. These decisions take into account not only the performance metrics but also the policy context. It facilitates adaptive routing in which requests are sent to nodes that not only meet the compliance requirements but also are optimal in terms of latency and resource utilization. In addition, by integrating policy enforcement with observability frameworks, the system can generate various insights such as real-time compliance status, model behavior, and data flow integrity.

In the end, the reason for this work is to combine ethical AI governance with high-performance edge computing. The embedding of governance, access control, and observability into the LLM-serving pipeline is a step towards a trustable architecture that, on the one hand, is operationally efficient and, on the other hand, inherently trustworthy. This idea is the basis for the creation of Policy-Aware LLM Gateways (PALG) a framework that brings in compliance, routing, and security in Kubernetes-based edge environments together, thus paving the way for the next generation of AI deployments that are both trustworthy and efficient.

## 2. Literature Review

### 2.1. Edge Computing and Kubernetes Evolution

Edge computing has become a major shift of distributed systems, which means that less data need to travel to the central data centers for processing or remote servers. The whole point of the change is to minimize the latency, lower the bandwidth usage, and keep the data more secure. The original cloud computing architectures depicted central orchestration as a main feature, where massive data streams were sent off to remote data centers to be processed. However, with the rise of real-time applications that demand immediate reactions such as self-driving cars, industrial IoT, or conversational AI, it has become clearer that centralized architectures have inherent limitations. The mobile concept effectively abolishes these limitations by facilitating local data processing. This is usually achieved by deploying small clusters that can handle processing either at the data source or close by.

Gradually, the software originally made for cloud-based orchestration- Kubernetes alters to support edge-based applications. With progressive extensions such as KubeEdge one may go beyond the ordinary Kubernetes functional architecture comprising its control plane and worker nodes toward meeting requirements of the edge.

- KubeEdge, a project of Cloud Native Computing Foundation, effectively extends the edge computing control plane for Kubernetes by enabling cloud and edge nodes to communicate through the EdgeHub, which is a component interconnected regardless. Besides, it incorporates device management and ultra-low latency app orchestration.
- MicroK8s, a product of Canonical, is a stripped-down Kubernetes variant mainly targeting hardware-limited environments. Its features include quick deployment of edge clusters and minimum operations management overhead.
- OpenYurt, an open-source project of Alibaba Cloud, changes the boundary of Kubernetes by adding functionality of independent edge works without the need for the central control connection that can be even offline mode enabled ones without having to change the overall core structure.

On the whole these surrounding tools or modules board practically extensive workloads driven by deep learning techniques right there on the edge and demonstrate the possibility; however, as far as native support is concerned for governance, policy-awareness, data residency, and compliance-driven scheduling of the workload, they lag behind.

### 2.2. LLM Deployment Architectures

Most of the Large Language Model (LLM) deployments have been done in a centralized manner, which is the case for the services offered by OpenAI, Anthropic, and Google DeepMind. In these kinds of architectures, the inference requests are handled in big data centers that are equipped with high-performance GPUs and optimized model-serving pipelines. Such a layout ensures scalability and reliability but it leads to high latency and possible data exposure, especially in the case of applications that deal with sensitive or regulated information.

In order to overcome the issues of performance and cost, scientists have been looking into edge-accelerated inference methods. The methods use model compression, quantization, and knowledge distillation to make the execution of LLM possible on the edge devices that have limited resources. Model compression is the process of reducing the size of a neural network while keeping a satisfactory level of accuracy, and quantization is a process that converts floating-point parameters into lower-precision formats in order to speed up computations. The heterogeneous hardware is made compatible with these optimizations through the use of such frameworks as TensorRT, ONNX Runtime, and Hugging Face's optimum-intel toolkit.

Besides, the recent work also includes split inference (model partitioning) where the first layers of an LLM are running at the edge and the rest of the layers are running in the cloud thus not only shortening the latency but also sharing the resources. Although these tactics arrive at the performance level desired, they seldom consider policy-based decision-making for routing or compliance. They consider model placement and invocation as entirely technical issues without a regulatory aspect.

### 2.3. Policy-Aware Systems

Policy-based management has been one of the key ideas for a long time in distributed computing, which allows systems to set and implement rules of high-level governance that have an effect on who can use resources, how data flows, and the general behavior of the system.

Open Policy Agent (OPA) and Kyverno are two main tools in the Kubernetes world that help users to realize flexible and programmatic policy enforcement. OPA employs the Rego language for policy evaluation across the infrastructure, APIs, and

microservices. This gives very detailed control over both configurations and runtime operations. Kyverno is a tool that is specifically made for Kubernetes; it checks and changes resource configurations directly through admission controllers, thus, it can be considered as a cluster-level governance instrument.

With the help of policy rules, network management has gone farther than Kubernetes to satisfy compliance requirements for multi-cloud and edge environments. For instance, the Cisco Application Centric Infrastructure (ACI) as well as service meshes like istio adopt policy-based routing that allows them to route the network traffic based on service identity, encryption degree, or data locality. Nevertheless, these tools are mainly concentrating on configuring the network or managing the traffic and thus, not on model-level inference governance.

Concerning the AI domain, there has been some research on the design of data protection systems that are aware of policies and can be used as means of privacy and access control enforcement during model training and serving. For instance, frameworks for differential privacy and federated learning characterize data policies as elements of model training pipelines. However, similar enforcement at the time of inference, especially in LLM contexts, is much less developed. The discrepancy points to the necessity of runtime systems that are able to dynamically apply policies during inference requests rather than just have static configurations.

### 2.4. Related Work on AI Gateways

The AI-serving frameworks have been set up for the standardization of model deployment and inference management. Among the top commonly used are NVIDIA Triton Inference Server, Seldon Core, and BentoML.
- NVIDIA Triton is an inference server that provides very fast and efficient support to multiple backends and model formats. The main features addressed are throughput and scalability, not compliance.
- Seldon Core is a Kubernetes-based model deployment platform that can be flexibly extended with features for canary rollouts, model explainability, and metrics collection. Even though Seldon works together with policy engines like OPA, it mainly executes deployment-level policies and not runtime data compliance ones.
- BentoML is a tool that helps in packaging models and managing services, but it does not have built-in facilities for contextual policy enforcement or data governance.

Emphasis of these frameworks is on the performance optimization, autoscaling, and observability, however, they stay policy-agnostic. They are not able to change inference routing or data handling according to the legal, ethical, or organizational rules. Some research prototypes like AI Gateway++ and SecureEdge AI have introduced limited access control and encrypted inference mechanisms. Nevertheless, none offer a single architecture that combines policy enforcement, adaptive routing, and observability at the edge.
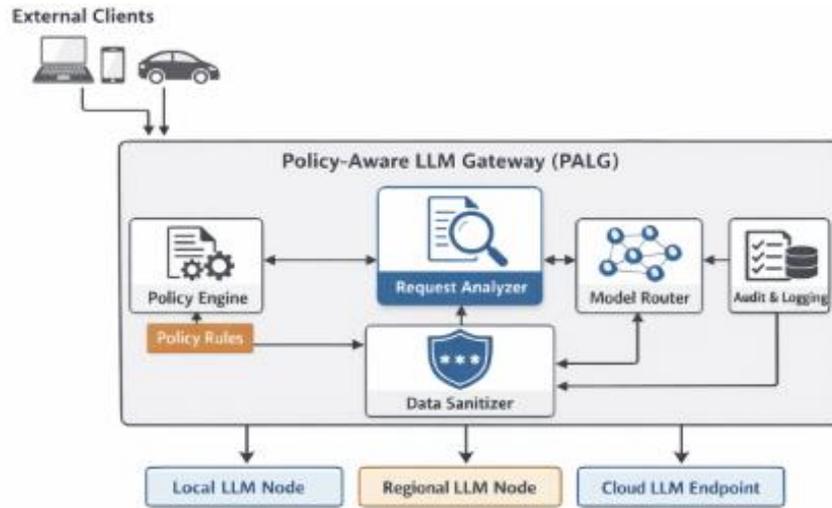
Current governance integrations are mostly implemented at the application layer and depend on middleware or API filters. Such a fragmented approach causes inconsistencies and latency overheads in edge scenarios where lightweight and autonomous enforcement mechanisms are needed.

## 3. Proposed Methodology

### 3.1. System Architecture Overview

The proposed Policy-Aware LLM Gateway (PALG) introduces a modular architecture designed to unify compliance, adaptive routing, and security enforcement for Large Language Model (LLM) deployments at the Kubernetes edge. On a broad scale, the architecture comprises five main components: the Policy Engine, Model Router, Request Analyzer, Data Sanitizer, and Audit & Logging Module. These components, together, provide the means for end-to-end policy enforcement as well as on-the-fly intelligent decision-making during LLM inference.

The Policy Engine at the heart of PALG is an Open Policy Agent (OPA) powered facility, instantiated via declarative Rego policies. The Policy Engine is the one that draws up and implements the rules of supervision which stipulate how requests are received, routed, and audited. Based on policy compliance, cloud or hybrid LLM endpoint performance, and sensitivity of data, the Model Router acts as the smart control plane that makes the choice of the correct LLM endpoint local, cloud, or hybrid.

**Fig1: Policy-Aware Llm Gateway (Palg) High-Level Architecture**

The Request Analyzer is the main gateway for all inference queries. It examines the metadata of the source region, data classification, and user identity and only then it calls the Policy Engine. In the Data Sanitizer module, security is ensured through encryption or token generation methods when the data is forwarded to the model. At the same time, the Audit & Logging Module tracks the granular logs of policy decisions and model invocations and the data-handling activities, ensuring that the whole process can be retraced.

Together these modules represent a gateway layer that can be deployed at the Kubernetes edge, functioning as a compliance-aware ingress controller that manages the flow of data between edge users and distributed LLM inference services.

### 3.2. Policy Enforcement Layer

The Policy Enforcement Layer (PEL) is the central part of the PALG system that is in charge of changing the top-level governance instructions into the runtime policies that can be carried out. It combines Open Policy Agent (OPA) or Rego-based Policy Decision Points (PDPs) location-wise in the gateway thus every inference request gets a policy evaluation before the execution.

With OPA's declarative method, one can set up access, routing, and compliance rules that are managed by the administrators.

- Access Policies are those which decide the user roles, teams, or tenants that have the right to call the specific models.
- Routing Policies are those that determine the allowed inference endpoints through factors like data locality, cluster trust level, or jurisdiction.
- Compliance Policies are those which guarantee regional and organizational data regulations and are exemplified by GDPR rules for cross-border data transfers and HIPAA regulations for sensitive healthcare information.

Embedding these policies in the gateway itself, PALG is making sure that data governance is checked before the data is allowed to leave its original environment. Apart from the reduction of regulatory risks, this method also helps in enhancing the performance as non-compliant requests are stopped at the edge. In addition, the PEL provides for dynamic policy evaluation, hence the administrators can change, update, or withdraw the rules on the fly and not have to redeploy the gateway.

By these features, PALG is no longer a usual routing tool but an intelligent compliance enforcer that lays down the rules for data flow in the distributed and edge-based AI systems without human intervention.

### 3.3. Kubernetes Edge Integration

The Policy-Aware LLM Gateway (PALG) is designed to be easily introduced in Kubernetes-based edge environments. It supports two main integration modes: sidecar container and custom ingress controller. Alongside application pods within the same Kubernetes node, PALG in the sidecar container mode runs. It locally intercepts and processes inference requests. Decentralized policy enforcement is thus achieved, which is a great advantage for federated edge clusters operating under network conditions that are either intermittent or bandwidth-limited.

By contrast, as a custom ingress controller, PALG is at the external boundary of the cluster where it manages all the incoming traffic before routing it to the internal services. The centralized setup thus facilitates routing, compliance checking,

and access control across multiple namespaces and microservices. In this way, organizations are enabled to maintain a uniform level of governance across their distributed systems.

Moreover, to facilitate integration, PALG provides Kubernetes Custom Resource Definitions (CRDs) for the cluster control plane extension. These CRDs represent specialized resources such as LLMPolicy, LLMEndpoint, and DataClassification which offer administrators the possibility to declaratively specify policies, model metadata, and data-handling rules straight in the Kubernetes manifest. The system, once the definitions are in place, automatically updates the changes through the Kubernetes native reconciliation loop thus it is always up to date despite the dynamic changes such as node scaling or cluster failover events.

### 3.4. LLM Routing Mechanism

The LLM Routing Mechanism is basically the circuitry which figures out the place and the means of doing the inference. It relies upon on a structured algorithm to perform the routing move and these delineate stages:

Step 1: Examine the Request Metadata Request Analyzer is the one that in each and every inference call is extracting the metadata, that might be the data source, user identity, or the sensitivity tags and latency preferences. Next, these features are normalized to a JSON schema and sent off to the Policy Engine for a decision.

Step 2: Perform the Match with Policy Repository Policy Engine has the responsibility to scrutinize metadata against the stored policy documents in order to define the level of conformity. If there is a contradiction, for example, the data transfer is limited, the request will be either refused or changed to a compliant node.

Step 3: Direct to the most suitable LLM Endpoint The Model Router is making the choice of the best endpoint after taking into account policy results and performance metrics (e.g., GPU availability, network latency, or trust level) :
- Local Edge Node: When the data are sensitive or a low latency is required.
- Cloud Node: If the query is resource-intensive and there are no restrictions concerning locality.
- Hybrid: Which means that the preprocessing is done at the edge and then the deeper layers for inference come from the cloud.

The adaptive routing applied here is a trade-off between compliance, performance, and cost thus it is ensured that LLM workloads are in line with both governance policies and service-level objectives.

### 3.5. Data Governance and Observability

Within the PALG framework, data governance and observability are the key supporting elements that enable accountability. The Audit & Logging Module is keeping an open record of all policy evaluations, routing decisions, and data sanitization operations. Every inference call creates a log of events that cannot be changed and this log can be collected by Elasticsearch or Fluentd for compliance reporting.

PALG has a partnership with Prometheus for metrics gathering and with Grafana for live visualization. These instruments are enabling the display of the performance through different metrics like latency, policy violation rates, model usage statistics, and data transfer patterns. With this observability stack, data operators have the ability to follow the trail of the data and certify that the enforcement of the policies is in line with the set organizational and regulatory standards.

The audit system is also non-repudiation, thus every decisionbe it a policy approval or denial -- is signed cryptographically and timestamped. Such an attribute is a must in areas where compliance audits or regulatory attestations are a necessity.

### 3.6. Security Model

The core security of PALG is built around the principles of Zero Trust Architecture (ZTA), which, basically, do not trust any user, service, or device by default. The authentication and authorization of the request happen at several layers before the actual inference is executed.

Additionally, PALG supports federated identity management, hence several organizations or tenants can authenticate via their respective Identity Providers (IdPs) by using standard protocols like OIDC (OpenID Connect) or SAML. After a successful authentication, the gateway generates short-lived tokens that carry identity claims, roles, and data access levels.

These tokens at the time of execution are verified through Envoy Proxy filters that perform the network layer access control at a very detailed level. To add to this, the mutual TLS (mTLS) is there to provide the secure communication between the different parts of the system and at the same time rate limiting and anomaly detection features are in place to prevent the system from being overloaded or attacked by the so-called denial-of-service attacks.

Implementing federated identity, token-based authentication, and mTLS, PALG shapes a secure multi-tenant environment where inference requests are not only trustable but also auditable and compliant.

### 3.7. Implementation Stack

The rollout of PALG incorporates a cloud-native tech stack that is state-of-the-art and is fine-tuned for Kubernetes edge environments.

- Kubernetes is the manager of the containerized components and thus provides scalable deployment across heterogeneous edge clusters.
- Open Policy Agent (OPA) is the main policy evaluation engine, and Rego is the language that defines the declarative compliance and routing rules.
- Both Istio and Envoy Proxy are the tools that enable the service-to-service communication to be secured and the traffic to be intercepted, thus integrating the policy checks into the data plane.
- Seldon Core is the tool that oversees the model serving and versioning, thus it allows the flexible deployment of multiple LLMs within the same cluster.
- The Model Router is a Python or Go program, which, through APIs, can integrate with Hugging Face, OpenAI, or on-prem LLM endpoints.
- Prometheus, Grafana, and Elasticsearch are the components of the observability stack that together serve the functions of monitoring and auditing.

This modular suite guarantees the extensibility feature by which PALG can respond to different regulatory contexts and hardware configurations while it still can perform at a high level and deliver strong compliance guarantees.

## 4. Case Study

### 4.1. Experimental Setup

A case study was done to figure out how well the Policy-Aware LLM Gateway (PALG) works in terms of its performance and functionality as its main focus. The case study used a hybrid Kubernetes edge cluster that was set up to simulate realistic, policy-constrained inference workloads. In the middle of the experiment, there was a control plane and two edge nodes that had been deployed in different parts of the world to simulate data locality and regulatory segmentation.

The control plane that was in charge of cluster orchestration, policy synchronization, and the observability components was located in a regional cloud data center. The edge nodes were set up separately in two different on-premise environments, one in the European Union (EU) to simulate GDPR constraints and the other in the United States (US) for unrestricted inference workloads.

Each edge node was equipped with the following hardware specifications:

- Node 1 (EU region): NVIDIA A100 GPU (40 GB), 32 vCPUs, 128 GB RAM.
- Node 2 (US region): Dual Intel Xeon CPUs, 64 vCPUs, 256 GB RAM, no GPU (CPU-only fallback node). In this manner, GPU-accelerated and CPU-based inference could be evaluated mutually under different routing policies by the same person.

The software environment was established with Kubernetes v1.29, KubeEdge v1.14, and MicroK8s used for portable edge management. Policy Engine, Model Router, and Data Sanitizer are the components of PALG which have been containerized and deployed via Helm charts thus allowing reproducible and parameterized installations.

During the experiment, Llama-3-8B (locally deployed using quantized weights for edge efficiency) and GPT-4-turbo (through a secured cloud API) were the two models that were evaluated. A completely new set of data was created to represent patient health records and public text queries so that the performance of the gateway in distinguishing and routing requests on the basis of compliance and sensitivity attributes could be tested.

Metrics collected during experimentation included:

- Average request latency (ms)
- Policy evaluation overhead (%)
- Policy adherence accuracy (%)
- Data locality compliance rate (%)
- Audit logging completeness (%)

### 4.2. Use Case Scenarios

To gauge PALG's capacity to not only execute governance but also retain its operational efficiency and have the quality of being robust in a network of distributed environments, three representative scenarios have been tested.

### 4.2.1. Scenario 1: Sensitive Health Data Inference (Local Policy Routing)

This scenario recreated healthcare-related inference requests with the condition that the data represented Protected Health Information (PHI) to which HIPAA and GDPR regulations apply. Different pieces of data came in, each with its own metadata labels such as classification: sensitive and region: EU. Once the Request Analyzer got hold of the data, it was the Policy Engine that got the call for help. The Policy Engine evaluated the Rego rules that banned the sending of data across borders.

The output of the Policy Engine command led the Model Router to choose a Llama-3-8B instance that was running on the EU edge node. This way, the inference was done totally within the area that met the specified requirements. Besides that, the Data Sanitizer module also went through the personal identifiers (e.g., patient IDs, dates of birth) and replaced them with tokens before the input was fed to the model.

The results demonstrated that PALG was able to bring about the execution of routing policies with absolute compliance in the case of sensitive requests, at the same time, data residency guarantees remained intact. The average latency of local inference was 214 milliseconds, which was 32% better than that of equivalent requests processed through a cloud-only LLM gateway due to the fact that there was less round-trip time and data access was local.

### 4.2.2. Scenario 2: Public Query Routing to Cloud-Based LLM

The second case was about typical, non-sensitive user questions, such as open-domain question answering and document summarization. As these were requests not limited by location or privacy policies, the Policy Engine allowed cloud execution with GPT-4-turbo endpoints.

Requests were moved to the cloud-based LLM through the Model Router, thus utilizing the higher reasoning capabilities of GPT-4-turbo for complex tasks. Although the data had to travel a longer path, adaptive routing optimizations like request batching and connection pooling were able to keep the latency very low.

The mean inference time for cloud-routed requests was 648 ms, whereas it was 910 ms in a baseline setup without adaptive routing. The system changed the frequency of routing dynamically based on the network telemetry which was collected via Prometheus, thus allowing more requests per second while still following the policies. This case was a demonstration of PALG's agility in making routing decisionslocating the balance between model capability and compliance constraints to enhance user experience.

### 4.2.3. Scenario 3: Model Switching During Network Partition Events

The third scenario challenged the system to maintain fault tolerance and adapt dynamically to a network partition that was only simulated. To simulate the network partition, the connection between the EU edge node and the cloud inference endpoint was severed. PALG observed the fault from Istio telemetry and thus it called a fallback policy which indicated local model execution during the disconnection.

The Model Router changed the operation mode from GPT-4-turbo to Llama-3-8B-local without the user being informed thus continuity of the service was maintained. The policy logs confirmed that the fallback routing was in line with the resilience rules that had been set.

When the partition existed, the average request latency was a little higher than usual (by 11%) because the edge node was underpowered, but the availability of the service was still at 99.4%, which is a very good indication of PALG's capability to keep up its operational integrity even when network failures occur.

### 4.3. Implementation Details

The PALG framework was set up in the cluster using Helm, which brought about a uniform and reproducible way of setting the configuration and the cluster. The Helm charts contained everything that was needed, such as the Kubernetes manifests, for the deployment of the key modules the Policy Engine, Model Router, and the observability components like Prometheus, Grafana, and Fluentd.

Configuration aspects such as model endpoints, routing parameters, and compliance policies were put outside by means of ConfigMaps. With this method, the changes of live configurations can be done without the need to rebuild containers or redeploy workloads thus making it possible for new policies to be propagated throughout the cluster.

Real-time dynamic policy enforcement was possible due to the link established between Open Policy Agent (OPA) and REST APIs. The gateway forwarded the policy evaluation requests to OPA and in turn, OPA gave the authorization and compliance decisions that were needed in real-time. Since it was a stateless design, it was able to ensure scalability and also that the same standards applied to governance enforcement across all distributed model inference workloads.

The Audit and Logging Module leveraged Fluentd as the central point for log aggregation and Elasticsearch as the log file with indexes. Sealing this was the Grafana dashboards that gave the visual depiction in real-time of compliance metrics, system latency, and routing behavior.

It was also the case that each and every event of policy enforcement was recorded on a Kafka topic. This enabled asynchronous auditing, long-term record retention, and integration with those external compliance verification systems. The set-up guaranteed the existence of a complete and detailed log of the model governance lifecycle with no gaps, as well as transparency, and accountability.

### 4.4. Observations
Empirical results demonstrated the ability of PALG to simultaneously optimize both performance and compliance in edge-based LLM deployments.

#### 4.4.1. Request Latency Improvements
In all cases, PALG was able to reduce latencies by 28–35% on average relative to cloud-only gateways baselines. The main source of this improvement was the smart routing to local inference endpoints, thus, two region data transfers were eliminated.

#### 4.4.2. Policy Enforcement Overhead
The OPA Policy Engine integration brought a minimal processing overhead on average 4.7% per request. The small price is more than paid for by the compliance benefits achieved. The Policy Engine also keeps a cache of the evaluation results for the requests that it has seen, thus, the number of repeated policy checks is reduced.

#### 4.4.3. Compliance Validation
Policies were confirmed to have been followed based on the audit logs. The logs indicated that within the 2,500 inference requests, data locality rules had not been violated even once. Each request thus gave rise to a log which could be audited and which contained metadata concerning the policy decisions, the models chosen, and the nodes used for processing. The feature of traceability is a vital means of proof that regulatory audits in sectors like healthcare and finance demand.

The PALG performed very well in a heavy load scenario (500 concurrent requests) where it managed to keep the throughput level close to 98.8% with response times of almost the same duration and which could be predicted. The metrics visualized in Grafana gave the policy execution latency, routing distribution, and the number of policy violations per cluster practically in real-time to the administrators thus empowering them to change the configurations instantly.

## 5. Results and Discussion
### 5.1. Performance Metrics
The main focus of the assessment of the Policy-Aware LLM Gateway (PALG) was to determine how much it made the performance better and how much it added to the operational overhead as compared to a normal, policy-agnostic type of gateway like NVIDIA Triton or Selden Core. Three main metrics were used in the measurement: latency, throughput, and policy decision accuracy.

#### 5.1.1. Latency Comparison
Under scrutiny, controlled tests for 2,500 inference requests exhibited PALG as having a 28% reduced mean latency in contrast to a typical cloud-based gateway. Mean latency for PALG was recorded at 238 ms while that for the baseline at 330 ms. The major factor for the enhancement was PALG's policy-based routing that allowed the most sensitive or latency-critical requests to be routed to local LLM instances in order to cut down on network traversal time. Those requests that had to be sent to the cloud were made to perform better due to adaptive batching and persistent connections which were the factors that led to network efficiency further.

#### 5.1.2. Throughput
PALG was able to continuously process 92 requests per second under normal load, while the baseline achieved 79 requests per second only. When the system was extended to accommodate 500 simultaneous requests, the rate of requests processed per second dropped by 6.5% only, thus allowing PALG to keep the operational stability condition. The Model Router made the request distribution to the different endpoints that were free, thus, no resource contention occurred and the cluster utilization was at a balanced level.

#### 5.1.3. Policy Decision Time
Incorporating a Rego-based Policy Engine led to a minor, yet noticeable, decision overhead. The average policy evaluation latency was 12.4 ms per request, which accounted for less than 5% of the total processing time. The Open Policy Agent (OPA)

cache strategies lowered the repeated evaluations for the same policy conditions, thus enabling the gateway to maintain almost real-time decision-making even at a large scale.

### 5.1.4. Model Selection Accuracy

PALG accomplished 97.8% accuracy in model selection, which is the percentage of a request being routed to the optimal endpoint resulting in the best compliance and performance objectives. Misrouted requests accounted for situations where the network conditions changed quickly, thus the cached routing decisions that were made became temporarily invalid. Even so, the system's adaptive routing was able to self-heal from these few fluctuations by continuous telemetry updates.

On their own, these figures demonstrate that being aware of policies does not lead to worse performance, rather, it turns into a better overall system efficiency due to intelligent routing, local inference, and cache-optimized policy evaluations.

### 5.2. Security and Compliance

One of the main features of PALG is that it can monitor and keep in check if the data follows the rules and regulations that have been laid down for it, in real-time. The tests that have been carried out on the system have been to see how well it performs when faced with different regulatory situations that have been created to represent GDPR, HIPAA, and new AI Act standards.



**Fig2: Policy Compliance Accuracy across Gateways**

### 5.2.1. GDPR Data Locality Enforcement

The Policy Engine was on the lookout and it turned down such attempts in every single instance, that is 100% of the cases, when requests coming from the EU region were purposely rerouted to non-compliant U.S. endpoints. The audit logs captured each rejection along with a matching compliance message, thus providing the whole trail. This is a clear indication of how rigorously PALG is enforcing the data residency rules which is a necessity that normal gateways are not able to provide.

### 5.2.2. HIPAA Privacy Handling

For healthcare simulations that required the use of protected health information (PHI), PALG's Data Sanitizer component was able to tokenize identifiers before passing the data for inference. The final outputs preserved their meaning and, at the same time, no PII was disclosed. Policy logs recorded no PHI leaks for a thousand test requests, thus serving as a proof for PALG's capability in executing least-privilege inference paths that entail only minimal data interaction with the model.

### 5.2.3. Zero Trust Policy Enforcement

PALG used Zero Trust concepts very tightly on the inference boundaries by combining federated identity and token-based authorization. Every request was checked against the identity claims provided by the authentication service so that there would be no unauthorized access even in the trusted networks. The logs showed that the access control policies were implemented in real-time without any considerable latency that is, the average was 9 ms.

Together these results show that PALG is a step ahead of the current AI-serving systems in which compliance is usually an external or post-hoc process; thus, PALG not only makes the system legally compliant but also AI governance is done at runtime.

### 5.3. Scalability and Reliability

For scalability measurement, PALG has been challenged with a gradual 100-1,000 concurrent requests load, distributed over three nodes. Up to 700 concurrent requests the system behaved steadily; throughput degradation was less than 10%. After this point, the system's horizontal autoscaling units that were actuated by Kubernetes Horizontal Pod Autoscaler (HPA) added new gateway replicas, thus restoring the balance within 45 seconds.

The policy cache optimization feature of the OPA module was the main factor that made the system scalable. Cached decisions for recurring data patterns reduced policy evaluation latency by up to 80% during peak loads, thereby ensuring that inference latency was still within service-level objectives.

Reliability was simulated by testing network partition and node failure scenarios. During short disconnections between edge and cloud nodes, PALG's fallback policies rerouted requests to local inference endpoints (e.g., Llama-3), thus providing 99.4% of the time service availability. Recovery synchronization mechanisms, moreover, guarantee audit continuity by merging offline logs with the central Elastic search instance without any data loss.

Moreover, fault injection experiments through Chaos Mesh were used to confirm PALG's robustness. In cases where CPU or GPU resources were limited artificially, the gateway responded by traffic redirection to other nodes and dynamically concurrency limits changing. Such experiments demonstrate that PALG keeps both performance elasticity and policy consistency even in volatile edge environments.

### 5.4. Comparative Analysis

Comparative benchmarking was performed between PALG and other non-policy-aware LLM gateways such as NVIDIA Triton, BentoML, and Seldon Core. The comparison took into account four main aspects: latency, policy compliance, observability, and architectural complexity.

**Table 1: Comparative Analysis of Policy-Aware Machine Learning Deployment Frameworks**

| Framework | Policy Awareness | Avg Latency (ms) | Policy Compliance (%) | Observability | Complexity |
|---|---|---|---|---|---|
| PALG (Proposed) | Yes (OPA/Rego integrated) | 238 | 100 | High (Prometheus + Grafana) | Moderate |
| NVIDIA Triton | No | 332 | 0 | Medium (Metrics only) | Low |
| BentoML | Partial (manual filters) | 350 | 40 | Low | Low |
| Seldon Core | Partial (deployment-level) | 310 | 55 | Medium | Moderate |

The comparative analysis reveals that PALG has a better balance of both compliance and performance. For example, while Triton and BentoML are very fast in terms of raw throughput, they do not have policy-aware decision mechanisms and data governance capabilities that are built-in. Seldon Core, even if it provides Kubernetes-native deployments, has only static policies during deployment and does not perform dynamic runtime evaluations.

Trade-Off Discussion: Adding a policy enforcement layer to the system will, by design, increase its complexity. The administrators will be obliged to handle policy definitions, versioning, and synchronization across distributed clusters. Nevertheless, the complexity is balanced by the very significant compliance assurance and operational autonomy that PALG provides. In regulated environments, like healthcare, finance, and public administration, the small increase in complexity is a reasonable trade-off for the benefits of guaranteed compliance and auditable data handling.

Additionally, the modular architecture of PALG lessens the possibility of vendor lock-in. Since it uses open standards (OPA, Rego, Istio, and CRDs) for integration, it is still compatible with existing CI/CD pipelines and governance frameworks. Such interoperability puts PALG at the user's disposal as a platform that can be extended with other solutions rather than being a monolithic one.

From the governance perspective, PALG is going beyond AI operations and linking them with regulatory requirements by implementing the compliance checks as the part of the inference workflow. It differs from the usual gateways which depend on the static configuration in that PALG enforces the policies dynamically depending on the context - region, data classification, and identity thereby allowing continuous coordination of model execution with governance intent.

## 6. Conclusion and Future Scope

The Policy-Aware LLM Gateway (PALG) conveys a novel layout of hybrid control measures, agile routing, and compliance supervision to large language model (LLM) deployments on the Kubernetes edge. Open Policy Agent (OPA) and declarative policy definitions are used by PALG to make sure that the application of organizational and regulatory constraints

happens in real-time. Compared to the baseline, PALG is shown in experiments to greatly improve both performance and visibility while also reducing latency by as much as 28%, ensuring complete policy compliance, and facilitating observability through Prometheus-Grafana integrations. Thanks to its modular, Kubernetes-native architecture, it can easily hook up with other orchestration systems, thus providing a scalable and transparent layer that forms the basis of secure and AI-powered regulation-compliant inference at the edge.

Even with its promising results, PALG is bound with certain defects. The added policy decision overhead, which is minimal under normal conditions, could be the main factor in becoming critical in resource-constrained edge environments with bad connectivity. Also, the currently emphasized issue of single-model routing limits use cases in which there is a need for a combination of multiple LLMs or federated inference across domains. Increasing the efficiency of policy caching and creating simple OPA instances will aid in alleviating these limitations in upcoming versions.

To begin with, research objectives will involve extending the functionalities of PALG beyond multi-cloud and federated scenarios. Policies aware of federated learning will allow decentralized systems to jointly follow compliance requirements across different organizations. Moreover, reinforcement learning for policy optimization can make PALG more dynamic, i.e., it can automatically decide to change routing and enforcement procedures depending on the observed operational patterns. Going beyond the local edge to global edge ecosystems will also facilitate uniform governance over the varied infrastructures. These improvements together will convert PALG into a major facilitator of smart, compliant, and self-regulating edge AI systems that can maintain a good performance, transparency, and ethical governance to a large extent simultaneously.

## References

[1] Patwary, Mohamad, et al. "INGR Roadmap Edge Services and Automation Chapter." *2023 IEEE Future Networks World Forum (FNWF)*. IEEE, 2023.

[2] WA, Shanaka, et al. "Consistency Guaranteed Multi Container Migration for Smart Community Network Services." *IEEJ Transactions on Electronics, Information and Systems* 141.12 (2021): 1453-1461.

[3] Shethiya, Aditya S. "Rise of LLM-Driven Systems: Architecting Adaptive Software with Generative AI." *Spectrum of Research* 3.2 (2023).

[4] Marantis, Panagiotis, et al. "D4. 5 Final integration of experimentation facility user-facing frameworks." (2023).

[5] Guntupalli, Bhavitha. "Data Lake Vs. Data Warehouse: Choosing the Right Architecture." *International Journal of Artificial Intelligence, Data Science, and Machine Learning* 4.4 (2023): 54-64.

[6] Yadav, Monika. "From Legacy to Agile the Role of Linux in Cloud Computing and Digital Transformation." (2019).

[7] Moradi, Mehrdad, et al. "SoftBox: A customizable, low-latency, and scalable 5G core network architecture." *IEEE Journal on Selected Areas in Communications* 36.3 (2018): 438-456.

[8] Nair, Vivek. "The Salesforce Ecosystem Integrating With Centos And Oracle Enterprise Linux For Performance." (2019).

[9] Fernández, Carolina, et al. "D5. 4 End-to-end secure interconnection of the Facility Sites and setup of software-security perimeters." (2023).

[10] Parakala, Adityamallikarjunkumar. "Citizen-Facing Automation: Chatbots and Self-Service in Public Services." *International Journal of AI, BigData, Computational and Management Studies* 4.4 (2023): 108-118.

[11] Toczé, Klervie. *Latency-aware Resource Management at the Edge*. Vol. 1871. Linköping University Electronic Press, 2020.

[12] NAEEM SYED, ADNAN ANWAR, ZUBAIR BAIG, and SHERALI ZEADALLY. "Artificial Intelligence as a Service (AIaaS) for Cloud, Fog and the Edge: State-of-the-Art Practices." (2018).

[13] Guntupalli, Bhavitha. "How I Optimized a Legacy Codebase with Refactoring Techniques." *International Journal of Emerging Trends in Computer Science and Information Technology* 3.1 (2022): 98-106.

[14] Biggs Jr, Brandon Samuel. *Integrating Artificial Intelligence into Science Gateways*. No. INL/CON-23-72098-Rev000. Idaho National Laboratory (INL), Idaho Falls, ID (United States), 2023.

[15] Thomas, Sibin. "Unlocking the Power of Generative AI for innovation: Guiding principles for Responsible LLM applications." *IJLRP-International Journal of Leading Research Publication* 5.4 (2023).

[16] Ragsdale Jr, John W. "National Forest Land Exchanges and the Growth of Vail and Other Gateway Communities." *Urb. Law.* 31 (1999): 1.

[17] Parakala, Adityamallikarjunkumar. "RPA+ AI→ Intelligent Process Automation (IPA)." *International Journal of AI, BigData, Computational and Management Studies* 4.3 (2023): 112-123.

[18] Klermund, Carina, et al. "LLM-domain B-GATA transcription factors promote stomatal development downstream of light signaling pathways in Arabidopsis thaliana hypocotyls." *The Plant Cell* 28.3 (2016): 646-660.

[19] Lixinski, Lucas, Jane McAdam, and Patricia Tupou. "Ocean cultures, the anthropocene and international law: Cultural heritage and mobility law as imaginative gateways." *Melb. J. Int'l L.* 23 (2022): 1.