# Optimizing SASE for Low Latency and High Bandwidth Applications: Techniques for Enhancing Latency-Sensitive Systems

Arunkumar Velayutham   Cloud Software Development Engineer
and Technical Lead at Intel, Arizona, USA

## Abstract

The increasing reliance on digital transformation and the demand for real-time processing in applications like video conferencing, real-time data analytics, and Internet of Things (IoT) networks have highlighted the need for optimized networking and security frameworks. Secure Access Service Edge (SASE) has emerged as a prominent solution to address the growing requirements of secure cloud and edge computing infrastructures. However, its effectiveness in environments with stringent latency and bandwidth requirements has been questioned. This paper examines the techniques for optimizing SASE to handle latency-sensitive applications without compromising security. The paper begins by providing a detailed overview of SASE's architecture and its core components, including SD-WAN, security services, and the cloud-based framework. It then addresses the challenges posed by latency-sensitive and high-bandwidth applications, with a special focus on video conferencing, real-time data analytics, and IoT environments. Several optimization strategies are discussed, including edge computing integration, packet optimization techniques, traffic prioritization mechanisms, and the use of artificial intelligence for dynamic path selection. The exploration also includes advancements in SD-WAN technologies such as multipath routing and dynamic bandwidth allocation. These techniques are useful for achieving seamless, low-latency performance while preserving the security features inherent to SASE. The paper argues for achieving a balance between security and performance in modern, distributed network architectures.

# 1   Introduction

The adoption of cloud-based services and the expansion of edge computing are re-shaping enterprise networks in a profound and multi-dimensional manner. In particular, Secure Access Service Edge (SASE) emerges as a significant paradigm by integrating wide area networking (WAN) capabilities with robust, multi-layered security functions. SASE presents a unified solution that addresses the need for both scalable networking and stringent security in increasingly decentralized enterprise environments. As organizations migrate towards distributed architectures—largely driven by the proliferation of remote work, mobile connectivity, and the Internet of Things (IoT)—the ability to optimize network performance while ensuring strong security becomes imperative. This need is exacerbated in the case of latency-sensitive applications such as video conferencing, real-time data analytics, and IoT-based services, all of which demand low-latency and high-bandwidth connections to function effectively (Chandramouli and Chandramouli, 2022).
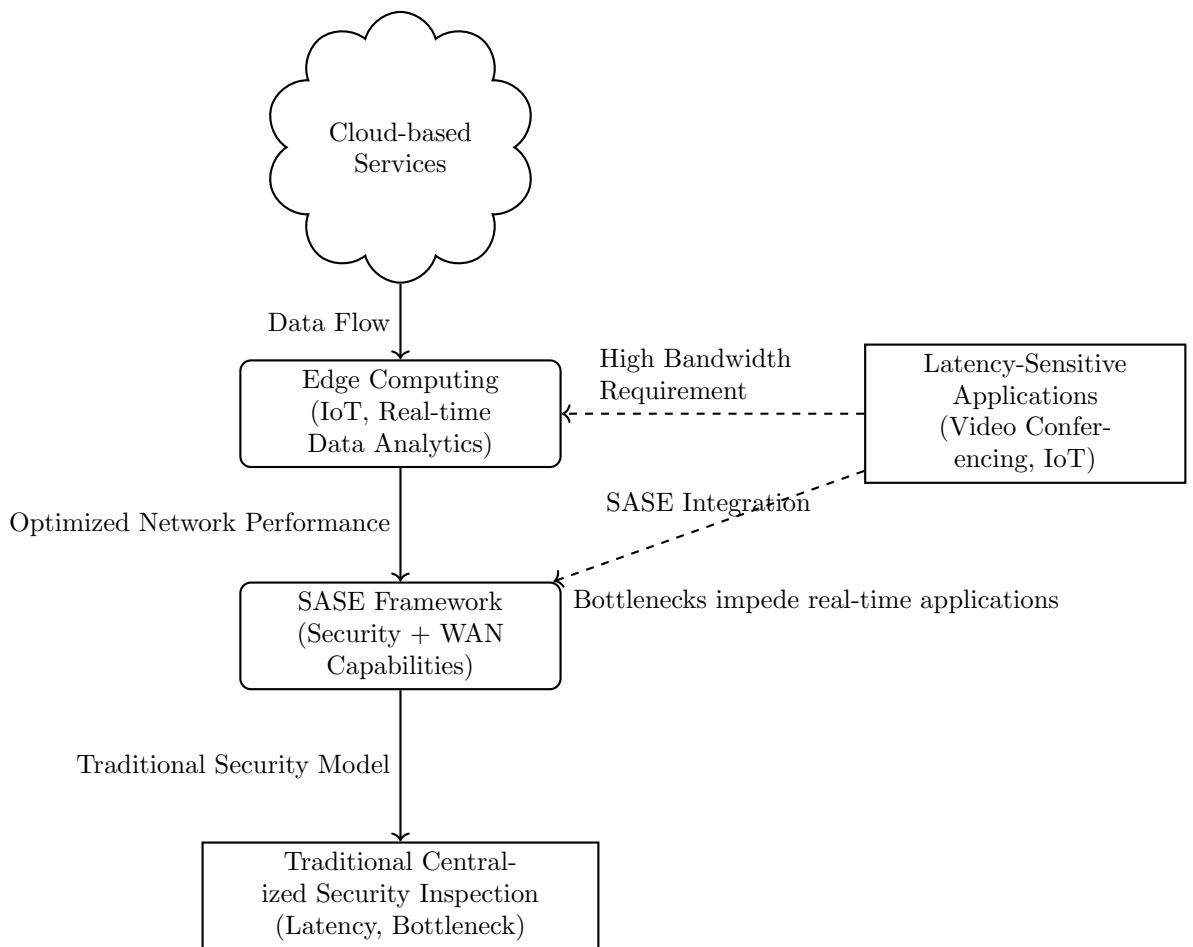


Figure 1: Optimizing SASE for Latency-Sensitive Applications in Cloud and Edge Computing Environments

Traditional security models those that rely on centralized inspection points, often introduce undesirable latency and bottlenecks into the network architecture. These centralized points of control and inspection can choke the data flow, especially when scaling to support a distributed workforce or large-scale IoT deployments. The inefficiencies in these legacy models are becoming increasingly apparent as enterprises strive to maintain performance benchmarks for real-time applications, which are sensitive to network latency and throughput limitations. Consequently, there is an urgent need to rethink and optimize SASE solutions, focusing on minimizing latency and maximizing bandwidth utilization without compromising the security integrity of the network. This delicate balance between performance and security poses several technical challenges, but also opens

up opportunities for innovative solutions.

Edge computing plays a critical role in this optimization process, as it enables data processing closer to the point of data generation. By moving compute and data analytics capabilities closer to the network edge, edge computing reduces the distance that data must travel, thereby minimizing latency. In the context of SASE, edge computing allows security functions—such as data encryption, threat detection, and traffic inspection—to occur closer to the user or device, reducing the need to backhaul traffic to a central data center for inspection. This decentralized approach to security inspection is beneficial for latency-sensitive applications, where even minor delays can lead to degraded performance. Moreover, edge computing supports a more efficient distribution of network load, alleviating potential congestion at central hubs by distributing traffic across multiple edge locations.

Packet optimization is another pivotal technique for enhancing SASE performance in the context of latency-sensitive and high-bandwidth applications. Packet optimization techniques, such as data compression, deduplication, and loss recovery, are instrumental in reducing the amount of data that must be transmitted across the network. By reducing packet overhead, these techniques can improve bandwidth efficiency and reduce latency. Compression algorithms, for instance, can significantly decrease the volume of data sent across the network, allowing for faster transmission and lower latency. Deduplication further enhances efficiency by ensuring that redundant data is not sent repeatedly across the network, thus preserving bandwidth for more critical data transmissions. Additionally, loss recovery techniques—such as forward error correction (FEC)—help to mitigate the impact of packet loss on real-time applications, which are sensitive to such disruptions.

Traffic prioritization in environments where multiple types of data flows coexist, also plays a crucial role in optimizing SASE for real-time performance. Not all traffic is created equal: real-time applications such as video conferencing or real-time analytics require near-instantaneous data transmission, whereas other traffic types, such as email or file transfers, can tolerate some delay. SASE architectures can leverage traffic prioritization techniques to ensure that latency-sensitive applications receive preferential treatment over less time-critical applications. Quality of Service (QoS) policies, for example, can be configured to prioritize traffic based on its sensitivity to delay, jitter, and packet loss. By prioritizing certain types of traffic, organizations can ensure that their most critical applications receive the bandwidth and low-latency connections they require to function effectively.

The emergence of Software-Defined Wide Area Networking (SD-WAN) further enhances the capabilities of SASE by providing dynamic, application-aware routing of traffic across multiple WAN links. SD-WAN allows organizations to intelligently route traffic based on real-time network conditions, such as bandwidth availability, latency, and packet loss. In the context of SASE, SD-WAN can be used to dynamically adjust network paths to ensure optimal performance for latency-sensitive applications. For example, if a particular WAN link experiences high latency or congestion, SD-WAN can automatically reroute traffic to a less congested link, thereby maintaining low-latency performance. Additionally, SD-WAN's ability to aggregate multiple WAN connections provides a level of redundancy and failover that further enhances network resilience and performance.

The integration of artificial intelligence (AI) and machine learning (ML) into SASE architectures offers yet another avenue for optimizing performance and security. AI and ML algorithms can be leveraged to improve routing decisions, enhance traffic management, and detect security threats in real-time. In terms of routing, AI and ML can analyze historical and real-time network data to predict traffic patterns and preemptively adjust routing decisions to minimize latency and congestion. These algorithms can also be used to identify anomalies in network traffic that may indicate a security threat, allowing for faster detection and response. In addition, AI and ML can be employed to optimize the configuration of security policies, ensuring that they are both effective and minimally intrusive to network performance. For example, AI-driven analysis can help determine which types of traffic are most likely to pose a security risk, enabling more targeted inspection and reducing the overall burden on network resources.

Despite these advancements, there are still several challenges associated with

optimizing SASE for latency-sensitive and high-bandwidth applications. One of the primary challenges lies in balancing the competing demands of security and performance. As security policies become more complex and data volumes continue to grow, the overhead associated with security inspection and enforcement can become a significant performance bottleneck (Jani, 2021). To mitigate this issue, SASE solutions must employ intelligent traffic inspection mechanisms that can distinguish between traffic that requires deep inspection and traffic that can be passed through with minimal scrutiny. Furthermore, as encryption becomes more pervasive across enterprise networks, decrypting and inspecting encrypted traffic introduces additional latency. SASE solutions must therefore include efficient encryption and decryption processes, or adopt selective decryption techniques that minimize the performance impact of inspecting encrypted traffic.

To illustrate the importance of balancing performance and security in SASE architectures, we provide a comparative analysis of latency and bandwidth consumption across different optimization techniques. The following table highlights the performance impact of packet optimization, traffic prioritization, and SD-WAN in a typical SASE deployment.

Table 1: Comparison of Performance Optimization Techniques in SASE

| Optimization Technique | Latency Reduction | Bandwidth Efficiency | Security Impact |
|:---:|:---:|:---:|:---:|
| Packet Optimization | Moderate | High | Low |
| Traffic Prioritization | High | Moderate | Low |
| SD-WAN | High | High | Moderate |
| Edge Computing | Very High | High | Low |

As the table demonstrates, different optimization techniques vary in their effectiveness at reducing latency and improving bandwidth efficiency. Edge computing, for example, provides significant latency reduction by processing data closer to the user or device, while packet optimization focuses on improving bandwidth efficiency with minimal impact on security. Traffic prioritization and SD-WAN offer high levels of performance optimization by dynamically managing network traffic based on real-time conditions, though each comes with its own set of trade-offs in terms of security and implementation complexity.

Looking forward, the continued evolution of SASE will likely involve greater integration of AI and ML technologies to further enhance performance and security capabilities. As AI and ML algorithms become more sophisticated, they will be able to provide even more granular insights into network traffic patterns and security threats, allowing for highly optimized routing and security inspection processes. Additionally, the use of blockchain technology in SASE architectures has been proposed as a means of enhancing trust and security in decentralized networks in environments where multiple organizations or entities share a common infrastructure. While still in the early stages of exploration, blockchain has the potential to provide an additional layer of security by ensuring the integrity and authenticity of network transactions.

Table 2: Future Research Directions in SASE Optimization

| Research Area | Potential Benefits | Challenges |
|:---:|:---:|:---:|
| AI and ML Integration | Improved traffic management | Complexity of implementation |
| Blockchain in SASE | Enhanced security and trust | Scalability and overhead |
| 5G Integration | Ultra-low latency for real-time applications | Infrastructure costs |
| Selective Encryption | Reduced latency in encrypted traffic | Security trade-offs |

## 2   Background on SASE Architecture

Secure Access Service Edge (SASE) is a network architecture designed to address the challenges of providing secure, scalable, and flexible network access in increasingly distributed enterprise environments. It combines the functionality of Software-Defined Wide Area Networking (SD-WAN) with integrated security
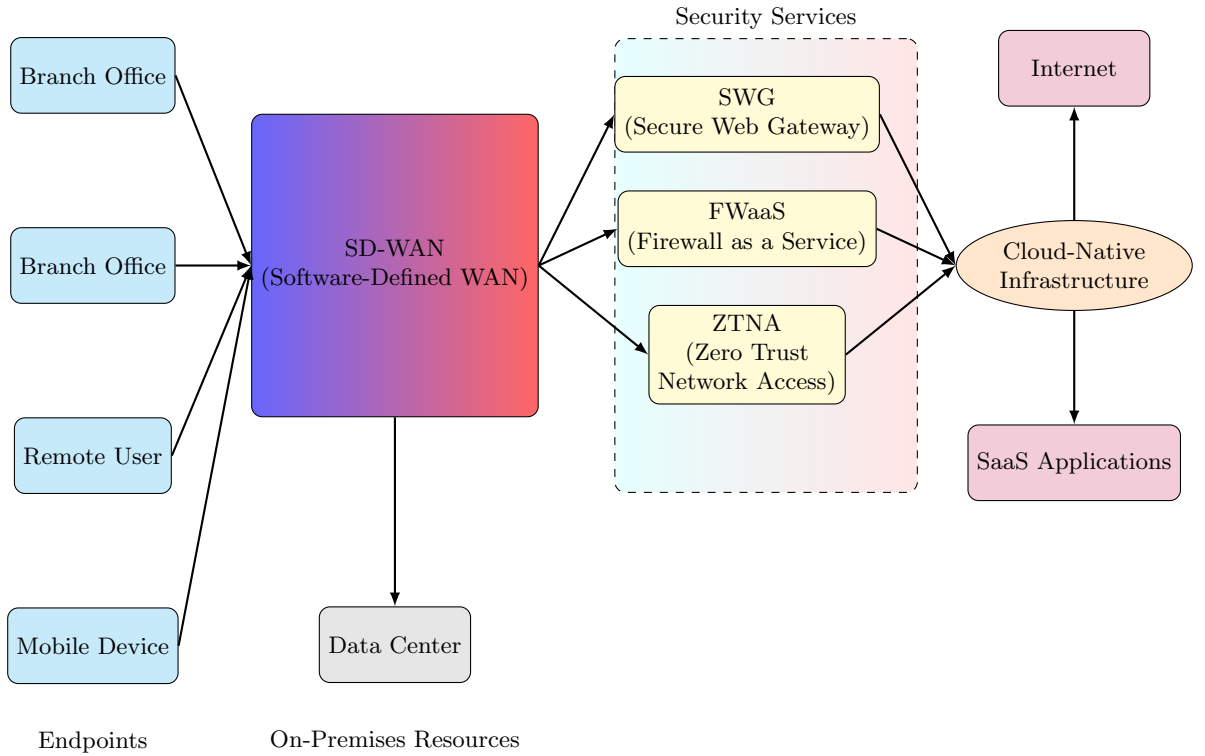
Figure 2:   SASE Architecture Diagram

services, including Secure Web Gateway (SWG), Firewall-as-a-Service (FWaaS), and Zero Trust Network Access (ZTNA). This integration allows for a more efficient and cohesive approach to managing both network performance and security. The shift towards decentralized architectures, driven by the rise of cloud services, remote work, and IoT deployments, requires solutions that can adapt to a broader range of access points without compromising security or network efficiency. SASE offers a model that supports these requirements by providing cloud-native, scalable network management and security enforcement.

One of the primary components of SASE is SD-WAN, which plays a crucial role in managing network traffic efficiently across various wide-area network paths. SD-WAN offers the capability to optimize traffic routing based on real-time network conditions, including metrics like latency, jitter, and packet loss. This allows enterprises to maintain performance levels for critical applications by dynamically adjusting how traffic is routed. For example, SD-WAN can prioritize traffic for latency-sensitive applications such as video conferencing or real-time analytics, ensuring that these services receive the bandwidth and low-latency connections they require to function effectively. The flexibility of SD-WAN in managing multiple network links, including broadband, MPLS, and LTE, also provides redundancy and failover options, improving network resilience and minimizing disruptions.

SASE's integration of security services directly into the network infrastructure distinguishes it from more traditional, centralized security models. In particular, Secure Web Gateway (SWG) is responsible for filtering and inspecting web traffic to prevent exposure to threats like malware or phishing attacks. SWG allows for real-time security inspection at the point of access, reducing the need for traffic to be backhauled to a central data center for security enforcement, which can introduce latency. This approach to distributed security is especially relevant for enterprises with users or devices spread across different geographic locations. Firewall-as-a-Service (FWaaS) extends this security functionality further by offering cloud-delivered firewall capabilities, enabling enterprises to apply consistent security policies across both on-premises and cloud environments. By integrating FWaaS into the SASE architecture, organizations can enforce security measures at the edge of the network, ensuring that traffic is inspected and policies are ap-

plied as close to the user or device as possible. This eliminates the performance penalties that would arise from centralized traffic inspection for high-bandwidth or geographically dispersed users.

Zero Trust Network Access (ZTNA) is another critical element within SASE, designed to enforce a strict access control model that does not automatically trust any user or device, even those within the organization's network. ZTNA operates on a principle of least privilege, granting users access only to the specific resources they require based on their identity and context, such as location or device type. This model enhances security by ensuring that access is continuously authenticated and authorized, reducing the risk of unauthorized access to sensitive resources. In the context of SASE, ZTNA's policies are applied at the edge, allowing for granular control over user access regardless of location, whether remote or on-site.

The cloud-native nature of SASE supports the scalability and performance required for today's distributed enterprise networks. By decentralizing both network management and security enforcement, SASE reduces the dependency on central data centers, which have traditionally been a source of latency in global networks. The cloud-native architecture allows for the distribution of network and security functions across a global infrastructure, reducing the distance that data must travel for security inspections and routing decisions. This distributed approach not only supports scalability but also improves the responsiveness of the network for real-time applications that are sensitive to delays. As the demand for high-performance connectivity grows with the rise of cloud-based services and remote work, the ability of SASE to dynamically scale and adapt to varying network loads becomes a crucial advantage.

While traditional network models rely heavily on centralized security appliances and data centers, SASE's architecture enables enterprises to push security enforcement closer to the user or edge device. This shift is beneficial for organizations that need to support geographically dispersed users, as it minimizes the latency introduced by routing traffic through central inspection points. By processing security policies and routing decisions at the edge, SASE reduces the overhead associated with backhauling traffic to centralized locations, leading to more efficient and timely data flow. This capability is essential for applications such as IoT deployments and real-time analytics, where any delay can impact performance and functionality.

## 3   Challenges of Latency-Sensitive Applications

Latency-sensitive applications, such as video conferencing, real-time data analytics, and IoT networks, present specific challenges that strain both network performance and security frameworks (Islam et al., 2021). These applications demand high throughput, real-time data transmission, and robust security, all of which must be balanced to avoid performance degradation. As these applications become more prevalent across industries, optimizing network architectures to accommodate their unique demands is critical. Secure Access Service Edge (SASE) has emerged as a framework to address these needs by integrating networking and security into a unified architecture. However, several challenges arise in deploying SASE for latency-sensitive applications, especially when low-latency and high-bandwidth communication is paramount (van der Walt and Venter, 2022).

Video conferencing serves as a prime example of a latency-sensitive application that requires both low latency and high reliability. Real-time communication tools video conferencing platforms, are highly sensitive to fluctuations in network performance. Delays of even a few milliseconds can lead to significant issues such as frame drops, video stuttering, audio desynchronization, and frequent interruptions in communication. These disruptions can be damaging in professional environments where seamless communication is critical. In a SASE architecture, network traffic associated with video conferencing must be securely routed and monitored without impacting performance. Security services like encryption, deep packet inspection (DPI), and real-time threat detection are essential to maintaining the integrity of communication (Yiliyaer and Kim, 2022). However, the real

challenge is ensuring that these security measures do not introduce additional latency or degrade the user experience. Encryption, for instance, while necessary for maintaining privacy, can increase processing overhead, potentially affecting the smooth flow of audio and video streams.

Real-time data analytics is another critical application area in industries such as finance, healthcare, and industrial automation, where timely insights are necessary for informed decision-making. In these environments, data is generated and consumed in real-time, requiring low-latency processing to extract actionable insights. Any delays introduced by the network can disrupt workflows, leading to suboptimal outcomes. In financial services, for example, milliseconds can mean the difference between a successful trade and a lost opportunity. Similarly, in healthcare, real-time analytics may be used for patient monitoring, where network-induced delays could have serious consequences for patient outcomes. SASE environments must apply security policies, such as encryption and data inspection, to protect sensitive data while avoiding performance bottlenecks. The challenge lies in performing deep security inspections at the edge or in the cloud, where traffic can be inspected without impacting the real-time processing requirements of the application. This is especially difficult when large data volumes are processed continuously, as in real-time analytics, where the need for security must be weighed against the need for speed (Chen et al., 2022).

IoT networks represent a unique challenge, as they often involve large-scale deployments of devices that generate continuous streams of data and require both high-bandwidth and low-latency communication channels. IoT systems span multiple sectors, including autonomous vehicles, smart grids, industrial IoT, and healthcare monitoring systems. The devices in these ecosystems often have minimal built-in security and may operate on low-power, low-compute platforms, making them vulnerable to cyber threats. For instance, autonomous vehicles require near-instantaneous communication with nearby sensors and control systems to ensure safe operation. Any delay in the transmission or processing of data could result in safety risks or system failures. In a SASE framework, ensuring secure communication for IoT devices is paramount, as these devices frequently lack native security capabilities. SASE's role is to provide end-to-end encryption, threat detection, and policy enforcement while maintaining the low-latency performance required by these systems. The need to balance security with performance in IoT networks is complex because of the sheer volume of data generated and the distributed nature of these devices, which operate in environments with diverse network conditions (Kaur, 2018).

The following table highlights the unique challenges presented by each of these latency-sensitive applications when integrated into a SASE architecture.

Table 3: Challenges of Latency-Sensitive Applications in a SASE Environment

| Application Type | Performance Requirements | Security Requirements | Challenges in SASE |
|---|---|---|---|
| Video Conferencing | Low latency, high reliability | Encryption, DPI, threat detection | Ensuring security without degrading audio/video quality |
| Real-Time Data Analytics | Near-instant data processing | Data inspection, encryption | Balancing real-time data processing with security overhead |
| IoT Networks | Low-latency, high-bandwidth communication | End-to-end encryption, threat detection | Securing devices with minimal native security while maintaining performance |

The common thread across these applications is the tension between performance and security. For video conferencing, maintaining low-latency communication while ensuring privacy and security is a balancing act that requires careful consideration of network optimization techniques. Similarly, in real-time data analytics, the rapidity of data processing must not be compromised by the introduction of security mechanisms that could introduce delay. IoT networks, on the

other hand, face unique challenges due to the diversity and scale of devices, which often operate with minimal inherent security but require robust communication channels to function effectively (Gandhi et al., 2022).

In addition to these individual challenges, the integration of these applications into a SASE framework poses broader architectural issues. Since SASE is designed to deliver both network and security services through a cloud-native infrastructure, optimizing its components to handle such diverse, latency-sensitive applications becomes essential. SASE's reliance on edge computing, where network and security functions are distributed closer to the user or device, offers a potential solution. By processing data and enforcing security policies closer to the source, edge computing reduces the amount of data that needs to traverse the network, thereby lowering latency. However, ensuring that edge devices and network nodes are equipped with the necessary processing power to handle both security and performance demands is a key challenge.

To better understand the impact of latency and bandwidth requirements on SASE performance, the following table compares key metrics such as latency tolerance, data volume, and security overhead for video conferencing, real-time data analytics, and IoT networks.

Table 4: Comparison of Latency and Bandwidth Requirements in Latency-Sensitive Applications

| Application Type | Latency Tolerance | Data Volume | Security Overhead |
| --- | --- | --- | --- |
| Video Conferencing | Milliseconds | Moderate to high | Moderate |
| Real-Time Data Analytics | Milliseconds to seconds | Very high | High |
| IoT Networks | Sub-millisecond in critical cases | High | Low to high (depending on the use case) |

This table (4) shows the variability in latency and bandwidth demands across different applications, as well as the varying levels of security overhead that SASE architectures must accommodate. Video conferencing requires a low-latency connection with moderate to high data volumes, but security overhead must be carefully managed to avoid disruptions in the user experience. Real-time data analytics, on the other hand, involves the processing of very high volumes of data, often with stringent security requirements, but the tolerance for latency may be slightly higher than for video conferencing. IoT networks, especially in critical applications like autonomous vehicles, demand extremely low-latency communication with highly variable data volumes, presenting another layer of complexity in securing these devices without introducing unacceptable delays (Zhang, 2019).

## 4   Optimizing SASE for Latency-Sensitive and High-Bandwidth Applications

### 4.1   Edge Computing Integration

The integration of edge computing within a Secure Access Service Edge (SASE) architecture addresses the critical need for reducing latency for applications where real-time performance is essential. As enterprises increasingly adopt cloud-based services and operate in distributed environments, the ability to process data closer to its source becomes paramount. Edge computing achieves this by decentralizing data processing, shifting it from traditional centralized data centers to nodes located nearer to the data origin. This reduction in the physical distance that data must travel directly translates into lower latency, an important factor for latency-sensitive applications such as video conferencing, real-time data analytics, and IoT systems.

Edge computing is relevant in a SASE environment due to the convergence of networking and security functions in this architecture. Traditionally, security services—such as data encryption, traffic inspection, and threat detection—are
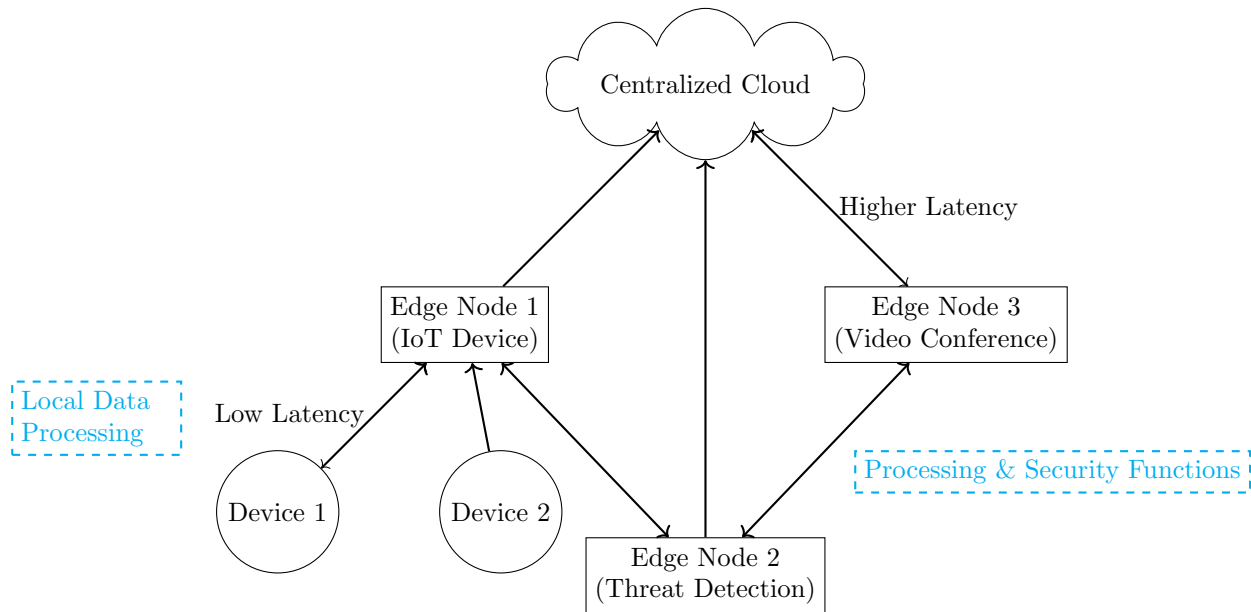
Figure 3: Local Processing and Security within Edge Nodes Minimizing Latency

applied at centralized points, often located at distant data centers or cloud infrastructure. This centralized model introduces delays, as data must traverse the network to undergo security processing before it is allowed to proceed. These delays can significantly affect performance for real-time applications, where even minor latency can degrade the user experience or the functionality of the system. Edge computing within SASE allows these security functions to be implemented locally, at the edge of the network, closer to where data is generated and consumed. This local processing reduces the round-trip time required for data to be secured, inspected, and transmitted, thus minimizing latency while still ensuring robust security.

In practical terms, the implementation of edge computing in SASE involves deploying edge nodes or devices that possess sufficient computational resources to perform security operations autonomously. These edge nodes can be positioned at strategic points, such as near branch offices, remote locations, or IoT hubs, allowing them to handle network traffic generated in their vicinity. This decentralization of security functions is essential in scenarios involving geographically dispersed users or devices, where routing traffic to a central data center would introduce unacceptable delays. For example, in an IoT deployment involving autonomous vehicles, edge nodes placed at roadside infrastructure can process critical vehicle-to-infrastructure (V2I) communication in real time, applying security measures such as traffic encryption and anomaly detection locally. This reduces the time taken for data transmission and processing, ensuring that the vehicles receive the information they need to operate safely and efficiently.

To understand how edge computing works within SASE, it is important to examine the specific components and technologies that enable this integration. The first major component is local data processing, where compute resources at the edge perform security and network functions. Edge devices such as gateways, routers, and specialized edge servers are typically used to execute security services like packet inspection, encryption, and threat detection. This processing capability is often powered by lightweight, containerized applications or virtualized network functions (VNFs), which allow edge nodes to handle specific security tasks without requiring the full capabilities of a traditional data center.

Another critical component of edge computing in SASE is distributed data storage. In many cases, security operations such as traffic analysis or data inspection require access to historical data or security policies, which are often stored in the cloud. However, in a distributed edge environment, some of this data can be cached or replicated locally at the edge nodes, reducing the need for constant

back-and-forth communication with the cloud. This ensures that security policies are applied in real time, even in environments with intermittent or limited connectivity. Moreover, by storing critical security data at the edge, organizations can minimize the impact of data breaches or attacks targeting central data repositories.

In terms of security, the ability to process data locally through edge computing also brings the advantage of localized threat detection and response. Threats can be detected and mitigated at the point of origin, rather than being sent across the network for centralized analysis. For instance, a branch office connected to an enterprise network via a SASE deployment might encounter a localized malware attack. In this case, the edge node deployed at the branch can detect the threat through deep packet inspection and neutralize it immediately, without introducing the delays that would occur if the traffic had to be inspected at a remote security gateway. This capability enhances the overall security posture of the enterprise by enabling faster responses to potential threats while keeping the network performance intact.

Edge computing also enables bandwidth optimization in a SASE architecture. By processing and filtering data at the edge, unnecessary traffic can be removed before it traverses the network. This is useful in IoT environments, where large volumes of raw data are often generated, much of which may not require transmission to the cloud. For example, in an industrial IoT setting, edge devices could analyze sensor data locally, filtering out irrelevant or redundant information and sending only actionable insights to a central processing hub. This reduces the overall bandwidth consumption, freeing up network resources for more critical tasks. Moreover, the ability to offload processing tasks to edge nodes ensures that centralized resources are not overwhelmed, further improving the scalability and efficiency of the network.

A key technological enabler for edge computing in SASE is containerization and microservices architectures. By leveraging containers, organizations can deploy specific security functions as microservices on edge nodes, ensuring that these functions are lightweight, modular, and scalable. Containers allow for the rapid deployment and updating of security services across distributed edge locations, providing flexibility in adapting to evolving network conditions or security requirements. Additionally, container orchestration platforms like Kubernetes can be used to manage the deployment and scaling of these microservices at the edge, ensuring that resources are allocated efficiently across the network.

To assess the impact of edge computing on latency and security performance, we can look at the following table, which outlines how different security functions are impacted when shifted to the edge in a SASE environment:

Table 5: Impact of Edge Computing on Security Functions in SASE

| Security Function | Traditional Centralized Processing Latency | Edge-Based Processing Latency | Impact on Security |
|---|---|---|---|
| Packet Inspection | Moderate to High | Low | Localized inspection reduces latency, maintains security |
| Encryption / Decryption | High | Low to Moderate | Reduces round-trip latency for encrypted traffic |
| Threat Detection | Moderate to High | Low | Faster threat detection and response at the edge |
| Traffic Filtering | Moderate | Low | Optimized bandwidth usage, real-time filtering |

As shown in the table 5, the shift from traditional centralized processing to edge-based processing significantly reduces latency across key security functions. For instance, packet inspection, which is resource-intensive in a centralized model, can be performed more efficiently at the edge, reducing delays for time-sensitive applications. Similarly, encryption and decryption processes benefit from local execution, as data does not need to be sent back to a central hub for security verification. This not only improves performance but also ensures that sensitive

data remains secure, even in distributed environments.

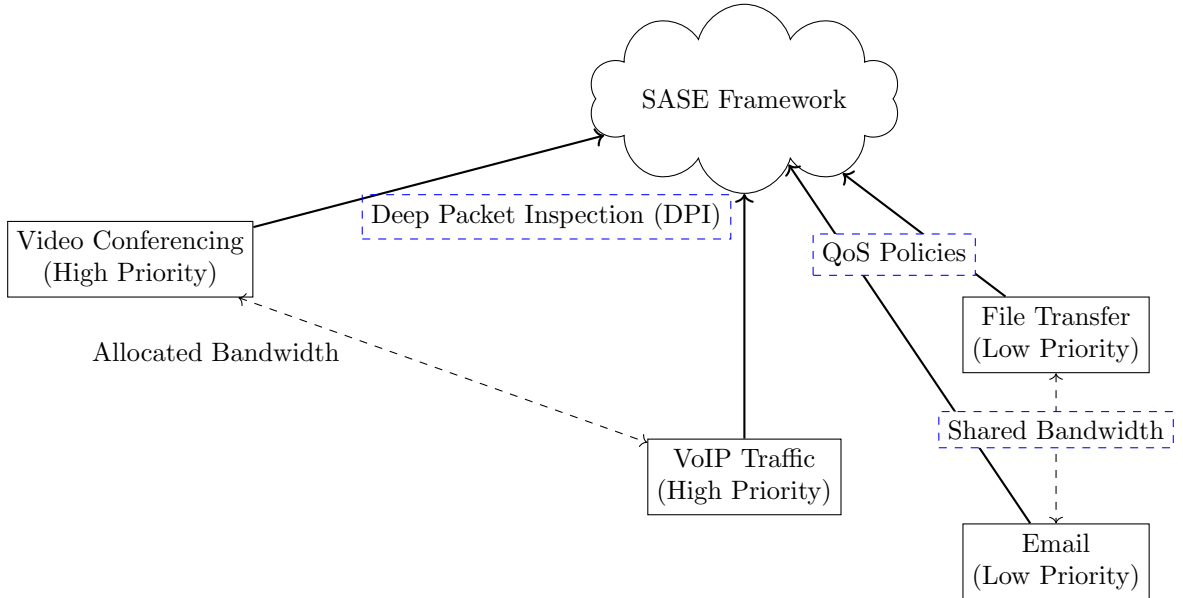## 4.2   Traffic Prioritization and Quality of Service (QoS)



Figure 4: QoS Implementation in SASE: Prioritization of Real-time Traffic using Deep Packet Inspection

Traffic prioritization through Quality of Service (QoS) is an essential technique for managing network performance in a Secure Access Service Edge (SASE) environment. By implementing QoS, network administrators can allocate bandwidth to specific types of traffic based on the needs of the applications in use. This allows real-time traffic, such as video conferencing and voice communications, to receive preferential treatment over less time-sensitive data, such as file downloads or email, ensuring that the network resources are used efficiently without introducing unnecessary delays.

QoS operates by classifying and marking network traffic according to predefined policies, which specify the priority for different applications or services. For example, voice or video traffic can be assigned a higher priority, ensuring minimal delay and jitter, while other data, like bulk transfers, is placed into a lower priority queue. This enables the network to deliver a more consistent quality of service for applications that are sensitive to latency and jitter. In a SASE environment, these policies can be enforced dynamically across a distributed architecture, allowing traffic to be prioritized at the edge or across various network paths, depending on the real-time conditions of the network.

Deep packet inspection (DPI) is often used in conjunction with QoS to identify the types of traffic flowing through the network, enabling more granular traffic classification and prioritization. With DPI, the SASE architecture can recognize traffic patterns and classify applications based on their specific requirements. For instance, real-time traffic, such as video conferencing or remote desktop sessions, can be identified and prioritized to ensure smooth operation, while less time-sensitive traffic is delayed or queued appropriately. This approach prevents congestion and ensures that high-performance applications receive the necessary bandwidth and low-latency paths required to function efficiently.

By optimizing traffic flow using QoS, SASE systems can maintain a balance between performance and resource allocation, especially in environments where bandwidth is constrained or multiple applications are competing for the same network resources. This is useful in distributed architectures, where network paths may vary in quality and capacity depending on the geographic location or the underlying network infrastructure. Implementing traffic prioritization in

such environments ensures that critical tasks are completed without unnecessary delays, maintaining the overall performance of the network.

The following table outlines how different types of traffic can be prioritized using QoS policies in a SASE environment, based on their sensitivity to delay and bandwidth requirements.

Table 6: Traffic Prioritization in a SASE Environment

| Traffic Type | Priority Level | Latency Sensitivity | Bandwidth Requirement |
|---|---|---|---|
| Video Conferencing | High | High | Moderate to High |
| VoIP (Voice over IP) | High | High | Low |
| File Downloads | Low | Low | High |
| Email | Low | Low | Low |
| Real-Time Data Analytics | High | Moderate to High | Moderate |
| Web Browsing | Medium | Low to Moderate | Moderate |

As shown in Table 6, applications like video conferencing and VoIP are highly sensitive to latency and jitter, requiring higher prioritization in a SASE deployment to maintain real-time performance. In contrast, activities such as file downloads or email transmission, which can tolerate delays without significant impact, are assigned lower priority, ensuring that high-demand services are not disrupted. By dynamically managing the distribution of network resources, QoS ensures that time-sensitive applications are able to function smoothly even under variable network conditions.

QoS implementation in SASE also extends beyond traffic prioritization to include mechanisms for traffic shaping and bandwidth reservation, which further optimize network utilization. Traffic shaping involves regulating the flow of data into the network to prevent congestion, while bandwidth reservation allocates a minimum amount of network capacity to high-priority applications. These methods, combined with QoS, provide a comprehensive approach to managing network traffic, ensuring that applications receive appropriate resources based on their operational requirements.

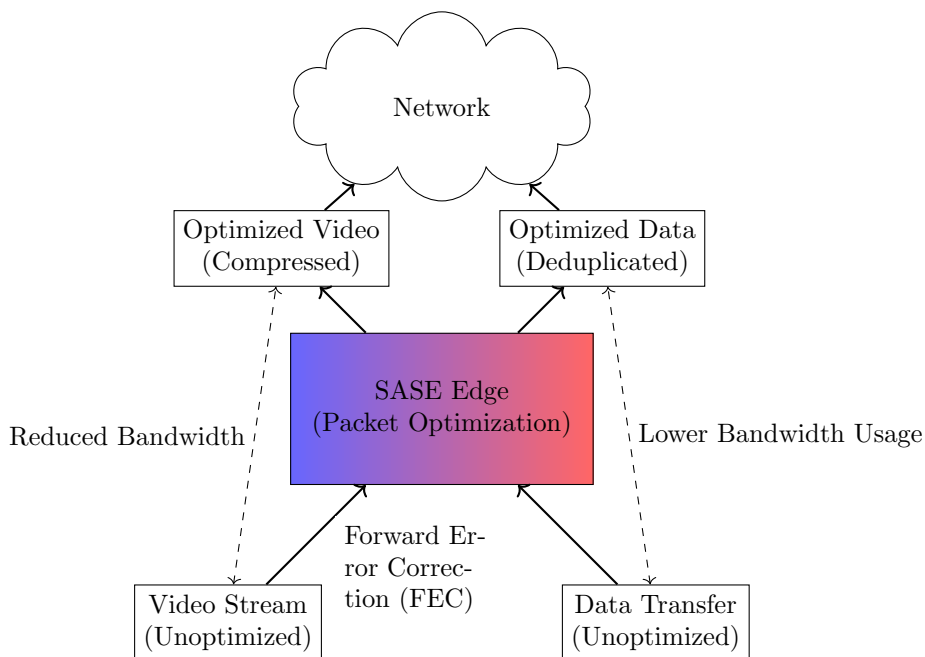## 4.3   Packet Optimization and Compression Techniques



Figure 5: Packet Optimization Techniques in SASE: Compression, Deduplication, and Forward Error Correction

Packet optimization and compression techniques help improve the efficiency of data transmission in a Secure Access Service Edge (SASE) environment, especially for high-volume applications. Techniques such as packet deduplication, compression, and forward error correction (FEC) reduce the size of data being transmitted, lowering the bandwidth needed without significantly impacting performance. This becomes important in networks handling large volumes of traffic, as it helps manage resource usage while maintaining service quality.

Packet deduplication is a straightforward method to reduce unnecessary traffic. In many cases, identical packets are sent repeatedly, consuming network bandwidth. Deduplication identifies and removes these redundant packets, ensuring that only unique data is transmitted. By applying this technique at the edge, where remote users or IoT devices often send repetitive data, SASE can reduce the volume of traffic before it enters the broader network. This reduces the strain on network resources without sacrificing data integrity.

Compression, another widely used optimization technique, reduces the size of transmitted data by encoding it more efficiently. For instance, video and audio streams can be compressed to lower their bit rates while still maintaining adequate quality. In a SASE environment, compression can be applied at the edge, minimizing the size of data before it is sent across the network. This is especially useful for bandwidth-intensive applications like video conferencing or real-time analytics, where even small reductions in data size can significantly lower bandwidth usage and improve transmission speeds.

Forward error correction (FEC) is a method that helps correct data transmission errors without requiring retransmissions, which can introduce delays. In environments where packet loss is common, such as over unstable network links, FEC can improve the reliability of data transmission by adding redundant information that helps recover lost packets. This reduces the need for retransmissions, which would otherwise increase bandwidth consumption and delay. FEC is useful in SASE environments where network performance can vary, as it helps maintain data integrity without overwhelming the network with repeated transmissions.

Table 7: Packet Optimization Techniques in SASE

| Technique | Effect on Bandwidth | Impact on Latency | Use Case |
|---|---|---|---|
| Packet Deduplication | Reduces redundant traffic | Minimal | High-volume IoT data, remote users |
| Compression | Lowers data size | Low | Video conferencing, large file transfers |
| Forward Error Correction (FEC) | Reduces retransmissions | Can increase slightly | Unreliable networks, packet loss-prone environments |

As shown in Table 7, these techniques, when applied appropriately, help reduce bandwidth usage without significantly affecting latency. Deduplication removes unnecessary traffic, compression reduces data size, and FEC improves transmission reliability. Applying these techniques at the edge in a SASE architecture ensures that data is optimized before being sent across the network, helping to manage bandwidth while still maintaining performance for various applications.

## 4.4   SD-WAN Advancements: Multipath Routing and Dynamic Bandwidth Allocation

Advancements in Software-Defined Wide Area Networking (SD-WAN) offer valuable improvements for optimizing Secure Access Service Edge (SASE) architectures for applications that demand low-latency and high-bandwidth connectivity. Techniques such as multipath routing and dynamic bandwidth allocation enable SASE to provide more reliable and efficient network performance by responding
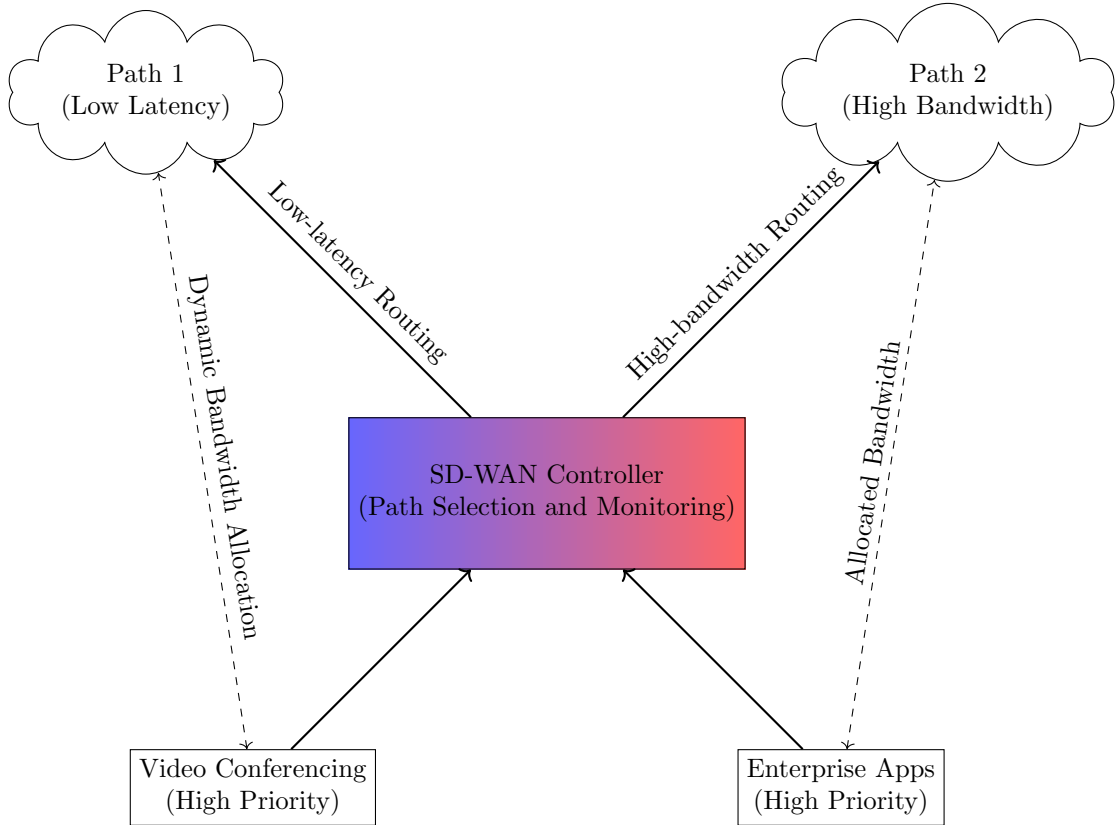
Figure 6: SD-WAN Optimization for SASE: Dynamic Path Selection and Bandwidth Allocation for High-Priority Applications

in real-time to varying network conditions. These techniques are especially relevant in distributed environments where network paths may differ significantly in terms of latency, jitter, and packet loss, impacting the quality of service for critical applications.

Multipath routing, a core feature of modern SD-WAN, allows traffic to be dynamically routed over multiple available paths based on real-time performance metrics. By continuously monitoring network conditions such as latency, jitter, and packet loss, SD-WAN can select the optimal path for traffic to ensure that high-priority, latency-sensitive applications, such as video conferencing or real-time analytics, are transmitted over the fastest and most reliable links. This is beneficial in environments where network conditions can fluctuate, such as across diverse geographic regions or when using multiple types of WAN connections, including broadband, MPLS, or LTE. Multipath routing improves not only the speed but also the reliability of data transmission by ensuring that traffic can be rerouted dynamically in the event of link degradation or failure.

The ability of SD-WAN to adapt to changing network conditions is central to maintaining the performance of real-time applications. For instance, in a video conferencing scenario, even slight increases in latency or jitter can degrade the quality of the experience, causing audio or video to become out of sync. By using multipath routing, SD-WAN ensures that such traffic is continuously routed over the path that offers the lowest latency and the highest stability. Should a particular link become congested or experience high packet loss, the system can immediately switch to an alternative path with better performance, maintaining the quality of service without requiring manual intervention.

Dynamic bandwidth allocation further enhances the capabilities of SD-WAN by ensuring that network resources are distributed efficiently based on the needs of different applications. SD-WAN continuously evaluates bandwidth usage and can adjust the allocation dynamically to ensure that high-priority applications have sufficient bandwidth, even during periods of heavy network usage. For ex-

ample, during peak usage periods, when many users are accessing cloud services or streaming high-definition video, SD-WAN can prioritize bandwidth for critical applications like video conferencing or VoIP, ensuring they receive the resources necessary to maintain smooth, uninterrupted performance. Meanwhile, less urgent traffic, such as bulk file transfers or background processes, can be deprioritized to prevent network congestion.

Dynamic bandwidth allocation is useful in scenarios where bandwidth availability is limited or fluctuates, such as in remote locations or branch offices relying on mixed network types (e.g., broadband and satellite links). In such cases, SD-WAN can dynamically balance bandwidth across available paths, ensuring that critical applications are prioritized while less important traffic is throttled or rerouted. This capability ensures efficient utilization of network resources, minimizing the risk of congestion that could otherwise lead to performance degradation for latency-sensitive applications.

The combination of multipath routing and dynamic bandwidth allocation in modern SD-WAN technologies aligns well with the goals of SASE, where performance optimization is necessary for distributed and cloud-based applications. These features allow SASE to maintain a high level of flexibility and responsiveness to changing network conditions, ensuring that applications continue to function optimally, regardless of location or connection type.

The following table summarizes the benefits of these SD-WAN advancements in a SASE environment:

Table 8: SD-WAN Advancements in Multipath Routing and Dynamic Bandwidth Allocation

| Technique | Benefit | Use Case |
|---|---|---|
| Multipath Routing | Optimal path selection based on real-time conditions | Video conferencing, real-time analytics |
| Dynamic Bandwidth Allocation | Efficient distribution of bandwidth during peak usage | Cloud applications, VoIP, file transfers |
| Failover | Ensures continuity in case of link failure | Mission-critical applications |
| Traffic Prioritization | Prioritizes latency-sensitive traffic | High-priority business applications |

Table 8 illustrates how SD-WAN's multipath routing and dynamic bandwidth allocation directly support the performance needs of latency-sensitive and high-bandwidth applications. By enabling the network to automatically adjust to current conditions, SD-WAN ensures that SASE can provide the necessary levels of service quality and reliability, especially for real-time communication and cloud-based applications that are sensitive to variations in network performance.

## 4.5   AI and ML-Driven Traffic Management

---

**Algorithm 1:** AI and ML-Driven Traffic Management in SASE

---

**Data:** Real-time network traffic data, historical traffic patterns,
application priorities

**Result:** Optimized traffic routing and bandwidth allocation

**Initialize:**

    $MLModel \leftarrow$ Trained ML model on historical network traffic data;

    $TrafficQueue \leftarrow$ Current real-time traffic;

    $PriorityList \leftarrow$ List of critical applications;

    $Bandwidth \leftarrow$ Available bandwidth in network;

    $NetworkPaths \leftarrow$ Set of all available paths for routing;

**Function** *TrafficManagement()*:

    **for each** $TrafficFlow\ in\ TrafficQueue$ **do**

        $Prediction \leftarrow MLModel.predict(TrafficFlow)$;

        **if** $Prediction == HighUsage$ **then**

            Allocate additional bandwidth to priority traffic;

        **else**

            Maintain standard bandwidth allocation;

        **end**

    **end**

    Monitor real-time network conditions;

    **for each** $TrafficFlow\ in\ TrafficQueue$ **do**

        $Anomaly \leftarrow DetectAnomaly(TrafficFlow)$;

        **if** $Anomaly == True$ **then**

            Adjust routing for $TrafficFlow$;

            Reprioritize traffic if necessary;

        **end**

    **end**

**Function** *DetectAnomaly(TrafficFlow)*:

    $Baseline \leftarrow$ Established normal behavior for $TrafficFlow$;

    **if** $TrafficFlow \notin Baseline$ **then**

        **return** True;

    **else**

        **return** False;

    **end**

**Function** *AdjustRouting(TrafficFlow)*:

    $BestPath \leftarrow SelectOptimalPath(NetworkPaths, TrafficFlow)$;

    Route $TrafficFlow$ through $BestPath$;

**Function** *SelectOptimalPath(NetworkPaths, TrafficFlow)*:

    Evaluate each path for latency, jitter, and packet loss;

    **return** Path with lowest latency and highest reliability;

**do**

    TrafficManagement();

    Continuously update $MLModel$ with real-time data;

**while** *network is active*;

---

Artificial intelligence (AI) and machine learning (ML) provide advanced capabilities for enhancing traffic management within a Secure Access Service Edge (SASE) architecture for real-time and bandwidth-sensitive applications. By leveraging AI and ML algorithms, SASE can dynamically adjust network configurations, enforce security policies, and optimize traffic routing based on real-time traffic patterns and network conditions. This enables the network to operate more efficiently, while also ensuring that critical applications maintain the necessary bandwidth and low latency required for optimal performance.

One of the key advantages of AI and ML-driven traffic management is the ability to predict network usage patterns and adjust resources preemptively. For example, by analyzing historical data on traffic loads, AI models can identify periods of high demand, such as peak hours during a workday, and allocate additional bandwidth to critical applications like video conferencing or real-time analytics

before congestion occurs. This predictive capability allows the network to remain proactive rather than reactive, ensuring that high-priority applications are not impacted by sudden surges in traffic. For instance, during a company-wide virtual meeting, AI algorithms can anticipate the increased bandwidth demand and dynamically allocate resources to ensure smooth communication and avoid disruptions caused by network congestion.

Moreover, AI-driven traffic management can analyze traffic flows in real-time to detect anomalies, such as unusual spikes in latency or packet loss, which may indicate potential network bottlenecks or performance issues. In such cases, the SASE architecture can automatically adjust routing and traffic prioritization to mitigate these issues before they significantly impact application performance. For example, if AI-driven analysis detects increased jitter or packet loss on a primary network path, the system can reroute traffic over a less congested path, minimizing disruptions for latency-sensitive applications. This automated approach to traffic optimization reduces the need for manual intervention and improves the overall efficiency of the network.

Machine learning models can also enhance security in SASE environments by detecting abnormal traffic patterns that may indicate malicious activity. ML algorithms can analyze vast amounts of traffic data to identify deviations from established baselines, such as a sudden increase in traffic from a particular device or location that deviates from normal behavior. This type of anomaly detection enables SASE systems to respond in real-time to potential threats, dynamically adjusting security policies, blocking suspicious traffic, or isolating affected network segments without compromising the performance of legitimate traffic.

Another significant advantage of AI and ML-driven traffic management is the ability to continuously learn from changing network conditions and optimize network performance over time. As ML models are exposed to more data, they can refine their predictions and traffic management strategies, leading to more accurate bandwidth allocation and improved routing decisions. This continuous learning process ensures that SASE systems can adapt to evolving network environments, including shifts in user behavior, application usage, and network capacity, all while maintaining high levels of performance and security.

The integration of AI and ML within SASE also supports adaptive traffic prioritization. Instead of relying solely on predefined QoS rules, AI can make dynamic prioritization decisions based on real-time application needs and network conditions. For instance, if the system detects that a video conferencing session is experiencing degraded performance due to network congestion, it can prioritize that traffic temporarily, ensuring a smoother experience. Similarly, less urgent traffic, such as background software updates or file transfers, can be deprioritized automatically during periods of high demand, allowing more critical services to function without interruption.

As shown in Table 9, the application of AI and ML in traffic management introduces several key improvements, from more accurate predictions of network demand to better handling of real-time anomalies. These capabilities allow SASE systems to be more responsive and adaptable, ensuring that performance-sensitive applications receive the resources they need without sacrificing overall network efficiency or security.

## 5  Conclusion

The ongoing evolution of enterprise networks, driven by the increasing adoption of cloud-based services and the growing prominence of edge computing, necessitates a comprehensive reevaluation of how organizations approach both network performance and security. Traditional network security models, which relied on centralized inspection points, often led to latency and bottlenecks, thus significantly impeding the performance of real-time, latency-sensitive applications such as video conferencing, real-time data analytics, and Internet of Things (IoT) systems. Secure Access Service Edge (SASE) represents a critical step forward in addressing these challenges by merging advanced wide area networking (WAN) capabilities with integrated security functions into a unified framework. This ar-

Table 9: AI and ML-Driven Enhancements in SASE Traffic Management

| Function | AI/ML Capability | Benefit |
| --- | --- | --- |
| Traffic Prediction | Forecast high-usage periods | Proactive bandwidth allocation to prevent congestion |
| Anomaly Detection | Identify unusual traffic patterns | Real-time detection of potential bottlenecks or security threats |
| Dynamic Routing | Optimize path selection based on real-time data | Reduced latency and improved application performance |
| Adaptive Prioritization | Adjust traffic priorities dynamically | Ensures critical applications maintain high performance during congestion |
| Continuous Learning | Refine predictions and optimizations over time | Improved accuracy in traffic management and resource allocation |

chitecture provides a foundation for addressing the growing demands placed on enterprise networks, while also maintaining the critical balance between performance and security for these highly sensitive applications (Sabella et al., 2021).

SASE combines Software-Defined Wide Area Networking (SD-WAN) technology with security services such as Secure Web Gateway (SWG), Firewall-as-a-Service (FWaaS), and Zero Trust Network Access (ZTNA). Through this architecture, organizations gain the ability to deliver secure, scalable network access across geographically dispersed environments. However, latency-sensitive applications, which require near-instantaneous processing of data, highlight several weaknesses in conventional implementations of SASE. Applications such as real-time video conferencing, rapid data analytics, and large-scale IoT deployments require both high throughput and minimal latency, which can be disrupted by security services that introduce overhead in the form of encryption, deep packet inspection, and traffic analysis. As such, optimizing SASE to support these types of high-performance applications is critical.

In exploring the complexities of optimizing SASE for these environments, a range of strategies and innovations come to the forefront. Among these, the integration of edge computing plays a pivotal role, as does the deployment of packet optimization, traffic prioritization, and dynamic network management techniques enabled by modern SD-WAN advancements. Artificial intelligence (AI) and machine learning (ML) have also begun to emerge as key tools in this space, offering unprecedented potential to enhance the intelligent management of traffic flows, dynamically adjusting to changing network conditions in real time.

At its core, SASE provides a cloud-native architecture that decentralizes network management and security enforcement. This decentralized model contrasts sharply with traditional centralized architectures, which often necessitated backhauling traffic to a central location for security inspection. Such backhauling can introduce significant latency, which is problematic for real-time applications that cannot tolerate delays. SASE's architecture, by operating closer to the edge of the network, minimizes this latency while still enforcing the necessary security controls. However, ensuring the seamless integration of security functions with high-performance networking remains a key challenge.

Real-time applications like video conferencing are acutely sensitive to network performance, requiring minimal latency and stable connections to deliver high-quality user experiences. Even small delays can cause degraded audio-video synchronization, dropped frames, and communication interruptions, significantly impairing usability. The challenge in a SASE environment arises from the need to apply security measures—such as encryption, deep packet inspection, and threat detection—while maintaining the low-latency and high-bandwidth conditions that real-time video conferencing demands. Any additional processing overhead in-

troduced by security functions can disrupt the smooth flow of data, leading to performance issues. Consequently, optimizing the performance of SASE for these applications requires techniques that can minimize security-induced latency while preserving the integrity and quality of the connection.

The demands placed on real-time data analytics are equally stringent, as these applications frequently operate in environments where delays in data processing can directly impact business outcomes. Industries such as finance, healthcare, and industrial automation rely heavily on rapid data processing to support decision-making, often in scenarios where timing is critical. In a SASE framework, security measures such as encryption, data inspection, and intrusion detection must be implemented without introducing bottlenecks that could compromise the timeliness of these analytics. As the volume of data continues to grow in large-scale analytics operations, optimizing SASE to handle both the security and performance requirements of these environments becomes an essential focus.

IoT networks add yet another layer of complexity to the equation, given the scale and diversity of devices generating continuous streams of data. In high-stakes IoT environments such as autonomous vehicle systems, smart grids, and healthcare monitoring, communication channels must remain low-latency and highly reliable. The challenge with IoT devices is further compounded by their inherent lack of built-in security features, placing a greater burden on the network infrastructure to ensure secure data transmission. SASE, in this context, must be optimized to handle not only the large volume of data but also the unique security challenges posed by IoT devices, balancing the need for secure communication with the high performance required for real-time operations.

One of the most effective strategies for reducing latency in a SASE environment is the integration of edge computing. By processing data closer to its source, edge computing dramatically reduces the need for data to travel across the network to a centralized data center, which is beneficial for latency-sensitive applications. In the context of SASE, edge computing enables security services to be hosted closer to the application, allowing security functions such as threat detection and data inspection to occur locally. This eliminates the need for the data to traverse long network paths for security processing, thus minimizing latency while still maintaining a high level of security. Edge computing is especially important for applications like video conferencing and IoT, where even milliseconds of delay can significantly affect performance.

In addition to edge computing, the implementation of traffic prioritization through Quality of Service (QoS) policies serves as another crucial optimization technique. QoS allows administrators to allocate network resources based on the importance of the traffic, ensuring that critical applications, such as video conferencing or real-time analytics, receive the bandwidth and low-latency conditions they require. SASE can support QoS by applying intelligent traffic management policies that prioritize performance-sensitive traffic over other, less urgent data streams. Deep packet inspection, a core feature of SASE, allows the system to identify and prioritize real-time traffic, ensuring that high-priority applications are not delayed by background processes or less time-sensitive communications.

Packet optimization techniques further enhance SASE's ability to support high-bandwidth, low-latency applications. Techniques such as packet deduplication, compression, and forward error correction (FEC) are critical in reducing the bandwidth required for data-heavy applications video streaming and large-scale IoT deployments. Compression allows video streams to be transmitted using less bandwidth without sacrificing quality, while packet deduplication removes redundant data from the transmission stream, thus conserving network resources. These optimizations, when applied at the network edge within a SASE framework, allow traffic to flow more efficiently, reducing both latency and bandwidth consumption.

The advancements in SD-WAN technology also contribute significantly to the optimization of SASE for real-time, high-bandwidth applications. Multipath routing is one such SD-WAN feature that enhances network performance by dynamically selecting the optimal path for traffic based on current network conditions. By continuously monitoring metrics such as latency, jitter, and packet loss, SD-WAN can ensure that latency-sensitive applications are routed over the fastest and most reliable network paths. Dynamic bandwidth allocation further complements

this by ensuring that high-priority applications have access to sufficient network resources during peak usage periods, thus preventing congestion and ensuring consistent performance.

Artificial intelligence (AI) and machine learning (ML) offer additional potential for optimizing SASE in terms of dynamic traffic management. AI-driven solutions can predict network traffic patterns and adjust routing and prioritization strategies in real time, ensuring that critical applications are always allocated the necessary resources. For example, AI can identify periods of peak network usage and allocate additional bandwidth to real-time applications in advance of potential congestion. Furthermore, AI can detect anomalies in traffic patterns, allowing the system to automatically adjust network configurations to mitigate potential performance issues before they impact users. These intelligent traffic management capabilities significantly enhance SASE's ability to support latency-sensitive and high-bandwidth applications in dynamic and unpredictable network environments. Edge computing, packet optimization, traffic prioritization, and AI-driven traffic management each operate with different mechanisms and objectives, which can lead to operational bottlenecks if not synchronized precisely. For instance, while edge computing reduces latency by processing data locally, its benefits may be counteracted by delays introduced during packet optimization procedures such as forward error correction (FEC) or real-time data compression. These packet-level optimizations introduce overhead that, if not executed in tandem with edge processing, could introduce additional latencies, negating the intended performance gains for latency-sensitive applications such as real-time analytics or video streaming. Moreover, the interaction between security enforcement and performance tuning complicates traffic prioritization. Techniques such as deep packet inspection (DPI) or encryption impose additional computational loads that can conflict with the performance needs of real-time applications, where minimal latency is critical. For example, when prioritizing video conferencing traffic, the additional latency from encryption or DPI could cause degradation in video quality, frame rate, or synchronization in high-resolution streams. The intricate dependencies among these optimizations require sophisticated orchestration mechanisms that can dynamically manage trade-offs between security and performance. This challenge is compounded when attempting to maintain a consistent Quality of Service (QoS) across a distributed and heterogeneous network, further limiting the agility of SASE solutions in highly variable environments.

The practical deployment of edge computing as a key latency reduction mechanism within SASE is severely constrained by uneven infrastructural distribution, especially across geographically disparate regions. Edge nodes, which are critical for minimizing the physical distance that data must travel, are often deployed in urban or metropolitan areas where infrastructure investment is economically justified. In contrast, rural and remote areas frequently lack access to robust edge infrastructures, which forces data from latency-sensitive applications to traverse greater distances, thereby introducing delays that edge computing was intended to mitigate. For example, real-time applications, such as industrial IoT monitoring or autonomous vehicle communication systems in remote locations, may experience significant delays if edge nodes are not available within a close proximity. Even within urban environments, the cost and complexity of deploying and maintaining sufficient edge resources to handle diverse and high-bandwidth applications pose substantial barriers. Each edge node must be equipped with computing power, storage, and security capabilities sufficient to handle real-time processing demands, such as high-definition video analytics or low-latency IoT telemetry. Additionally, securing these decentralized, geographically distributed edge nodes introduces further complications. Unlike centralized data centers, edge nodes are more vulnerable to physical attacks or localized cyber threats, and ensuring that security policies are uniformly enforced across these distributed systems creates additional overhead. This not only impacts performance but also complicates compliance with regulatory frameworks in industries like finance or healthcare where data locality and security are tightly controlled. As a result, the limitations in edge computing infrastructure deployment constrain the scalability and reliability of SASE optimizations in global or multi-regional networks.

The reliance on AI and machine learning (ML) to optimize traffic management

in SASE, while theoretically promising, introduces significant concerns regarding the robustness, interpretability, and reliability of these systems in real-world scenarios. AI-driven algorithms for path selection and traffic prioritization rely on large-scale data inputs, such as real-time metrics on latency, jitter, and packet loss, to dynamically adjust routing decisions. However, these algorithms are only as effective as the data they receive and the models on which they are trained. Inconsistent or biased training data can lead to suboptimal routing decisions in edge cases or highly dynamic network environments. For example, AI models that have been trained primarily on data from densely populated areas may struggle to accurately optimize traffic in rural or sparsely populated regions where traffic patterns and network conditions differ significantly. This misalignment could lead to inefficient path selection, where latency-sensitive traffic is routed through suboptimal network paths, exacerbating delay rather than minimizing it. Additionally, the black-box nature of many AI and ML models raises concerns about the interpretability of the decisions they make. Network administrators and engineers often lack visibility into how these algorithms prioritize traffic or allocate bandwidth, which creates challenges when diagnosing performance bottlenecks or security lapses. For example, if an AI system inadvertently deprioritizes critical traffic such as emergency communication in a healthcare network, it could be difficult to trace the algorithmic decision-making that led to this outcome.

# References

Chandramouli, R. and Chandramouli, R. (2022). *Guide to a Secure Enterprise Network Landscape*. US Department of Commerce, National Institute of Standards and Technology.

Chen, R., Yue, S., Zhao, W., Fei, M., and Wei, L. (2022). Overview of the development of secure access service edge. In *International Conference On Signal And Information Processing, Networking And Computers*, pages 138–145. Springer.

Gandhi, I., Barton, R., and Henry, J. (2022). Obtaining visibility into a secure access services edge (sase) network.

Islam, M. N., Colomo-Palacios, R., and Chockalingam, S. (2021). Secure access service edge: A multivocal literature review. In *2021 21st International Conference on Computational Science and Its Applications (ICCSA)*, pages 188–194. IEEE.

Jani, Y. (2021). The role of sql and nosql databases in modern data architectures. *International Journal of Core Engineering & Management*, 6(12):61–67.

Kaur, T. (2018). Secure access service edge (sase): Extending network security to client.

Sabella, D., Reznik, A., Nayak, K. R., Lopez, D., Li, F., Kleber, U., Leadbeater, A., Maloor, K., Baskaran, S. B. M., Cominardi, L., et al. (2021). Mec security: Status of standards support and future evolutions. *ETSI white paper*, 46(1):26.

van der Walt, S. and Venter, H. (2022). Research gaps and opportunities for secure access service edge. In *International Conference on Cyber Warfare and Security*, volume 17, pages 609–619.

Yiliyaer, S. and Kim, Y. (2022). Secure access service edge: A zero trust based framework for accessing data securely. In *2022 IEEE 12th Annual Computing and Communication Workshop and Conference (CCWC)*, pages 0586–0591. IEEE.

Zhang, Z. (2019). Enterprise networking with secure acess service edge.

AFFILIATION OF ARUNKUMAR VELAYUTHAM
Arizona, USA

CLOUD SOFTWARE DEVELOPMENT ENGINEER AND TECHNICAL LEAD AT INTEL, ARIZONA, USA