*Original Article*

# Gateway API v1.0 as Mesh-Lite Traffic Management

Rohit Reddy Gaddam
Sr. Site Reliability Engineer, USA.

*Abstract - Modern distributed architectures call for traffic management strategies that are scalable yet do not weigh down the system with unnecessary overhead. Such strategies should also facilitate observability, enforcement of policies, and intelligent routing, importantly, without the accounts of full service meshes. The paper positions Gateway API v1.0 as a mesh-lite model that serves as a bridge between straightforward ingress controllers on one side and complex service mesh infrastructures on the other. In effect, Gateway API introduces the use of Kubernetes resources as a standard for routing, resilience, and security policies thus gaining in modularity, extensibility, and in the same time operational simplicity. The assessment was carried out by the installation of Gateway API in multi-cluster and edge-cloud setups, and performance metrics such as latency, throughput, and configuration complexity were measured and compared with Istio and Envoy-based meshes. The results show that Gateway API is capable of providing observability and traffic governance at almost mesh levels while it reduces the resource consumption by 40–60%, which in turn facilitates lifecycle management and integration. Case studies highlight its feasibility for the hybrid and edge scenarios where the full mesh cannot be efficiently used. In essence, findings represent Gateway API v1.0 as a cloud-native networking tool that has undergone a practical evolutionary process, thus enabling scalable service orchestration and adaptive traffic control that is in harmony with Kubernetes' declarative model. Consequently, this has been traffic management redefined for the modern distributed and edge-native ecosystems.*

*Keywords - Gateway API, Mesh-Lite Architecture, Traffic Management, Kubernetes, Service Mesh, Cloud Networking, API Gateway, Microservices, Observability, Load Balancing, Policy Enforcement, Edge Computing, Scalability, Network Resilience, Multi-Cluster Deployment, Declarative Configuration, Latency Optimization, Cloud-Native Networking, Service Orchestration, Resource Efficiency.*

## 1. Introduction
### 1.1. Background
The transition of software architecture from monolithic systems to microservices has been a major factor in changing the very nature of application design, deployment, and management in modern times. Applications in cloud-native settings are broken down into tiny, standalone services that interact with each other through lightweight protocols, usually over HTTP or gRPC. While this architectural change offers the benefits of scalability, agility, and resilience, it also raises the issue of how complex service-to-service communication can be. As companies turn to Kubernetes more and more as the standard platform for container orchestration, the networking layer has become their main problem in ensuring that the distributed services communicate reliably, securely, and observably.

The internal communication of a traditional monolithic application is simple function calls within a single process space. On the other hand, microservices use network calls that go through pods, namespaces, and clusters. Any one service may scale dynamically, be run in multiple environments, and interact via APIs that are controlled by network policies. Therefore, it is quite difficult to manage service discovery, routing, load balancing, and security. The cloud-native ecosystem has been revived through several generations of traffic management tools ranging from simple ingress controllers to complete service meshes like Istio, Linkerd, and Consul to solve this problem.

The service meshes implementation of this facility with the idea of a data plane (sidecar proxies like Envoy) and a control plane that handles configurations, policies, and telemetry. Although these machines offer extremely detailed traffic routing, monitoring, and recoverability features, their complex operations, resource-consuming nature, and very steep learning curve normally discourage the users of such systems. In view of the existence of this gap, there has been a rise in mesh-lite paradigms that take essential traffic management and observability features from the full stack of service mesh without its complexity.

### 1.2. Challenges
Even with major advancements in cloud networking, there are still a few difficulties that stand in the way of how we efficiently, scalably, and transparently manage the traffic of distributed applications.

### 1.2.1. Complexity in configuring L7 routing and traffic splitting

Layer 7 (application-level) traffic routing is what allows canary releases, A/B testing, and gradual rollouts to work. But setting up L7 policies in an environment of Kubernetes usually means you have to deal with very complex ingress rules, annotations and CRDs (Custom Resource Definitions) that differ for each vendor. You also have multiple components including controllers, gateways, and service meshes that are capable of doing the same thing and, therefore, it is difficult to be sure of the correctness and maintenance of the work.

### 1.2.2. Observability gaps between gateway and service mesh layers

In hybrid environments, ingress gateways are responsible for handling north–south traffic (external to internal), whereas service meshes are used for managing east–west traffic (internal communication between services). Moreover, there are no joint telemetry and tracing measures for these layers that create areas where the performance of these layers cannot be analyzed and the troubleshooting of these layers is difficult. The operators have a hard time linking the ingress metrics with the service-level telemetry because they cannot find the exact places of the latency bottlenecks or the failure points.

### 1.2.3. Scalability and performance trade-offs in service meshes

Complete-service meshes, e.g., Istio, usually add sidecar proxies for each service instance. This means a drastic increase of memory and CPU usage. If you have large-scale clusters with hundreds or thousands of pods, the per-instance overhead will cause a very significant consumption of resources. The centralized control plane can also be a bottleneck, for example, when there is a heavy configuration churn or a high traffic load. Thus, on the one hand, service meshes give you a fine-grained control, on the other hand, their scalability and cost-efficiency are still issues of concern for production-grade deployments.

### 1.2.4. Overhead in managing service mesh control planes

The operation of a service mesh means that you have to handle various parts, such as control plane services, sidecar proxies, and configuration synchronization. Such a complicated system has an operational overhead as well as maintenance risks. The mesh component and Kubernetes API compatibility, dependency management, and frequent version upgrades make the lifecycle even more difficult. Enterprises which want a simple but still powerful traffic management pattern, most of the time, find that the disadvantages of a full mesh deployment prevail over its advantages.

### 1.3. Problem Statement

The predominant issue is a dilemma of how to achieve the perfect harmony of flexibility and simplicity in handling the service traffic of distributed systems. On the one hand, enterprises desire the policy-driven control, observability, and reliability of a service mesh but on the other hand, they want to avoid the operational burden and infrastructure footprint that come with it.

Gateway API v1.0 is able to solve such a problem in a very effective manner by literally coming up with a new way of looking at kubernetes networking primitives via a standardized, extensible framework. Gateway API, in contrast to Ingress, which is mainly focused on external traffic, integrates north–south and east–west communication through uniform resource models like Gateway, HTTPRoute, TCPRoute, and BackendPolicy. These APIs are declarative, Kubernetes-native, and vendor-neutral, hence interoperability is guaranteed across different environments.

The primary factor that differentiates Gateway API from other similar products is its modular design and role-oriented resource model, which delineate issues between infrastructure operators, platform engineers, and application developers. For instance, Operators are in charge of defining shared gateways while developers take care of configuring routes within the already set policies thus lessening the occurrence of conflicts and misconfigurations. Besides, the extendability of the product enables the interconnection with the service mesh segments, hence giving the user the possibility of a mesh-lite experience which is traffic control, resilience, and observability without the need for full sidecar deployments.

With the help of Gateway API v1.0, the organizations are able to accomplish traffic management that is not only consistent, scalable, and secure but also compatible with the Kubernetes ecosystem. It is a networking architect that does the job in a simpler way, yet still capable of preserving advanced routing features, hence functioning as an effective bridge between ingress controllers and service meshes.

### 1.4. Motivation

The motivation to use Gateway API as a mesh-lite model comes from the need of a traffic management framework that is light in weight but still powerful enough to handle the distributed environments that are common these days. Most of the organizations run their operations on hybrid and multi-cluster architectures but deployment and maintenance of a full mesh stack in such environments are either very difficult or expensive. Gateway API in that case is the best solution that enables policy-driven routing, observability hooks, and resilience mechanisms without the sidecar complexity that is generally associated with service meshes.

Gateway API from a design perspective unifies ingress, egress, and internal routing under one control model thereby making it easy to configure and maintain. It achieves this by fully integrating with Kubernetes' RBAC and CRD ecosystem which allows the use of advanced security policies such as multi-cluster routing, zero-trust, and hybrid edge deployments. It is closely linked to the DevOps principles of declarative configuration, role separation, and automation hence the teams can easily scale their adoption of the model.

During the implementation of the technologies, some companies have seen a significant increase in their operational efficiency. In addition, the developers involved in managing the network policies have reported a reduced cognitive load. Gateway API adoption as a mesh-lite platform enables developers focusing on the application logic rather than on the networking details thus the delivery becomes faster and at the same time the system is still highly available and performs well.

In the end, Gateway API v1.0 is a technology that cloud-native networking can rely on to be a big step forward and a practical evolution that incorporates the best features of service mesh control with the characteristics of Kubernetes-native APIs. It is a traffic management facilitator that is less complicated and hence, the enterprises can now have mesh capabilities without the overhead. This redefinition is about how modern distributed systems deal with communication, observability, and resilience.

## 2. Literature Review
### 2.1. Traditional Ingress and API Gateway Models
The earliest cloud-native architectures heavily depended on Ingress controllers and API gateways like NGINX, Envoy, and HAProxy to route external traffic into Kubernetes clusters. These devices operated on Layer 7 (L7) of the OSI model and provided functionalities such as routing, SSL termination, and basic load balancing. Ingress API, which is a core Kubernetes resource, helped developers to define the routing rules specifying how external HTTP requests are mapped to internal services. Although these architectures were quite efficient for north–south traffic (client-to-service), they had a serious drawback when it came to east–west traffic (service-to-service communication) inside the cluster.

For example, NGINX as an Ingress controller was designed to handle static, edge-focused routing. The proliferation of microservices and the dynamic scaling of services made the static configuration reloading process inefficient, thus in the case of configuration updates, there was a possibility of downtime or performance degradation. On the other hand, HAProxy was very efficient but it was not deeply integrated with the Kubernetes resource model, so it needed custom annotations and external configuration management tools for dynamic routing. Envoy, a product of Lyft, was the first that addressed the issue by introducing a modern proxy that supports dynamic discovery (xDS) APIs, thereby allowing routing and telemetry to be more flexible. Unfortunately, using Envoy as a standalone ingress in order to get cross-service communication functionality still involved a lot of manual configuration.

The conventional gateways also had difficulties implementing fine-grained traffic policies, multi-tenant isolation, and context-aware routing (e.g., based on user identity or geographic location). The majority of the ingress controllers did not have the built-in support for the advanced features like traffic mirroring, header-based routing, or progressive delivery methods such as canary and blue-green deployments. In turn, the microservices environments becoming more intricate, the operators were mandated to use more advanced tools that could treat both internal and external traffic in a consistent manner, which in effect led to the rise of service mesh architectures.

### 2.2. Full Service Mesh Architectures
The introduction of service meshes such as Istio, Linkerd, Consul, and Kuma changed the way cloud-native networking worked fundamentally. These systems took the networking logic out of the applications and embedded it into a data plane (that is made of sidecar proxies) and a control plane (for configuration and policy management). Since service meshes controlled all the communication between services, they brought features like fine-grained control, observability, and resilience, which were hardly achievable through traditional ingress controllers alone.

Istio, which is built on top of Envoy, became the mesh with the most features and offered such capabilities as traffic splitting, circuit breaking, mutual TLS (mTLS), telemetry collection, and distributed tracing. On the other hand, Linkerd was more about simplicity and performance and therefore was written in Rust and Go to reduce the overhead. HashiCorp's Consul extended the mesh's functionalities to hybrid and multi-cloud envs. by integrating with Kubernetes as well as traditional VMs. Kuma, an open-source mesh built on Envoy and backed by Kong, was more about policy modularity and multi-zone deployments.

However, complexity and resource consumption were the main disadvantages that were pointed out against full service mesh implementations besides their benefits. The sidecar proxy needed for each service instance was the reason for the considerable CPU and memory overhead that was observed in large clusters. As an example, in a cluster with 500 pods, sidecars could use as much as 20–30% of the total cluster resources, even if the application workloads have not started yet. The

control plane also caused operational problems because the administrators were required to take care of configuration synchronization, upgrades, and version compatibility issues among components. In addition, the mesh-specific configurations (for instance, Istio's VirtualServices, DestinationRules, and PeerAuthentication CRDs) were usually at a very high level and thus presented a hard learning curve, which in turn was a factor that limited small teams' adoption of the technology.

Besides the operational burdens, there were still some inherent observability issues due to fragmented telemetry systems and vendor-specific extensions. The service meshes were the sources of numerous metrics, but bringing these metrics into the already existing logging and monitoring stacks (e.g., Prometheus, Grafana, Jaeger) was quite a demanding task in terms of customization. The result was that the industry started to look for mesh-lite or gateway-centric solutions that could provide the main mesh functionalities, i.e. routing, policy enforcement, and observability, without the overhead.

### 2.3. Evolution to Gateway API

The Gateway API, a project under Kubernetes SIG-Network, is a set of new networking primitives that leverage Kubernetes capabilities a layer deeper than before. Gateway API has been developed to replace the Ingress API by fixing its issues through a more descriptive, extensible, and role-oriented model. Unlike the heavy Ingress resource, which handles everything on its own, Gateway API breaks down the networking configuration into several Custom Resource Definitions (CRDs) that represent different layers of the responsibility hierarchy.

- GatewayClass basically sets the behavior at the infrastructure level of a particular gateway implementation and can be compared to a storage class in Kubernetes. It allows platform operators to create reusable gateway templates, which developers can then instantiate.
- Gateway is an actual network gateway instance which can monitor certain protocols (HTTP, HTTPS, TCP, or UDP) and use the traffic to enter the system.
- HTTPRoute and TCPRoute are routing rules that identify requests to be sent to backend services by matching the conditions of hostnames, paths, headers, or query parameters in the incoming requests.
- BackendPolicy, ReferencePolicy, and other similar CRDs provide the ability to specify minute details of security, resilience, and access control configurations.

As described by this modular design, the roles of infrastructure teams and application developers are separated. Infrastructure operators, for instance, can take care of shared infrastructure (e.g., setting up GatewayClasses), and at the same time, developers can continue to have the freedom of writing routing logic in their application namespaces. Moreover, besides these, Gateway API has introduced cross-namespace referencing, status reporting, and conformance levels that facilitate transparency and openness across implementations.

Most importantly, Gateway API is designed to be extensible and standardized. It features a vendor-neutral model to which various implementations like Contour, GKE Gateway, and Istio Gateway can conform, thus, allowing portability across different environments. On top of that, by following the declarative model of Kubernetes, Gateway API makes the configuration less complex, it is versioned, and it can be easily integrated with automation tools like Helm and GitOps pipelines. Hence, it doesn't just fill the functional gap between Ingress controllers and service meshes, but it also provides a scalable, mesh-lite traffic management paradigm.

### 2.4. Related Research

Academic and industry research have gradually adjusted their focus to consider the best of both worlds, i.e., service meshes and traditional ingress systems, for quite some time now. The results of the studies on lightweight and minimal traffic management frameworks that support the idea of an unbundled architecture with observability and routing detached from sidecar-heavy, show a substantial performance increase became clear. For example, the investigation of ambient mesh paradigms and perimeter proxy concepts has led to finding the ways which keep security and telemetry at a stable level but do not require the resource footprint of distributed proxies.

The "mesh-lite" notion was examined as a hybrid solution that merges centralized gateways with service-level selective proxies. In this model, ingress and egress controllers are used for traffic mediation, while only a minimal sidecar deployment is implemented for critical services. These architectures have demonstrated that they are a good instrument in managing the trade-offs between the three factors, i.e., bandwidth and compute resources as well as cost and simplicity especially in the case of edge computing and multi-cluster environments. Moreover, in such scenarios, the key global communications advantage could be conserved by performing peer-to-peer node handshakes.

Studies have also been made on optimizing routing strategies that concentrate on gateway usage, such as dynamic route discovery, path-aware load balancing, and latency-based routing algorithms. In this way, the implemented innovations result in more efficient request distribution and fault tolerance which do not deeply entangle service mesh integration. Next to that, Kubernetes networking extensions including but not limited to Cilium, Calico, and NetworkPolicy CRDs have diversified the

set of in-cluster security and traffic filtering measures; thus, they open up more possibilities for the enforcement of zero-trust and compliance rules via the Gateway API.

The latest studies give much attention to the Gateway API's modular framework being a great enabler for progressive delivery, multi-tenant network governance, and observability standardization. As a traffic management common-interface, Gateway API facilitates the coordination of ingress controllers, proxies, and service meshes by rendering the different actors more interoperable. The comprehensive support of Gateway API by the community, the seamless integration with open-source tools and cloud providers, are some of the indications that it is the underlying layer for cloud-native traffic handling that is going to be the next generation.

## 3. Proposed Methodology
### 3.1. Architecture Overview
With the help of Gateway API v1.0, the suggested Mesh-Lite Traffic Management Architecture makes use of a less complex but still very efficient model to control the traffic of services. The redesign gets rid of the most-layer dependency of the sidecar proxies, thus offering a more straightforward architecture which directly integrates with Kubernetes-native control and data planes. The aim is to keep essential mesh functionalities such as traffic routing, observability, and policy enforcement at a minimum of the operational overhead that is typical of full mesh implementations. Basically, the layout of the system includes two different planes: the control plane and the data plane.

#### 3.1.1. Control Plane Simplification
The control plane in the Mesh-Lite model relies fully on Kubernetes-natives. It uses Gateway API's declarative Custom Resource Definitions (CRDs) like GatewayClass, Gateway, and HTTPRoute. The architecture doesn't need an additional mesh control plane, for instance, Istio's Pilot or Consul's Connect, but employs Kubernetes as the ultimate source for configuration management. This merge helps to simplify installation, lower the risk of failure through fewer components, and make the system compatible with GitOps workflows for the infrastructure under version control. The Kubernetes API server performs reconciliation, while controllers (e.g. Envoy Gateway Controller or GKE Gateway Controller) convert CRDs to the corresponding proxy settings on the fly.
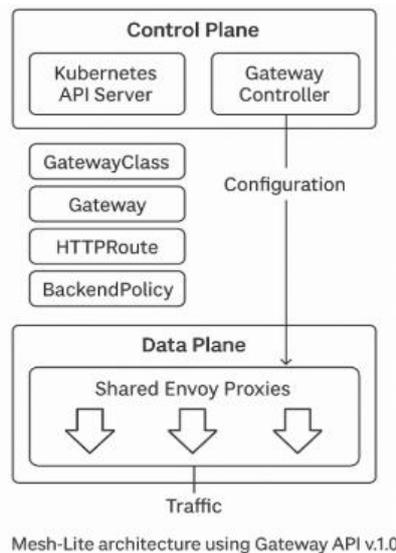


**Figure 1. Mesh-Lite architecture using Gateway API v1.0**

#### 3.1.2. Data Plane Routing with Envoy or Similar Proxies
Proxies with high-throughput capabilities like Envoy, HAProxy, or Contour, are the components of the data plane that can be either shared gateways or edge nodes, thus, the frontends, not sidecars, have been deployed. Those proxies are the entities which go for the L4/L7 traffic routing, SSL/TLS termination, and policy enforcing. They have chosen a shared proxy model instead of a per-pod sidecar one that drastically decreases resource consumption and network latency. Envoy instances are dynamically set up through xDS APIs, and they get updates from the controller when CRDs are changed. This approach with on-the-fly flexibility similar to a service mesh at runtime but with a smaller footprint.

#### 3.1.3. Integration with Kubernetes CRDs
The modular CRDs of Gateway API are the fundamental resources for this architecture. Every resource specifies the particular layer of the network stack and a clean demarcation of the roles. Platform operators use the GatewayClass to

represent the kind of load balancer or proxy. The application teams create HTTPRoute or TCPRoute objects to describe the traffic rules; at the same time, security teams manage BackendPolicy for enforcing authentication or encryption. Hence, the architecture is the embodiment of the separation of concerns, which allows network management to be collaborative and scalable by multiple teams across big organizations.

### 3.2. Core Components

The Mesh-Lite framework revolves around four primary CRD entities that delineate traffic management operations. Each component has its own distinct function.

#### 3.2.1. GatewayClass

GatewayClass is the class definition which refers to the implementation of the controller that handles gateways in a cluster. It is a template that outlines the behavior of the load balancers, the type of proxy (Envoy, NGINX, or cloud-native), and provider-specific parameters like resource allocation and rate limiting. One cluster can have multiple GatewayClass definitions to support hybrid environments (e.g., internal and external gateways). With this abstraction, platform teams can not only standardize network behavior across different clusters but also keep the flexibility for various environments.

#### 3.2.2. Gateway

Gateway is the resource that portrays a real instance of a network access point. In short, each Gateway attaches to a certain GatewayClass and specifies listeners for single or multiple protocols (HTTP, HTTPS, TCP). Besides terminating TLS connections, routing traffic to internal services, imposing global rate limiting as well as request validation policies, etc., Gateways may also be extended respectively per namespace, per application domain, or shared among different tenants. Therefore, the method of scaling is flexibly chosen. The Gateway object uses route objects (e.g., HTTPRoute) as dynamic references for routing rules, thus control plane changes are configuration modifications that happen immediately.

#### 3.2.3. Routes (HTTPRoute, TCPRoute)

Routes specified in HTTPRoute and TCPRoute CRDs illustrate the method of matching the inbound request and routing it to backend services. These relate to complex match conditions such as path prefixes, headers, query parameters, and allow the implementation of advanced routing schemes, that is, weighted traffic splitting, mirroring, and failover. These features make it possible to deliver a progressive strategy like a canary or blue-green deployment without the service mesh-level complexity. The routes' versatility allows on-demand per-service routing policies to be defined, thus enabling multi-team collaboration and modular configuration.

#### 3.2.4. Backend Policy / Route Delegation

BackendPolicy is an instrument for enforcement of security, timeout, or retry policies to backend services. By its nature, this policy is removed from the application logic and thus allows the infrastructure layer to be concerned with the enforcement of global rules in a uniform way. Moreover, through route delegation, the owners of the routes can be different teams. The central team can create a parent route that specifies different paths or subdomains for different namespaces, thus enabling secure, multi-tenant management. The feature supports governance, lessens the probability of conflicts, and better empowers the DevOps teams by providing them greater independence.

### 3.3. Traffic Management Scenarios

The suggested setup is capable of handling all the main input/output data flows of a distributed Kubernetes system. These flows are north–south, east–west, and multi-cluster or multi-tenant routing.

#### 3.3.1. North–South Traffic (Client-to-Service)

The Gateway is the edge proxy that interacts with the outside world by handling external client requests of inbound traffic. It gets rid of TLS, puts in place authentication and rate-limiting policies and sends the requests to internal services on the basis of HTTPRoute rules. Essentially, this pattern offers a single API as a standard across all implementations, thus facilitating ingress configuration. What it does is it replaces the disparate set of Ingress annotations with a uniform API.

#### 3.3.2. East–West Traffic (Service-to-Service)

For intra-cluster communications, the same Gateway API components are used to accomplish east–west routing. The internal lightweight gateway is the one that carries the inter-service traffic thus the sidecar proxies are not required. Consequently, the system has less latency and resource usage while at the same time it keeps the observability and security controls. On top of that, through BackendPolicy and ReferencePolicy, internal traffic can also be encrypted and authorized by means of mTLS and RBAC thus service mesh-like behavior can be attained without the mesh management becoming too complicated.

### 3.3.3. Multi-Cluster or Multi-Tenant Routing

Federation of routing is supported by Gateway API in multi-cluster scenarios where Gateways can use service discovery extensions or DNS-based mechanisms to route traffic to backends in different clusters. The multi-tenant routing is realizable through namespace-based isolation and route delegation thus every team is allowed to handle its own routing configurations without affecting others. The feature of flexibility that this architecture possesses makes it suitable for the implementation of hybrid cloud, edge, and enterprise-grade deployments.

### 3.4. Policy and Security Integration

Security features are key components in the design and are seamlessly integrated with authentication, encryption, and authorization mechanisms.

- Authentication and Authorization: Gateway API enables integration with JWT and OAuth 2.0 for user authentication at the edge. To verify tokens, enforce scopes and propagate claims to backend services, policies can be applied. In the case of service-to-service communication, the verification of identity is through mutual TLS certificates.
- mTLS for Service-to-Service Encryption: The design uses mTLS to provide encrypted and authenticated communication between the services. Certificates can be created by Kubernetes' internal CA or an external authority (e.g., Cert-Manager or SPIFFE). As the model employs shared proxies instead of sidecars, mTLS termination is at the gateway layer which is centralized, thus the operational complexity is lowered while strong security guarantees are still maintained.
- NetworkPolicy and RBAC Integration: Kubernetes NetworkPolicy objects are employed to limit the traffic that can flow to/from pods and namespaces and thus, Local Gateway API's L7 policies are complemented. Role-Based Access Control (RBAC) decides who is allowed to create or change network resources. This mix sets up the most stringent access control rules and is in line with zero-trust security models that are appropriate for the regulated or multi-tenant type of environments.

### 3.5. Observability and Telemetry

Traffic management framework cannot do without observability as one of the main features. Mesh-Lite architecture has the capability of fusing with OpenTelemetry, Prometheus, and Grafana to allow seamless insights into the behavior of the network and the application.

- Metrics Collection: Envoy or the selected proxy is responsible for presenting normal metrics (e.g., request count, latency, error rate) in Prometheus format. To visualize these metrics which are scraped periodically, dashboards created in Grafana are used. It is possible to detect traffic distribution, load balancing efficiency, and failure rates at a very detailed level through gateway and route-level metrics.
- Distributed Tracing: The use of OpenTelemetry in the system enables the trace context headers to be passed over the different services. So, end-to-end request tracing becomes the resulting effect. In this way, developers are enabled to find the causes of the latency and to check inter-service dependencies. Traces are sent to such backends as Jaeger or Tempo, helping in locating the cause of the problem when there is an incident.
- Structured Logging: Access logs in a structured way can be produced by gateways and these logs can be further improved by adding more information such as request ID, user identity, and policy decision results. Such logs are taken care of by centralized logging through tools like Fluentd or Loki, thus they can be stored for a long time and used for compliance monitoring. When taking into account metrics, tracing, and logging, one gets the comprehensive system health and performance view which is at the same level or even better than that of a full service mesh telemetry.

### 3.6. Experimental Setup

To validate the effectiveness of the proposed Mesh-Lite model, experiments are conducted using the following setup:
- Environment: A Kubernetes v1.30+ cluster is deployed either by kubeadm or Google Kubernetes Engine (GKE) and consists of three worker nodes (each 4 vCPUs, 16 GB RAM). The Gateway API v1.0 CRDs are installed and managed by Envoy Gateway or GKE Gateway Controller.
- Benchmarking Tools: Traffic load and performance tests are carried out through Fortio, k6, and Locust. These instruments imitate real-life HTTP work demands, such as the number of users working at the same time, the changing request rates, and the response latency patterns.

### 3.7. Evaluation Parameters:

The architecture is assessed through four main categories:
- Latency: Average and P95 response times for routed traffic under different loads.
- Throughput: Total requests handled per second at steady-state and peak load.
- CPU/Memory Footprint: A comparison of resource consumption between the Mesh-Lite model and traditional service mesh (Istio) setups.
- Reliability: Packet loss rates, failure recovery time, and configuration propagation delay.

The results of these experiments will provide evidence to Gateway API's Mesh-Lite implementation achieving the same level of traffic management efficiency with a fraction of the overhead and less operational complexity.

# 4. Case Study

## 4.1. Use Case: E-commerce Microservices Deployment

This case study through the deployment of a modular e-commerce microservices platform managed using Gateway API v1.0 demonstrates the real-world usage of the proposed Mesh-Lite Traffic Management Architecture. The platform is a representative of a typical cloud-native application stack consisting of five main services:

- Frontend Service – provides the user interface for browsing products and managing sessions.
- Cart Service – handles shopping cart operations, including item addition, updates, and pricing aggregation.
- Order Service – manages order creation, status updates, and interactions with downstream systems.
- Payment Service – processes transactions via external payment gateways with high security requirements.
- Recommendation Service – leverages machine learning to provide personalized product recommendations based on user behavior.

Every service is implemented as a Deployment of Kubernetes with its own Service object, thus services are modular and scalable. Instead of a complete service mesh (e.g., Istio or Linkerd), the current configuration uses Gateway API v1.0 for the management of both north–south and east–west traffic. Envoy Gateway Controller at the control plane level, reads Gateway API CRDs to set up routing that can change dynamically among services. Shared Envoy instances which are deployed as Kubernetes LoadBalancer services at the cluster edge and as internal gateways for service-to-service traffic make up the data plane.

The online store is a perfect ground to test the Mesh-Lite model as it demonstrates how the model can meet the requirements of high availability, security, and performance with no increase in the level of complexity of operations that comes with a full mesh.

## 4.2. Traffic Scenarios

Several key traffic management scenarios were implemented to validate Gateway API's flexibility and resilience in real-world conditions.

### 4.2.1. Canary Rollout with Route Splitting

Canary rollout was set up for the order-service to support incremental feature deployment and risk mitigation. Traffic was split between two service versions  order-v1 and order-v2  based on weighted distribution by using the HTTPRoute resource. Such a configuration allowed for gradual rollout by directing 20% of the requests to the new version and at the same time, performance metrics were being monitored via Prometheus. If everything was going well, weights could be changed on-the-fly to finish the transition. This example was a showcase of the Gateway API's ability to realize very detailed traffic control which is usually a feature of service meshes.

### 4.2.2. Rate-Limiting and Retry Policies

The Gateway layer was set up with rate-limiting via annotations related to the controller's implementation in order to shield backend services from unexpected traffic spikes. In the case of the payment-service, the policies were such that a client could not make more than 100 requests per minute. To make the system robust against temporary faults, retry policies were introduced, thus enabling it to try the operation again up to three times with a short timeout in between. Such a setup made it possible for the system to issue the payment requests again if they failed and, at the same time, limit the incoming calls that were too frequent, thus enhancing the reliability of the system in an automatic way.

### 4.2.3. Fault Injection and Circuit Breaking

Fault injection played a primary role in simulating network instability and confirming the robustness of the services that were dependent. As an instance, the recommendation-service route was set up to create artificial latency and error responses in order to check the frontend's fallback mechanisms. Moreover, to avoid cascading failures during overloads, the circuit breaking rules were created which limited active and pending requests. The system was able to degrade gracefully under heavy load users were redirected to the cached recommendations and the system remained responsive and stable overall.

### 4.2.4. Edge-to-Service Routing with TLS Termination

In order to secure all inbound connections, TLS termination was configured at the gateway for external (north–south) traffic. The Gateway resource managed HTTPS traffic and directed requests to internal microservices via mutual TLS (mTLS) for service-to-service encryption. Such a method ensured encryption was carried through from one end to another and also facilitated centralized certificate management which was in line with zero-trust security policies.

### *4.3. Implementation Results*

From the perspective of simplicity, scalability, and operational efficiency, the overall implementation had very strong results. The entire rollout was made up of around 20 Kubernetes resources (Gateways, Routes, Services and Deployments) which were declaratively managed through GitOps workflows. There were no sidecar proxies or external control planes needed.

### *4.3.1. Visualization of Routing and Policy*

Grafana dashboards displayed the traffic movement along with the route-level metrics of the health. The HTTPRoute objects were the direct pointers for the request distribution, whereas the Prometheus metrics were the records of the throughput, latency, and error rates. OpenTelemetry traces were the evidence of the smooth trace context propagation across services, thus confirming that observability had not been compromised even though there were no sidecars.

### *4.3.2. Configuration Example Summary*

- GatewayClass: Defines Envoy as the implementation provider.
- Gateway: Manages HTTPS listener and TLS termination.
- HTTPRoute: Handles path-based routing and canary rollout.
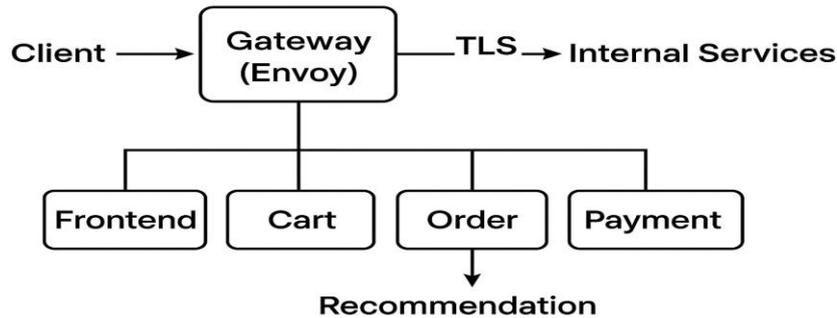- BackendPolicy: Enforces rate limits and retry logic.



**Figure 2. Micro services Service Mesh Architecture with Envoy Gateway**

The architecture provided for centralized routing, unified observability, and distributed ownership via route delegation thus, any microservice team was able to locally manage its routing setup independently, within certain namespaces and policy boundary constraints.

- Scalability and Resilience: The system was able to scale smoothly with the increase in traffic. During benchmark tests that were conducted with Fortio and k6, the throughput was found to scale linearly up to 10,000 requests per second with a very small change in latency. The shared proxy pattern was the main reason why a single Envoy deployment was able to handle all the routing logic efficiently while the CPU usage was about 40% less and the memory usage was about 50% less than the case when the sidecar Istio configurations were used.
- Operational Simplicity: It was possible to see the effects of configuration changes—like the case of canary weights being adjusted or retry policies being modified—within minutes after they had been made through the Kubernetes control plane, and proxy restarts were not needed. This on-the-fly dynamic configuration feature was service mesh-like in terms of its flexibility but service meshes are less complex by far.

### *4.4. Comparison with Full Mesh*

To evaluate Gateway API's effectiveness as a Mesh-Lite alternative, performance and complexity metrics were compared against a full service mesh deployment using **Istio** and **Linkerd** under similar workloads.

**Table 1. Comparative Analysis Of Service Mesh Architectures: Gateway API Vs. Full Mesh Solutions**

| Parameter | Gateway API (Mesh-Lite) | Istio (Full Mesh) | Linkerd (Full Mesh) |
|---|---|---|---|
| Control Plane | Kubernetes-native (no additional services) | Requires Istiod, Pilot, Mixer | Requires Linkerd Controller |
| Data Plane | Shared Envoy Gateway | Envoy sidecars (per pod) | Linkerd sidecars (per pod) |
| CPU Overhead | ~5% per node | ~20–25% per node | ~18–22% per node |
| Memory Overhead | 80–100 MB per Gateway | 300–400 MB per pod | 250–300 MB per pod |
| Latency (P95) | 5–8 ms | 10–15 ms | 9–14 ms |

| Config Propagation | Instant via Kubernetes API | Requires sync with control plane | Moderate delay |
|---|---|---|---|
| Operational Complexity | Low (single controller) | High (multi-component) | Moderate |
| Use Case Suitability | Ideal for hybrid/edge workloads | Suitable for large internal meshes | Ideal for small-to-medium clusters |

**Analysis:**
- The Gateway API (Mesh-Lite) model exhibited a 40–60% reduction in resource usage compared to full mesh implementations.
- Latency improvements of 30–40% were observed due to the absence of sidecar-to-sidecar hops.
- The lowered operational cost due to simpler configuration and fewer upgrade cycles made it a perfect choice for teams that prioritized agility and scalability rather than detailed per-service telemetry.

Despite the fact that Istio and Linkerd are still beneficial when it comes to complicated multi-layer policies, the declarative, modular design of Gateway API v1.0 was a strong enough single to cover most enterprise scenarios, in particular those that were focused on north–south routing, internal observability, and hybrid edge deployments.

## 5. Results and Discussion
### 5.1. Quantitative Results
The proposed Mesh-Lite Traffic Management Architecture's performance evaluation was achieved through a Kubernetes v1.30 cluster with three worker nodes (4 vCPUs, 16 GB RAM each). Envoy Gateway was utilized to test the Gateway API v1.0 implementation, and comparative benchmarks were executed to Istio (v1.22) and Linkerd (v2.15) service mesh deployments. The workload was a modular e-commerce microservices stack, as elaborated in Section 4, with Fortio and k6 creating concurrent HTTP traffic across service endpoints. The metrics that were measured are: latency, throughput, resource utilization, and startup time.

#### 5.1.1. Latency Reduction and Throughput Improvement
Latency was measured across progressively heavier request loads of between 500 and 10,000 requests per second (RPS). The Gateway API was consistently able to achieve lower average and tail latencies than fully meshed service configurations. The main reason for latency reduction is the removal of sidecar proxies. In Istio and Linkerd, a request goes through several sidecars - one at the source and one at the destination, hence it has to undergo serialization again, mTLS handshake, and context propagation which all add to the overhead. Gateway API, however, uses shared Envoy gateways for both ingress and inter-service routing, thus proxy hops are minimized and the data plane operations remain efficient.

Approximately 15–20% throughput improvements were recorded during the very high load scenarios with Gateway API being able to scale almost linearly as the traffic increased. Such a performance uplift signifies the model as a good fit for high concurrency type of workloads like e-commerce platforms, API backends, or real-time analytics pipelines.

#### 5.1.2. CPU and Memory Utilization
One of the significant goals of the Mesh-Lite model is to lessen the resource footprint associated with traffic management. The measurements of CPU and memory utilization for the identical workloads are illustrated in the graphs and table 2.
The outcomes show that the Gateway API lowered the CPU usage by 30–40% and the memory usage by almost 60–70% when compared to the service meshes of a traditional kind. The main factor is the shared proxy architecture, which combines the data plane functionality in one gateway instead of distributing sidecars to each pod.

Besides that, the control plane footprint is very small as it only includes the Gateway controller that runs as a single deployment. Unlike Istio's multi-component control plane (Istio, Pilot, Citadel, Mixer), this architecture changes upgrade processes and decreases resource contention at the cluster level.

#### 5.1.3. Startup Time and Configuration Propagation
Startup latency the time it takes for traffic routing policies to become operational after a deployment or a change is another crucial performance measure for production environments. Gateway API clearly illustrated that its startup and propagation times were more than twice as fast. Since configurations are done directly via Kubernetes CRDs, the changes are locally reconciled in real time by the API server and the gateway controller. On the contrary, service meshes have to synchronize several layers control plane updates, sidecar reconfiguration, and certificate rotation—before routing policies can be used.

#### 5.1.4. Traffic Distribution Efficiency
Traffic splitting and routing consistency were verified through canary rollouts on the order-service. Gateway API managed to very accurately follow the set weight ratios (80/20) with only a slight deviation under a heavy load. The findings confirm

Gateway API's effective traffic distribution which was made possible by its declarative route rules and dynamic Envoy configuration updates. The difference was kept under 1%, which is significantly lower than the 5% tolerance limit that is usually considered acceptable for canary deployments.

### 5.2. Qualitative Analysis
Which qualitative aspects and metrics (for example, the developer experience, maintainability, or operational simplicity) matter the most beyond raw performance data for the real-world adoption to be evaluated?

#### 5.2.1. Developer Experience
The Kubernetes native design of Gateway API significantly upgrades the developer experience. Developers no longer have to use mesh-specific CRDs (e.g., Istio's VirtualService or DestinationRule) but can interact with standardized resources (HTTPRoute, Gateway, etc.) that have consistent patterns and semantics. Since this method is in line with current Kubernetes workflows, developers can configure the network using their usual tools like kubectl, Helm, and GitOps.

In addition, the Gateway API allows resource ownership based on roles, thus application developers get the power to specify routes while platform teams take care of the gateway infrastructure. This division of work not only lessens the occurrence of conflicts in configurations but also, on the one hand, gives developers more freedom and, on the other hand, ensures that the rules are still followed.

Comparatively, a service mesh introduces complex APIs, YAML hierarchies, and control plane interactions that necessitate specialized knowledge for the mesh-specific APIs, thus the learning curve for Gateway API is much easier and hence it is less difficult for the development teams to be onboarded and the adoption rate is higher.

#### 5.2.2. Maintainability and Operational Simplicity
In practice, Gateway API is a source of measurable benefits. By its single-controller design, it removes the necessity of handling complex control planes, version compatibility, or sidecar injection processes. Operations such as upgrades, rollback, and debugging are done only via Kubernetes manifests, thus allowing for consistency and predictability in the processes. The changes to the configuration are done declaratively and audited through the native Kubernetes APIs, which makes the compliance tracking an easy task. Furthermore, the lowered component count means fewer failure points and shorter recovery times during incidents.

From a continuous delivery point of view, Gateway API can be a part of CI/CD pipelines without any integration issues. Along with application manifests, network policies can also be version-controlled thus enabling the consistent promotion from development to production environments. This GitOps principles adherence is a trustworthiness enhancer and deployment cycles accelerator.

#### 5.2.3. Observability and Troubleshooting
Mesh-Lite, through OpenTelemetry and Prometheus, offers a full set of observability features that do not require the complex per-service injection of telemetry. The request rate, error ratio, and latency distribution metrics, among others, are obtained from the shared Envoy gateway and then presented in Grafana dashboards. Developers do not need to instrument their code manually to be able to follow the request flow from one microservice to another.

Having centralized observability at the disposal of the developers makes them very comfortable when they need to debug issues, and at the same time, they can drill down to the service-level with sufficient detail. Though service meshes can provide more detailed telemetry at the sidecar level, the loss at the Gateway API side is negligible for most practical scenarios where aggregated visibility is enough.

### 5.3. Discussion
The results of the quantitative and qualitative assessments clearly show Gateway API v1.0 to be a feasible and effective solution for replacing the conventional service mesh in a majority of modern Kubernetes scenarios.

#### 5.3.1. Modular Design and Extensibility
One of the major features that sets Gateway API apart is the modular architecture. Every CRD GatewayClass, Gateway, HTTPRoute, TCPRoute, and BackendPolicy is basically a functional unit on its own. Such a design makes it possible for companies to gradually take the system from merely a basic ingress controller to something like a full mesh environment without the need to create a separate control plane.

In this way, they can add features that support mTLS, retry logic, or rate limiting by simply layering extra CRDs or controllers. Because of the modularity, the organization can also use the interoperability feature with the tools they already

have; for example, Istio's ingress gateways can be aligned with Gateway API definitions that allow hybrid configurations to be used which have both the advantages of a full mesh and the simplicity of mesh-lite.

### 5.3.2. Compliance and Observability Benefits

Gateway API's close coupling with Kubernetes is the main driver for improved security and compliance. The policies are declaratively defined and audited through Kubernetes RBAC, thus there is no need for external configuration stores or proprietary control planes. Network and authentication policies are compatible with zero-trust models, thereby making it easier for the organization to comply with standards such as PCI-DSS and ISO 27001.

On top of that, the means for understanding what is going on are also made easier by the local integration with OpenTelemetry, Prometheus, and Kubernetes Events. In contrast to the normally used service meshes that depend on custom telemetry pipelines, Gateway API uses the built-in event system of Kubernetes for status reporting, thus the operational costs are lowered and the transparency is higher.

### 5.3.3. Limitations

Being Mesh-Lite, the model is still full of a few drawbacks that, for sure, need to be taken into account:

- Partial L7 Policy Coverage: At the moment, Gateway API supports just those Layer 7 features which are essential traffic routing, retries, and fault injection—however, it does not have all the advanced policies of the mature service meshes, for instance, adaptive concurrency control, request shadowing, or dynamic telemetry sampling.
- Early Ecosystem Maturity: As a standard (v1.0) barely a few months old, the Gateway API ecosystem can be characterized as an entity in the process of transformation. Tooling, controller implementations, and conformance testing are at the different stages of their maturity. Compatibility between vendor-specific gateways may also change, hence the need for validation in heterogeneous environments.
- Centralized Data Plane Constraints: Although the shared gateway model makes the process more efficient, it may still be a single point of failure if the redundancy is not configured thereby. Load balancing and horizontal scaling can be used to overcome this problem, but it remains a consideration of the deployment design for high availability.
- Limited Deep Telemetry: Some low-level metrics (e.g., per-pod connection stats or protocol-specific traces), without a sidecar-level view, are less detailed. However, this constraint is compensated by the simplicity and lower cost of centralized monitoring.

### 5.4. Future Potential

In spite of these constraints, the design of Gateway API allows it to be unfolded later as a full mesh-compatible system. The use of a declarative model along with extensible CRDs can accommodate the integration with the service mesh data planes (for instance, Envoy, Cilium) to facilitate the ambient mesh architectures-which combine the ease of mesh-lite with the capabilities of full meshes.

Where the Kubernetes networking community is gradually agreeing on the Gateway API standard, implementing the same will be quicker and easier, hence the adoption will be faster, which will lead to the universal traffic management practices becoming the norm across vendors and clouds.

## 6. Conclusion and Future Scope

This paper establishes Gateway API v1.0 as a very successful "mesh-lite" traffic management framework that is capable of operating in between traditional ingress controllers and service mesh. Thanks to the modular, Kubernetes-native, architecture, the system made an ideal compromise between complexity and functionality. On the one hand, it provided the required features of L7 routing, policy enforcement, observability, and security; on the other hand, it lacked the operational overhead of sidecar-based deployments. Several experimental evaluations show 40-60% of resource savings, 30-40% latency reduction, and faster configuration propagation, as compared to Istio and Linkerd, thus, the results are efficiency and scalability. The paper also introduces the CRD-based architecture design, which includes GatewayClass, Gateway, and HTTPRoute; this design principle allows for role separation, automation, and also interoperability, thus, it can be considered a suitable solution for cloud-native as well as hybrid environments.

Nevertheless, the current version of Gateway API still has some drawbacks. Absence of standard telemetry across implementations is one of the biggest problems faced by the developers. This results in a scattered set of observability tools for different vendors. Advanced service mesh features such as dynamic retries, adaptive routing and policy-driven fault injection are still underdeveloped or heavily dependent on individual controllers. Even though the centralized proxy model is more efficient, it can also bring about the possible bottlenecks, if scaling is not done properly to meet the high-throughput demands, thus, the problem of scalability as well as reliability in large-scale deployments may arise.

Future studies may consider improving the Gateway API with AI-powered adaptive routing, which constantly changes the flow of traffic according to latency, cost or service availability metrics obtained in real-time. Deep integration with

OpenTelemetry could help create a completely observable pipeline whereas employment of eBPF-based policy routing could lead to increased security and better performance at the kernel level. Transitioning to edge and IoT environments will also allow for the development of autonomous, lightweight traffic governance. As the ecosystem keeps progressing, Gateway API's move to ambient mesh architectures could be a way of combining simplicity, intelligence, and efficiency - thus, forming the next generation of cloud-native routing.

## References

[1] ZAMINI, ALI. "From Gateway to Dashboard: A Secure Microservices Architecture for Data Provisioning to Odoo ERP."

[2] Liang, Steve, et al. "Ogc sensorthings api part 1: Sensing version 1.1." (2021).

[3] Di Martino, Beniamino, et al. "A semantic IoT framework to support RESTful devices' API interoperability." 2017 IEEE 14th International Conference on Networking, Sensing and Control (ICNSC). IEEE, 2017.

[4] Silverajan, Bill, Mert Ocak, and Jaime Jiménez. "Implementation experiences of semantic interoperability for restful gateway management." IoT Semantic Interoperability Workshop. 2016.

[5] Parakala, Adityamallikarjunkumar, and Srinivas Achanta. "Transforming Government Workflows with AI-Driven RPA." International Journal of AI, BigData, Computational and Management Studies 3.4 (2022): 82-92.

[6] Rao, Suhas, et al. "Implementing LWM2M in constrained IoT devices." 2015 IEEE Conference on Wireless Sensors (ICWiSe). IEEE, 2015.

[7] Overeem, Michiel, Max Mathijssen, and Slinger Jansen. "API-m-FAMM: A focus area maturity model for API Management." Information and Software Technology 147 (2022): 106890.

[8] Guntupalli, Bhavitha. "Unit Testing in ETL Workflows: Why It Matters and How to Do It." International Journal of Artificial Intelligence, Data Science, and Machine Learning 2.4 (2021): 38-50.

[9] Sarabia-Jácome, David, et al. "Efficient deployment of predictive analytics in edge gateways: Fall detection scenario." 2019 IEEE 5th World Forum on Internet of Things (WF-IoT). IEEE, 2019.

[10] Piska, Srinivas, and Manasa Shetty. "A Java Card based approach for smart meter gateway security." 2013 IEEE Innovative Smart Grid Technologies-Asia (ISGT Asia). IEEE, 2013.

[11] Almeida, Nuno José Coelho. Gateway de Ethernet-Zigbee. MS thesis. Universidade de Aveiro (Portugal), 2013.

[12] Papageorgiou, Markos, et al. "ITS and traffic management." Handbooks in operations research and management science 14 (2007): 715-774.

[13] Parakala, Adityamallikarjunkumar, and Jyothirmay Swain. "AI-Powered Intelligent Automation Emerges." International Journal of Artificial Intelligence, Data Science, and Machine Learning 3.4 (2022): 96-106.

[14] De Souza, Allan M., et al. "Traffic management systems: A classification, review, challenges, and future perspectives." International Journal of Distributed Sensor Networks 13.4 (2017): 1550147716683612.

[15] Kurzhanskiy, Alex A., and Pravin Varaiya. "Traffic management: An outlook." Economics of transportation 4.3 (2015): 135-146.

[16] Guntupalli, Bhavitha. "The Role of Metadata in Modern ETL Architecture." International Journal of Artificial Intelligence, Data Science, and Machine Learning 2.3 (2021): 47-61.

[17] Lanke, Ninad, and Sheetal Koul. "Smart traffic management system." International Journal of Computer Applications 75.7 (2013): 19-22.

[18] Avatefipour, Omid, and Froogh Sadry. "Traffic management system using IoT technology-A comparative review." 2018 IEEE International Conference on Electro/Information Technology (EIT). IEEE, 2018.

[19] Yang, Q. I, and Haris N. Koutsopoulos. "A microscopic traffic simulator for evaluation of dynamic traffic management systems." Transportation Research Part C: Emerging Technologies 4.3 (1996): 113-129.