

Universidade Federal de Campina Grande  
Centro de Engenharia Elétrica e Informática  
Coordenação de Pós-Graduação em Ciência da Computação

Resource-efficient management of stateful  
confidential applications

Clenimar Filemon Souza

Dissertação submetida à Coordenação do Curso de Pós-Graduação em Ciência da Computação da Universidade Federal de Campina Grande - Campus I como parte dos requisitos necessários para obtenção do grau de Mestre em Ciência da Computação.

Área de Concentração: Ciência da Computação

Linha de Pesquisa: Sistemas Distribuídos e Computação em Nuvem

Andrey Elísio Monteiro Brito  
(Orientador)

Campina Grande, Paraíba, Brasil

©Clenimar Filemon Souza, 27/10/2020

## Resumo

O alto nível de poder computacional oferecido por provedores de nuvem pública, combinado com a flexibilidade, eficiência e custo reduzido, fazem desta um ambiente atrativo para a implantação de serviços e aplicações. Contudo, num ambiente altamente dinâmico e sem controle sobre a localidade precisa de dados e código, ou quem pode acessá-los, surgem preocupações a respeito de confidencialidade e integridade. De fato, a maioria das preocupações relacionadas à adoção da nuvem pública é relacionada a segurança. Embora algumas aplicações e serviços possam tolerar estas preocupações em nome da redução de custo, da maior flexibilidade e simplicidade de gerenciamento, outras aplicações têm requisitos tão rigorosos de confidencialidade, integridade e isolamento que a maior parte dos recursos ofertados por provedores de nuvem pública é simplesmente inadequados.

Estas aplicações, a exemplo de *internet banking* e finanças, sistemas de saúde, de cidade inteligente (*smart grids*) ou serviços de documentos confidenciais, podem utilizar de tecnologias assistidas por hardware, como Ambientes de Execução Confiável, ou TEEs (do inglês, *Trusted Execution Environments*), para prover garantias de confidencialidade e integridade, mesmo em infraestruturas não-confiáveis. Os principais provedores de nuvem pública têm começado a oferecer instâncias com suporte a TEEs. Essas tecnologias, no entanto, comumente dependem de recursos de hardware especializados e escassos, o que se traduz em custos elevados e desempenho subótimo, especialmente para atender a uma larga escala. Este trabalho define um conjunto de requisitos de aplicação para confidencialidade, integridade e isolamento, e apresenta uma abordagem que visa à otimização do uso de recursos através do gerenciamento dinâmico de instâncias de aplicação em infraestruturas de nuvem confidencial. A abordagem é implantada em infraestruturas gerenciadas por Kubernetes, o gerenciador de contêineres padrão da indústria, e avaliada no contexto de uma aplicação de processamento de dados confidenciais providas por medidores de energia inteligentes (*smart meters*).

## **Abstract**

The high level of computing power offered by cloud providers, combined with the flexibility, efficiency and reduced costs make the public cloud an attractive environment for deploying services and applications. However, this highly dynamical environment and the uncertainty about the actual location of the application and data, or who has access to it, raise questions about confidentiality and integrity of running applications and services in such environments. In fact, most of the concerns that prevent an even broader adoption of public cloud providers are related to security. Although some applications can overlook such concerns in the name of the cost reduction and almost infinite resources, some other applications have such a strict set of requirements with respect to confidentiality, integrity and data isolation that most public cloud offerings are simply not suitable.

These applications, such as internet banking and finance, healthcare systems, smart grid or confidential document services, can rely on hardware-assisted technology, such as Trusted Execution Environments (TEEs), to provide confidentiality and integrity guarantees, even in untrusted infrastructures. The major public cloud providers have also started to offer TEE-enabled instances. However, these technologies usually rely on scarce hardware resources, that often translate to higher costs and subpar performance, especially when deploying for large scales. This work defines a set of application requirements w.r.t. confidentiality, integrity and data isolation, and presents a resource-efficient approach to dynamically manage large sets of application instances in confidential cloud infrastructures. This approach is deployed to infrastructures managed by Kubernetes, the industry-standard container orchestrator, and evaluated in the context of an application that manages sensitive smart meter data.

# Contents

<b>1</b>	<b>Introduction</b>	<b>1</b>
1.1	Problem . . . . .	3
1.2	Objective and research questions . . . . .	4
1.3	Contributions . . . . .	4
1.4	Document structure . . . . .	5
<b>2</b>	<b>Background</b>	<b>6</b>
2.1	Trusted Execution Environments with Intel SGX . . . . .	6
2.2	Kubernetes . . . . .	7
2.3	Related works . . . . .	10
<b>3</b>	<b>Use case: smart grids</b>	<b>12</b>
3.1	Architecture . . . . .	13
3.2	Supported operations . . . . .	13
3.2.1	LOGIN . . . . .	13
3.2.2	LOGOUT . . . . .	14
3.2.3	WRITE . . . . .	14
3.2.4	LONG_WRITE . . . . .	14
3.2.5	QUERY_AVG . . . . .	14
3.3	Implementation . . . . .	15
3.3.1	Server . . . . .	15
3.3.2	Client . . . . .	16
<b>4</b>	<b>Dynamic provisioning of stateful applications in TEEs</b>	<b>17</b>
4.1	Alternative approaches . . . . .	17
4.1.1	The trivial approach . . . . .	17

---

4.1.2	Dynamic provisioning: the on-demand approach . . . . .	18
4.2	Dynamic provisioning: the standby servers approach . . . . .	19
4.3	Integration with example application . . . . .	20
4.3.1	Application Wrapper . . . . .	20
4.3.2	Routing Manager . . . . .	22
4.3.3	Load balancer . . . . .	22
4.3.4	Data encryption and provisioning . . . . .	26
<b>5</b>	<b>Performance evaluation</b>	<b>27</b>
5.1	Experiment description . . . . .	27
5.2	Experiment environment . . . . .	28
5.3	Experiment tools . . . . .	29
5.4	Baseline performance . . . . .	29
5.5	Discussion . . . . .	33
5.6	Threats to validity . . . . .	39
<b>6</b>	<b>Final remarks</b>	<b>40</b>
6.1	Conclusion . . . . .	40
6.2	Future work . . . . .	41
6.3	Acknowledgments . . . . .	42

# List of Figures

4.1	Ingress controllers: latency in a 10-minute benchmark. . . . .	24
4.2	Ingress controllers: tail latency (p99) versus throughput in a 10-minute benchmark. .	25
5.1	Baseline: per-operation latency (Trivial approach). . . . .	31
5.2	Baseline: per-operation latency distribution (Trivial approach). . . . .	32
5.3	Baseline: average CPU utilization (Trivial approach). . . . .	33
5.4	Latency distribution for 56 clients, or 50% Active User Ratio (RF 10 seconds; WR 50%). . . . .	34
5.5	Latency distribution for 56 clients, or 50% Active User Ratio (RF 60 seconds; WR 50%). . . . .	35
5.6	Latency distribution for 56 clients, or 50% Active User Ratio (RF 10 seconds; WR 10%). . . . .	36
5.7	Latency distribution for 28 clients, or 25% Active User Ratio (RF 10 seconds; WR 50%). . . . .	37
5.8	Per-operation latency versus Active User Ratio (RF 10 seconds; WR 50%). . . . .	38

# List of Tables

3.1	Minimum enclave size requirements in SCONE. . . . .	16
3.2	List of execution parameters supported by the client application. . . . .	16
4.1	Server events and Application Wrapper actions. . . . .	21
4.2	Ingress controllers: total requests processed in a 10-minute benchmark. . . . .	25
5.1	Levels for baseline experiment with <i>Trivial</i> approach. . . . .	30
5.2	Latency percentiles for LOGIN operations. . . . .	30
5.3	Latency percentiles for LOGOUT operations. . . . .	30
5.4	Latency percentiles for QUERY_AVG operations. . . . .	30
5.5	Latency percentiles for WRITE operations. . . . .	31
5.6	Baseline: total memory allocated to application instances in the cluster. . . . .	33
5.7	Total memory allocated to application instances in the cluster. . . . .	38

# Chapter 1

## Introduction

The rise of public cloud providers is certainly one of the key factors that allowed the World Wide Web to transition to its current, almost ubiquitous state. The increasingly cheaper costs associated with the readily available computing power, more adaptable to the users needs, have long turned public cloud computing into the *de facto modus operandi* of the industry. In fact, a large amount of the internet is running on public cloud providers, such as Amazon Web Services, Microsoft Azure and Google Cloud Platform. Running on the cloud allows applications and services to be hosted at multiple different sites, providing not only high availability and fault tolerance, but also reducing the average response times of many people, contributing to an overall better user experience. In most cases, it is also more cost effective, as the costs of maintaining an owned data center (internet connection, electricity, cooling, operators, maintenance) tend to be far more expensive than paying a public cloud provider bill at the end of the month.

However, there are also many concerns, most of them directly related to security and/or privacy. Huge data centers and server farms help reduce costs and deliver availability and virtually unlimited computational power, but also make it nearly impossible to track the physical location of the applications and data being processed, or who has access to it.

Although many applications and services can overlook such concerns and choose the practicality and lower costs of public cloud providers, there are classes of application and services whose security and privacy requirements are not compatible with the standard public cloud offering. Such requirements are usually associated with private and/or sensitive data being processed. Data that must not be seen, leaked or modified by unauthorized parties. Another common issue is related to protection of intellectual property, where the client wants to protect the piece of software itself, be it the business logic, an algorithm, a machine learning model. Internet banking and other financial services, health



---

services (patient records and exams), social security, confidential document management, smart grid applications are examples that do require integrity and confidentiality and, in some cases, strong context isolation during runtime.

Public cloud providers usually offer some sort of security or privacy-aware features and services. More recently, with the term confidential computing gaining some traction, some providers have started to offer support to Trusted Execution Environments (TEEs). TEEs provide hardware-assisted confidentiality and integrity guarantees to the application execution. It is usually possible to attest the piece of software running inside protected regions of the memory to ensure it was not modified. Neither the operating system or a systems administrator, even with physical access to the machine, are able to read the contents of these regions. Since TEEs do require specialized hardware to establish trust and to allow the secure execution, through private memory regions and additional feature set, most chip manufacturers have hardware technologies that implement the TEE standard. Intel, for example, offers Security Guard eXtensions, or SGX [1; 9]; AMD offers Secure Encrypted Virtualization, or SEV [8]; and ARM offers TrustZone [17].

However, it is important to note that TEE technologies, such as the three mentioned before, have their downsides. They have inherently different security guarantees, which are based on the actual implementation of the technology. But regardless of the security guarantees and the feature set, these technologies impose some overhead, notably in performance [16], since the code runs in a protected, often limited, memory space, which also requires special setup or even attestation with external services. Another issue is development overhead, since these technologies often limit the scope of the application (e.g., limited system calls) or even require that the application be modified, or even rewritten. Intel SGX, for example, provides an SDK (Software Development Kit) that only supports C/C++ applications.

Another limitation is the hardware requirements, since most features are dependent on a specific feature set, often included in the chip. When the hardware is being offered by a cloud provider, the specialized hardware and the rather scarce amount of protected resources, as it is in most TEE technologies, translates directly into higher operational costs.

Therefore, it is a challenge to allow applications that leverage TEEs to ensure their strict set of confidentiality, integrity and isolation requirements to be deployed to third-party confidential cloud infrastructures. Depending on a specialized, scarce resources such as protected memory might be prohibitively expensive and slow, especially when running a large set of applications. Another challenge is to allow these complex architectures to benefit from well-established and ubiquitous infrastructure management tools.

## 1.1 Problem

The following set of requirements defines the class of applications that this work considers, and is shared among all use cases aforementioned, such as internet banking, financial applications, health care services, smart cities applications.

1. The application must be available at all times.
2. The application requires confidentiality and integrity guarantees at runtime, even when data is in use, as provided by a TEE technology.
3. The application runs multiple instances in separate processes, for better context isolation.
4. Each application instance serves only one client at a time, for the period that the client is active (e.g., within a session).
5. The client data is encrypted at rest.
6. All communications between the application and clients are end-to-end encrypted.
7. The application must be usable in an interactive form, with no noticeable interaction overhead.

This set of requirements gives a stateful quality to the application, since it holds a particular, unique state which belongs to the client being currently served. Moreover, the application is supposed to keep serving the same client until the end of the active period, which is henceforth denoted as a session.

Deploying a large set of stateful, isolated application instances within TEEs, where each application instance has a particular confidential state (e.g., the data of the user being served) raises a number of non-trivial challenges. How to distribute such instances across a large set of nodes? How to perform the correct routing of the messages if the communication uses end-to-end encryption? How to overcome the quality-of-service impact of having complex applications running inside TEEs, which are inherently slower than native, unprotected applications? How to monitor the state of this set of applications? How to optimize resource consumption in order to lower the operational cost of keeping this service running? How to enable a reasonable auto-scaling mechanism for this set of applications?

Each one of the challenges listed above is both non-trivial and relevant.

## 1.2 Objective and research questions

This work focus on a subset of the challenges previously posed. More specifically, this work evaluates the *Standby Servers approach* (Chapter 4) for dynamically provisioning a large set of isolated, stateful instances of an application that runs inside of TEEs. The research questions are defined as follows.

**Q1. Can this approach be implemented with off-the-shelf cloud orchestration tools (such as Kubernetes), minimizing the complexity of managing the infrastructure?**

**Q2. How to correctly and efficiently route messages that are end-to-end encrypted in a stateful setting?**

**Q3. What are the expected performance benefits when comparing such an approach with the trivial approach of having all instances running simultaneously?**

## 1.3 Contributions

The main contributions of this work are:

- An architecture that allows stateful, isolated applications that run in TEEs to be deployed in confidential cloud infrastructures<sup>1</sup>.
- An approach for routing end-to-end encrypted messages in a stateful application setting.
- An evaluation of the proposed architecture with respect to its viability, performance, quality-of-service and resource consumption.
- Insights on how to minimize the latency overhead imposed by TEEs and scarce specialized resources (e.g., protected memory).

This work assumes an adversary that wants to compromise the confidentiality or the integrity of the application processing sensitive data, by leveraging privileged access to the operating system (i.e., privileged software), to the network in which the message exchanges occur, or even to the physical machines that host the application. Our threat model does not consider availability, since the attacker is able to shutdown machines or the network, for example. In fact, most TEEs do not include availability in their security guarantees.

---

<sup>1</sup>This architecture is a joint work between UFCG and Scontain U.G., see Section 6.3.

## 1.4 Document structure

The rest of this document is organized as follows:

**Chapter 2. Background.** This chapter presents a set of concepts that are relevant to better understand the problem, the motivational use case and the proposed solution. It also overviews existing solutions to similar problems in the literature, and compares such solutions to the approach presented by this work.

**Chapter 3. Use case: smart grids.** This chapter introduces a use case to illustrate the problem. This example will be used throughout the work, and will serve as a base to illustrate the proposed solution and its evaluation.

**Chapter 4. Dynamic provisioning of stateful applications.** This chapter presents an approach for deploying stateful, isolated application instances that run inside of TEEs in a more efficient way. Two other alternative approaches are presented for comparison.

**Chapter 5. Evaluation.** This chapter presents the results of the analysis of the proposed architecture in terms of performance and resource consumption.

**Chapter 6. Final remarks.** This chapter presents final remarks and future work.

# Chapter 2

## Background

### 2.1 Trusted Execution Environments with Intel SGX

Trusted Execution Environments, or TEEs, provide confidentiality and integrity guarantees through special features, such as private regions of the memory, special CPU instructions or hardware-based attestation, which allow the underlying privileged software to be untrusted. There are several TEE technologies, and the most famous are Intel SGX, ARM TrustZone and AMD SEV.

Intel SGX [1; 9], or Secure Guard eXtensions, leverages a special set of CPU instructions to create private regions of the memory, called enclaves, and then protect the execution of sensitive applications in untrusted environments. An attacker with administrative privileges or physical access to the machine would not be able to see the contents of an enclave, which guarantees the confidentiality of the computation being executed. In addition, it is possible to attest the contents of an enclave, ensuring the integrity of what is running (i.e., it was not modified by anyone). Confidentiality and integrity are the key concepts of Intel SGX technology, which is available in most off-the-shelf Intel processors (e.g., Intel Core family has it since its 5th generation) and some server processors as well (e.g., Intel Xeon family).

The private region of the memory, which is called Enclave Page Cache, or EPC, and is part of the Processor Reserved Memory, has a limited size, usually 128 MB. All the enclave pages are stored in the EPC. If the enclave is bigger than the available EPC, old pages are encrypted and then evicted from the EPC to the main memory. This pagination process usually imposes a significant overhead to the application performance.

To run an application inside of Intel SGX enclaves, the developer must rewrite their code to use the Intel SGX SDK, which provides wrapper functions to run sensitive parts of the application inside

of the enclave. This requirement imposes a high effort to adapt applications to benefit from Intel SGX security guarantees. Not only the developer has to rewrite their code and use, in most cases, a completely new toolset, but also the toolset itself limits the possibilities, as Intel SGX SDK only supports C and C++ languages currently.

SCONE (Secure CONTainer Environment) [2] is a framework that allows entire applications to run inside of Intel SGX enclaves, thus providing the same integrity and confidentiality guarantees. SCONE supports a broader set of programming languages, and has also additional features, such a more robust attestation mechanism, transparent file system encryption and automated secret generation. However, having the whole application fit inside the enclave (whose EPC is limited in most cases) might impose some performance overhead, as the data needs to be encrypted before going to the main memory, and decrypted when coming back to the EPC.

SCONE relies on a set of tools, which include a modified GCC compiler and system libraries, such as *glibc* or *musl-libc*, to compile existing applications to run inside of Intel SGX enclaves. It supports dynamically and statically linked applications, and a broader set of high-level programming languages (such as C, C++, Python, Go, Java, JavaScript, C#).

Once inside the enclave, the application talks to the SCONE runtime, which handles all system calls through an asynchronous interface, with one or more queues, and two types of threads (*e-threads* and *s-threads*), one inside of the enclave and other on the outside (in contact with the operating system).

The SCONE runtime also enables two important features of the SCONE framework: remote attestation and transparent encryption.

## 2.2 Kubernetes

The grow in popularity faced by microservice architecture and container-based workloads [4] in the last decade have created the need for specialized tooling to manage and monitor large sets of production containers. This class of applications, known as *container orchestrators* act on top of an infrastructure (be it composed by baremetal machines, virtual machines or both), sometimes being able to let the operator act on a higher abstraction level, which avoid coupling between applications and infrastructure and, ultimately, increases overall portability.

Kubernetes [6; 7] is the de-facto container orchestrator of the industry. Released in 2014, the open-source project reused many concepts that had been applied at Google internally for years, helping manage millions of productions containers daily. The first release of Kubernetes also marked the

foundation of the Cloud Native Computing Foundation, or CNCF, together with Linux Foundation. The mission of CNCF is to define standards and spread best practices that aim at a portable, scalable and fault-tolerant world of applications, completely aligned with the ever-growing resource offer from cloud providers around the world.

Kubernetes manages an infrastructure by grouping a set of nodes in a *cluster*. Kubernetes clusters have two types of nodes: *masters* and *workers*. Master nodes are supposed to run the core components of Kubernetes, responsible for maintaining the cluster, such as its controller manager, scheduler and API server. It is a good practice to use the master nodes exclusively for workloads related to cluster management. It is common to see clusters with a single master node. However, in order to achieve more strict high-availability and fault-tolerance requirements, one should have at least 3 master nodes. All cluster state is persisted to *etcd*<sup>1</sup>, a distributed key-value storage.

One of the key aspects of Kubernetes is its application-centric approach. Its powerful declarative API offers objects, such as *Pods*, *Deployments*, *DaemonSets* and *Services*, that allow the operator to not care about infrastructure specifics, application scheduling or even resource allocation. This allows Kubernetes to manage fairly large infrastructures in an efficient way, almost completely transparent to the user. As of today, Kubernetes officially supports infrastructures with up to 5000 nodes. Another direct result of having a declarative, uniform API that takes care of infrastructure specifics is portability. An application that runs on a Kubernetes cluster will most likely run in any other Kubernetes cluster, regardless of the underlying infrastructure (for instance, on different cloud providers).

When a user deploys an application to a Kubernetes cluster, they specify the desired state of the application. A set of controllers then make sure that the state is the desired at all times, taking action whenever is necessary. This allows tasks such as replica management, failure recovery and application updates to be taken care automatically. Moreover, the Kubernetes API has been consistently evolved and extended over the years, resulting in a rich and powerful cloud-native ecosystem. Various tools and solutions for monitoring, logging, access control, storage and more have been created and integrate seamlessly to the Kubernetes API.

The main Kubernetes API objects and the behaviors associated with them are detailed below. They are *Pod*, *Deployment*, *Service*, *Ingress*, *DaemonSet* and *StatefulSet*. Kubernetes offers many more API objects, as well as provide an API for custom objects, called *CustomResourceDefinition*.

**Pod.** The smallest application unit. Represents one or more containers that share a common goal and can be seen as one. For instance, a database application, such as MariaDB or redis, bundled with an auxiliary container that acts a proxy or authorizer of said database. *Pods* have their own network

---

<sup>1</sup><https://etcd.io/>

namespace (IP address and ports), which is shared among its constituent containers. *Pods* are also considered to be ephemeral by the Kubernetes scheduler, which can terminate or move them for some reason (e.g., to free resources in a node under some kind of resource pressure). Consequently, *Pods* should always be created with an associated replication controller, which is offered by other objects, such as *Deployments* and *StatefulSets*.

**Deployment.** Represents a *pod* associated with a replication controller, which ensures the desired state is achieved. It is possible to define the amount of replicas and failure recovery policies. For instance, if a node of the cluster fails and becomes unavailable, causing a *deployment* to have less replicas than specified, the replication controller will act (i.e., spawn new replicas) until the desired state is established again.

**Service.** Allows *Pods* to be discovered and accessed by other entities. Since *Pods* are ephemeral, their IP addresses can vary constantly, which is not a reliable option for communication. A *service* is a layer 4 load balancer (according to the OSI model) for *Pods*. Through *labels* and *selectors* one can bind certain *Pods* to *services*. *Services* have an IP address that does not change and entry in the in-cluster domain name server (DNS), which can be accessed by other in-cluster entities (e.g., other *Pods*). To allow applications to be reached from outside the cluster, there are two distinct *service* types. *NodePort services* expose a random port within a limited range in all nodes of the cluster for the application associated with the *service*. This way, an outside client could talk to any node IP address on the specified port. *LoadBalancer services*, on the other hand, contact the underlying cloud provider to provision a managed load balancer for the application.

**Ingress.** *Ingress* are layer 7 load balancers (according to the OSI model). They expose *services* to outside of the cluster, through a public IP address. Being a layer 7 (i.e. application layer) load balancer gives *ingresses* a broader set of features, such as TLS termination, virtual host routing or canary deployments. It is important to note that Kubernetes does not implement the *ingress* logic, and the *ingress* API must be implemented by an external controller, chosen by the user. The *ingress controller*<sup>2</sup> is responsible for deploying the actual load balancer *Pods* and managing *ingress* objects. The most popular load balancers on the market, such as HAProxy, NGINX and Envoy, offer an *ingress controller* solution.

**DaemonSet.** Represents a *Pod* and a replication controller that ensures that there will be exactly one replica per cluster node. It is possible to filter out nodes through selectors and labels. Useful to deploy applications such as log collectors, storage, monitoring or attestation agents, and many more.

**StatefulSet.** Represents a *Pod* and a stateful replication controller, which allows each replica to

---

<sup>2</sup><https://kubernetes.io/docs/concepts/services-networking/ingress-controllers/>



have its identity persisted. Useful to deploy stateful applications, such as a Database Management System. It is important to note the Pods are still ephemeral, but now they have an identity that is unique and persistent, which can be combined with data persistence techniques, such as volume provisioning.

## 2.3 Related works

AMD SEV [8] is a TEE technology provided by AMD. Based on AMD Memory Encryption technology, it has recently faced an increase in popularity since it was announced as the default TEE for Google's new confidential computing offerings. Applications do not require software change, and the protected application can use as much memory as it is available on the machine. However, SGX offers a more strict control over the protected memory (including memory integrity protection and access control) and defines a clear boundary between the trusted and untrusted portions of an application. Therefore, for applications that deal with sensitive data that require confidentiality and integrity, SGX is most suitable, despite its performance overhead [10].

ARM TrustZone [17] focus more on embedded and mobile systems with a single purpose. In TrustZone architecture, the processor is partitioned in secure and insecure portions, offering no further isolation within these regions. Intel SGX supports multiple enclaves, which are isolated from each other, enabling multi-purpose, multi-tenant use cases.

For Intel SGX SDK, the need to rewrite the code to use the enclave interface and define secure and insecure functions is prohibitive in most cases. Also, enclave attestation becomes a non-trivial process with SDK. Graphene-SGX [15] and SGX-LKL [11] are shielded execution frameworks that allow unmodified applications to run inside Intel SGX enclaves. Unlike SCONE, however, they include the operating system library within the enclave, broadening the attack surface considerably.

Ataíde *et al.* [3] and Sampaio *et al.* [12] use a similar approach of having generic application instances in TEEs managed by an application wrapper in the context of batch processing in untrusted infrastructures. As both works target batch processing, they make design decisions that are not adequate for handling interactive applications, as is our case. Segarra *et al.* [13] evaluates stream processing of confidential batch jobs using a modified Spark inside of SGX enclaves. Although streaming favors latency, it is still a rather different approach as compared to interactive services.

Brenner and Kapizta [5] propose a generic confidential Function-as-a-Service (FaaS) platform that shares some goals with this work, such as minimizing the performance overhead (specially cold-start times) and optimizing resource usage. However, in order to achieve it, their application instances,

---

which are JavaScript functions, share the same secure interpreter, and thus the same Intel SGX enclave. Their isolation is provided by software. Our approach provides further isolation by providing each application instance its own enclave. We also support a broader range of applications, since we aim at deploying services, which are inherently more complex than standalone functions.

Teodoro *et al.* [14] propose a stateful session-based load balancing mechanism for clustered web servers. Their approach relies on a central component, called *session directory*, that assigns user sessions to back-end servers, similarly to our Routing Manager. However, their work assumes that the load balancing is performed by the application itself, which has the ability to migrate sessions to other servers based on the current server load. Our approach for message routing relies on battle-tested, off-the-shelf load balancing solutions, such as HAProxy and Kubernetes, which allows it to be seamlessly integrated with any application, not to mention the broader set of features that it inherently offers (e.g., SSL termination, SNI support, etc).

# Chapter 3

## Use case: smart grids

This chapter presents a use case and a sample application that follows the set of requirements described in Section 1.1. This use case and its implementation are used to better motivate the solution proposed in Chapter 4, and also to guide its evaluation in Chapter 5.

Consider the Electric Power Distribution Operator (EPDO) of a large city or region. The EPDO is responsible for distributing electricity from transmission lines to end customers through a distribution grid. End customers include urban, rural and also industrial facilities. The grid is equipped with smart meters and other sensors, which send metrics about the overall state of the network constantly. This data is used to detect failures, frauds and anomalies, to analyze the quality of the transmission, and to guide decision making, aiming at a more optimized system. It comprises important information about the distribution network, which is in itself sensitive, as it exposes points of failure, but also contains detailed metrics and consumption profiles of all end customers.

A large volume of data is processed by the system constantly, as the smart meters and sensors push their metrics. The service provides an interface to visualize and query the data. All communications are end-to-end encrypted, and unauthorized access must not be allowed. Data at rest is also protected through encryption. Only the owner of the data is able to access, visualize or query it. It means that the EPDO is not able to visualize end customers data, despite being the application provider. It only has access to metrics of the distribution and overall network state (pushed by sensors and smart meters that belong to the EPDO).

The system must be available at all times, and also provide a reasonable response time to the end customer, whilst guaranteeing the isolation, confidentiality and integrity of their respective private records.

## 3.1 Architecture

The system is composed by a server application and two types of client applications. One client application (Type 1) runs on smart meters and sensors, and its only function is to push new metrics to the server. The other client application (Type 2) allows the user to perform queries and visualize their data. All user data is saved to an embedded relational database, and each individual user has its own database, which is stored in a single file in the server file system. The database file is encrypted with an individual encryption key, which is provided to the server after a connection is successful established through a `LOGIN` operation (Section 3.2.1). All operations are performed through HTTPS requests to the server, and all communications between client and server are encrypted with TLS (Transport Layer Security).

## 3.2 Supported operations

The server application supports a set of operations that can be performed by the different client applications to transform and query the user data. All operations, except the `LOGIN` operation, are performed in the context of a *Session*. The `LOGIN` operation creates a *Session*, where cryptographic data is exchanged and the secure communication is established. The complete list of operations supported by the server application are listed below.

### 3.2.1 LOGIN

The `LOGIN` creates a user *Session* in the server, and must be performed by the client before any other operation, in order to authenticate the server and exchange cryptographic information to decrypt the user database. All further operations must be executed within a *Session* context. A *Session* has also a limited duration, after which it expires, requiring the user to perform another `LOGIN` operation and create a new *Session*.

After a secure mutual TLS connection is established, which means that both server and client authenticated each other, a *Session* is created, with a unique *authorization token*. After the client authenticates the server (e.g., with a TLS certificate that could only be provide after an initial SGX attestation), it shares its encryption key, which will be used by the server to decrypt the user records.

### 3.2.2 LOGOUT

A LOGOUT operation destroys the current session and closes all open connections and resources, committing the user data to the database. It can be performed by the user or automatically triggered after a *Session* expires. After a LOGOUT, the server is idle, i.e., not actively serving anyone.

### 3.2.3 WRITE

WRITE updates the user database with new information via an SQL INSERT statement. It is intended for Type 2 client actions (e.g., updating metadata, creating alarms and thresholds, or creating devices). In the context of the example application, it writes a few alarms to the `alarms` table. One WRITE operation might trigger multiple INSERT statements at a time, in which case the server application runs them inside of an SQL transaction, that is either commit, if all inserts are successful, or rolled back, if at least one of the inserts fail.

### 3.2.4 LONG\_WRITE

LONG\_WRITE also updates the user database with new information via an SQL INSERT statement. However, it is intended for Type 1 clients pushing new metrics collected from a smart meter or sensor. This actions are usually much larger, since the client accumulates a large number of measurements (e.g., the 900 one-second measurements for the last 15 minutes or 3600 from the last hour). For this reason, they are isolated in a separate type of operation. A LONG\_WRITE operation triggers dozens of INSERT SQL statements at a time, which are all executed by the server within an SQL transaction context.

### 3.2.5 QUERY\_AVG

The QUERY\_AVG operation allows the user to execute queries on their data through a SELECT SQL statement to calculate metric averages (e.g., consumed power, power factor, peak demand). The server processes such queries and return their results, which can then be used to plot visualization or feed analysis algorithms, for example. Only Type 2 clients are allowed to perform this operation.

## 3.3 Implementation

### 3.3.1 Server

The server application is implemented in C, which allows for an efficient memory management and overall small resource footprint. To manage the user database, the example application uses SQLite<sup>1</sup>, a library that implements a small, fast and highly reliable SQL database engine. The database lies in a single file, which makes SQLite the perfect choice for embedded or self-contained applications. In fact, SQLite is built-in in nearly all mobile phones today, making it arguably the most used database in the world<sup>2</sup>. In the example application, SQLite allows for a better performance, if compared to a full-fledged DBMS (database management system). Its file format also makes it easier to encrypt the whole database at once.

We use SCONE to shield the application execution. C applications provide the smallest enclaves when running on SCONE. The enclave size of SCONE applications does not vary at runtime, and the entire enclave size (defined by environment variable *SCONE\_HEAP*) is allocated at startup time, which increase cold-start times. Also, having the smallest possible enclave size helps with overall performance, even when the startup time overhead is not considered. Larger enclaves, due to the EPC memory constraint, result in more swapping from the protected memory to the main memory, which hurts performance severely. Go, for example, has an absolute minimum enclave size of 152 MB in SCONE. As the application grows and becomes more complex and adding dependencies, this minimum requirement can reach up to a few gigabytes, which would make it prohibitive to deploy many instances. Not only the swapping from EPC to main memory would be expensive, but also the usage of main memory is not efficient, especially when considering deploying a few dozen servers per node. In most servers, this would result in memory pressure. Table 3.1 shows the absolute minimum enclave size requirements in SCONE for C, Go and Python applications, alongside their average startup time. The test application is empty and returns immediately. Please note that interpreted Python code needs to be encrypted in order to provide the same integrity and confidentiality guarantees as C or Go applications.

The C server exposes an HTTP REST API, and uses TLS (Transport Layer Security) to encrypt the messages end-to-end. The server API has distinct endpoints for each operation. Since each server is supposed to serve only one client (due to the isolation requirements), the server certificate is issued to the client UID name (i.e., the server certificate's *Common Name* field will be the client UID).

---

<sup>1</sup><https://www.sqlite.org/index.html>

<sup>2</sup><https://www.sqlite.org/mostdeployed.html>

Language	Minimum enclave size	Avg. startup time
C	8 MB	0.116 second
Go	152 MB	62 seconds
Python	8 MB	0.23 second

Table 3.1: Minimum enclave size requirements in SCONE.

The SCONE configuration of the server uses only 1 queue, 1 e-thread and 1 s-thread, both unpinned. The enclave size (*SCONE\_HEAP*) is 64 MB.

### 3.3.2 Client

The client application is written in Go, and does not run inside of Intel SGX enclaves. It can simulate clients of Type 1 or Type 2, and supports different execution parameters that allow for different client workloads. The list of parameters supported by the client application is shown in Table 3.2.

Parameter	Description
<code>-server</code>	Server address
<code>-port</code>	Server port
<code>-sni</code>	Client UID, which is also the server name (as in the server certificate)
<code>-cert</code>	Path to client certificate to establish a TLS connection
<code>-key</code>	Path to client private key to establish a TLS connection
<code>-ca</code>	Path to client CA certificate to authenticate the server identity
<code>-keydir</code>	Path to the key used to decrypt user's records
<code>-mode</code>	Whether the client is simulating clients of Type 1 or Type 2
<code>-duration</code>	How long the session will be. Long sessions might require another LOGIN
<code>-frequency</code>	How often requests will be sent within a session
<code>-writeRatio</code>	The proportion of WRITE operations to be performed, on average
<code>-csv</code>	Print output in CSV format

Table 3.2: List of execution parameters supported by the client application.

# Chapter 4

## Dynamic provisioning of stateful applications in TEEs

Provisioning containerized applications with a strict set of privacy and isolation requirements, as posed in Section 1.1, is already a non-trivial problem. When such applications run on a third-party infrastructure, such as public cloud providers, it raises the complexity of any acceptable solution. In fact, when running applications in a public cloud infrastructure, the addition of TEEs might cover for the various privacy concerns, especially when assuming a threat model that does not trust the cloud itself. However, the addition of TEEs also brings some other concerns, especially related to performance (which, in a public cloud setting, incurs in more financial cost to the application owner).

This section presents an approach that aims at solving such concerns, by providing an architecture in which such strongly-isolated applications (i.e., each application instance is a separate, independent process) can be deployed to serve a large number of users. Two alternative approaches are also presented for comparison.

The three approaches are named *Trivial*, which is the simplest, *On-demand* and *Standby Servers*. The later will be used throughout the rest of this work, and compared against the trivial approach.

### 4.1 Alternative approaches

#### 4.1.1 The trivial approach

If the applications have strong isolation requirements (i.e., one application instance serves one user), as described in Section 1.1, the trivial approach would be to have all application instances running at



all times, ready to serve the user whose state they hold. This is the simplest approach, and in fact, probably the most responsive too (from the user perspective), since no additional work is needed to provision the requested state.

However, the trivial approach is not resource efficient, which translates into a much higher operational cost. At any time, all the states are available and ready, but only a fraction of them is being actively requested. For instance, if a certain service has 100 users, but the average user activity is 25% (i.e., 25 active users), the application would have 75% of the instances ready and running, but idle. This is clearly not efficient and, in fact, prohibitive when the user base start to grow over millions of users.

### 4.1.2 Dynamic provisioning: the on-demand approach

A second approach, which resembles serverless architectures, is to spawn application instances only when there is an active client. At any time, the number of running application instances is equal to the number of active users, and this number is dynamic. This approach introduces two new components to the architecture:

- **Routing Manager:** some sort of proxy responsible for processing the incoming client requests and deciding whether a new application instance must be started or not.
- **Application Wrapper:** the component that would actually start the application instance when requested by the Routing Manager. This component might or might not be under the application owner's responsibility. In a serverless environment, for example, it would be provided by the cloud provider itself (e.g., the serverless API).

This dynamic provisioning approach is more resource-efficient, since only the needed amount of resources is being used at any given time. However, it might come at the cost of performance and overall complexity of the solution. Now there are a few additional processes between the user request and the server response. For instance, when a user request a new state:

1. The user requests a new state (i.e. a new application must be spawned);
2. The Routing Manager decides that a new instance is needed;
3. The Routing Manager talks to the Application Wrapper to spawn a new instance;
4. The Application Wrapper spawns a new instance;

5. The requested state is provisioned, either by the Application Wrapper, in case of having some sort of dynamic data provisioning, or by the application instance itself, in case the state is pulled from an external entity (e.g., a database or a content delivery network);
6. The application instance receives the user request;
7. The application instance responds to the user request.

Some of these processes might happen over the network (e.g., 3 and 5), or depend on external entities (e.g., 5). Also, processes 4 and 5 are particularly more expensive if strict privacy requirements are in place. Starting applications that run inside of TEEs is way more expensive than running native applications. Also, if the state being provisioned is encrypted, an additional overhead is imposed to the application instance, and overall response times. All of these extra steps (2–5) combined point to a terrible user experience and prohibitively high tail-latency times.

If we consider the scalability aspect, the approach should scale well as long as the Routing Manager is able to consider multiple nodes to which it can route requests or request new application instances to be deployed.

## 4.2 Dynamic provisioning: the standby servers approach

The previous approaches have advantages and disadvantages. For instance, the Trivial approach, although the simplest of the three, does not provide a good solution from resource consumption and scalability points of view. The on-demand covers for the resource-efficiency bit, which is important when the infrastructure is run by a third-party. However, starting application instances on-demand also imposes a severe performance overhead, which may be prohibitive for the user experience and quality of service.

The approach presented below is very similar to the on-demand approach, having the same architecture and components. However, it tries to mitigate the overhead imposed by starting the application on demand within a TEE. It does that by keeping a number of application instances running, but without any state loaded. Such application instances are ready to serve any user. Whenever a new user arrives, the Routing Manager forwards their requests to one of the idle application instances. Such instance then assumes the right identity and responds to the user, who it will continue to serve until the end of the session.

In this approach, an idle application instance assumes the correct identity by retrieving the requested state from an external source after it performs a LOGIN operation. The external source can

be, as already described in the previous approach, a database, another service or website, a local file system, or even a dynamically-provisioned file system or volume.

The amount of application instances that run is fixed and should be calibrated based on the average load of the service. If all the idle application instances are already serving someone, the amount of servers can be increased, leading the Application Wrapper to spawn more idle application instances to accommodate the extra load.

This architecture also supports multiple nodes, as long as the Routing Manager is able to coordinate between all the Application Wrappers. Thus, assuming a fair ratio of application instances deployed by this approach, the latency increase caused by instantiating a new application within a TEE is attenuated.

## 4.3 Integration with example application

The example application (as described in Chapter 3) is integrated with the *standby servers approach* and the *trivial* approach (to provide a baseline performance for comparison). Both approaches (including all needed components) are defined in Kubernetes manifests, allowing them to be quickly deployed in any Kubernetes cluster with Intel SGX support. The following subsection discuss implementation aspects of the components.

### 4.3.1 Application Wrapper

The Application Wrapper is responsible for managing the life cycle of the server applications, and for reporting their status to the Routing Manager. It is implemented in Python and does not run inside of Intel SGX enclaves. Each application server is managed by a sidecar thread created by the Wrapper. The Wrapper is application-agnostic, and the communication between application and Wrapper is performed by the sidecar thread, which listens to specific events logged by the application, taking appropriate action upon them. Thus, the sidecar thread is the interface between Wrapper and application, and must be adapted to listen to the right set of events. In this example, the sidecar thread listens to events logged to the application standard output (Table 4.1).

The Application Wrapper has two operation modes: *standalone* and *dynamic*. If in *standalone* mode, the Application Wrapper has the needed parameters (such as the Routing Manager address, the number of servers to run, the server executable binary, whether to use dynamically provisioned volumes or not) in its environment, and starts all sidecar threads and applications right away. In *dynamic mode*, the Application Wrapper is remotely started or stopped by REST API calls, allowing

Server event	Wrapper action
WAITING_LOGIN	Report the instance status as <i>Ready</i> (i.e., accepting new connections) to the Routing Manager.
LOGIN <ClientUID>	Report the instance status as <i>Active</i> (i.e., currently serving someone) to the Routing Manager. Provision the encrypted user data for <ClientUID>.
LOGOUT	Report the instance status as <i>Stopped</i> to the Routing Manager. Deprovision the encrypted user data. Close all open connections and resources. Restart the server application.

Table 4.1: Server events and Application Wrapper actions.

it to be integrated to an external service, such as an auto-scaling controller. In *standalone mode*, the number of servers is fixed. In this work, we use the Application Wrapper in *standalone mode*, which means that the number of applications is fixed. Although using dynamic mode would provide a more flexible and production-like behavior (e.g., number of servers varying dynamically), it would also bring more complexity to the architecture and introduce the need for another component (e.g., an auto-scaling controller that manages the Wrappers via their API), thus being considered out of the scope of this work.

After starting, the Wrapper contacts the Routing Manager and advertises itself, sending its current status and how many applications it manages. Then, each sidecar thread also contacts the Routing Manager to advertise the status of the application it manages (i.e., *Ready*, *Active* or *Stopped*). Each status change in the application is also reported by its sidecar thread to the Routing Manager. Besides that, the Wrapper also reports periodically the state of all applications.

Wrappers are deployed to Kubernetes clusters as Deployment resources. The servers run inside of the Wrapper container, each one listening to a different port. Thus, although the Wrapper is seen as a single unit from the Kubernetes API perspective, it also contains all the applications that it manages running inside of the same container. Running each application as a process in the same container is more efficient, since spawning a Pod per application would take an additional round trip to the Kubernetes API, additional resources, and a much more noticeable overhead (from performance and management complexity perspectives). Each application is then exposed to the cluster through a Service, acquiring a local IP address and an entry in the in-cluster DNS server.

### 4.3.2 Routing Manager

The Routing Manager is responsible for the dynamic assignment of server instances to incoming clients. It holds the state of each application instance and updates it according to the reports sent to it by the Wrappers periodically. The Routing Manager exposes a Report API, used by the Wrapper and the sidecar threads, and a Query API, which is used by the load balancers in order to make a routing decision. Whenever a new client connects to the load balancer, the Routing Manager is queried and returns an unassigned application instance, to which the client request is forwarded. Further requests from the same client within the same Session can either have the load balancer query the Routing Manager again, or just bypass it via some sort of persistent session (e.g., session cookies or the *Keep-Alive* HTTP header).

The Routing Manager is deployed to Kubernetes as a Deployment resource, and is exposed by its own Service resource, which allows Wrappers, sidecar threads and load balancer to reach it via its local in-cluster DNS record.

### 4.3.3 Load balancer

The load balancer is the outermost component of the architecture and a key point of the end-to-end communication. It exposes a single endpoint through which the clients can contact the application instances transparently. It runs replicated, balancing and routing the incoming traffic based on queries made to the Routing Manager, which holds the state of the system. In this approach, the load balancer replicas are not supposed to run inside of Intel SGX enclaves. Therefore, they are not able to inspect the requests or any sensitive information: they should only route the incoming request to the appropriate server. The routing is possible through the SNI (Server Name Indication) TLS extension, which is set by the client. The SNI extension indicates the target server of a TLS request in the handshake process, allowing the load balancer to query the Routing Manager for the appropriate server. The content and the headers of the request are only accessible when the TLS is terminated, which only happens inside the application server (with the appropriate certificate and inside of Intel SGX enclaves).

The load balancer is deployed to Kubernetes as an Ingress resource (as discussed in Section 2.2), which exposes each Service for application servers. Since Kubernetes does not provide a default implementation for Ingress controller, it is required to pick one third-party controller among the many currently available. A preliminary investigation was conducted in order to define the most suitable Ingress controller solution.

## Picking an Ingress controller

The official Kubernetes documentation page for Ingress controllers lists at least 15 different options for Ingress controllers, each one with its own target audience and feature set. The most popular load balancers and proxy solutions on the market have their own Ingress controller implementation. In order to find the most suitable Ingress controller implementation, the following requirements were defined.

1. Support the TLS SNI extension. Required for the routing in end-to-end encrypted channels, so the load balancer does not need to terminate the TLS connection.
2. Support some sort of mechanism that allow the load balancer to be integrated with the Routing Manager.

The following three implementations met these requirements, and were further compared.

- **Contour.** A Cloud Native Computing Foundation incubating project, Contour is an Ingress controller based on Envoy Proxy, an open-source layer 7 (according to the OSI model) proxy created by Lyft. Envoy is a very popular, production-ready cloud-native proxy written in Go. It supports SNI routing natively. Even though Contour does not provide an easy way to integrate with external services (the feature that would allow this use case, called *External Authorization*, is currently under development), it is possible to leverage its dynamic route management to mimic that. In Contour, each route (i.e., the binding to a server application) is a separate Kubernetes CRD (Custom Resource Definition), which allows for more flexibility. The Routing Manager, however, would need to be extended in order to talk to the Kubernetes API and create new routes dynamically.
- **NGINX Ingress Controller.** Based on NGINX proxy, one of the most famous proxy solutions on the internet, created and maintained by the community. It supports SNI-based routing natively, and its Lua extension support allows for a straightforward integration with the Routing Manager through sockets.
- **Voyager.** Based on HAProxy, which is another famous load balancing solution, Voyager is an Ingress controller created and maintained by Appscore. It supports SNI routing natively and Lua extensions, similar to the NGINX alternative, which allow for socket use to contact the Routing Manager.

To compare the three Ingress controller solutions, a preliminary experiment was conducted. A simple web server application, written in C and built with SCONE, was deployed to a 10-node Kubernetes cluster. Each Ingress controller solution was deployed to manage 800 backend servers. A separate client machine ran one or two clients per backend server, with TLS and SNI, for a 10-minute benchmark consisting on HTTPS requests. The load balancers (i.e., the actual proxy Pods) were replicated (10 replicas, one per node). The output variables are **latency**, in milliseconds, **throughput**, in requests per second, and total requests. Metrics for CPU and RAM usage of the load balancers were also collected.

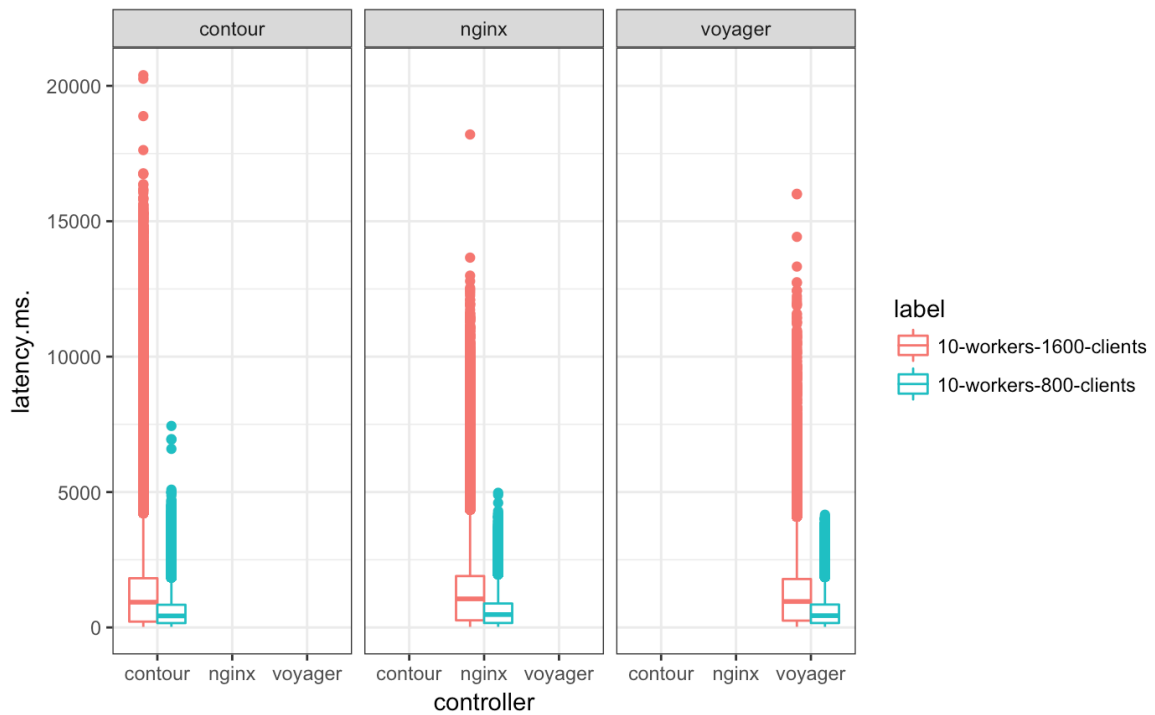


Figure 4.1: Ingress controllers: latency in a 10-minute benchmark.

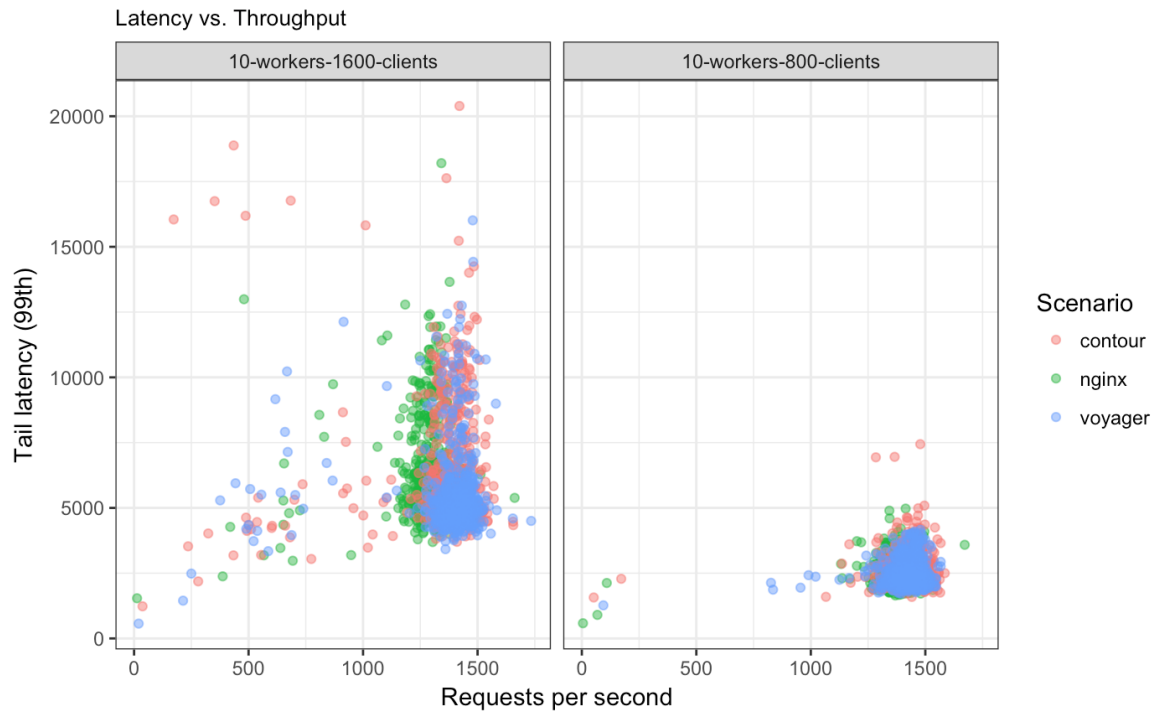


Figure 4.2: Ingress controllers: tail latency (p99) versus throughput in a 10-minute benchmark.

Voyager had a slightly better performance, presenting lower latencies and also smaller performance drops when doubling the benchmark clients (Figure 4.1). Figure 4.2 also puts Voyager in the lead when comparing tail latency and throughput. Finally, Table 4.2 confirms the slightly more stable behavior of Voyager, which also presented the smallest resource footprint of the three Ingress controllers.

Ingress controller	Avg. total requests, 800 clients	Avg. total requests, 1600 clients
Contour	859208	816512
NGINX	830716	771935
Voyager	854320	831612

Table 4.2: Ingress controllers: total requests processed in a 10-minute benchmark.

For this reason, Voyager was picked as the Ingress controller for this work. In addition, HAProxy, on which Voyager is based, is well-known for its reliability and efficiency, being used in an enormous number of production systems.



In Voyager, the integration with the Routing Manager is achieved by a simple Lua action that is triggered upon every incoming TCP connection. It extracts the SNI information and sends a request to the Routing Manager, which responds with the name of the appropriate backend. The socket operation performed by Lua is non-blocking, which attenuates the performance impact of an HTTP request to an external service.

#### 4.3.4 Data encryption and provisioning

The user data is provisioned through an auxiliary script, that leverages SCONE CLI to encrypt all the user data transparently (i.e., the server application does not have to deal with cryptography explicitly). The input to the script is the plain data, which is an empty SQLite database. The script outputs the encrypted data directory and the respective keys. The keys are distributed to clients, while the encrypted data is distributed to all of the cluster nodes. Servers access the encrypted user data via a local mount of the node file system. Ideally, the data would be provisioned by a dynamic volume provisioner, such as Ceph RBD. In order to simulate the delay imposed by a dynamic volume provisioner, the servers impose a time penalty of 100 milliseconds when opening the user data during a LOGIN operation.

# Chapter 5

## Performance evaluation

### 5.1 Experiment description

The proposed architecture was evaluated with respect to its efficiency and performance in a simulated environment, managing server instances of the example application described in Chapter 3. The proposed approach is compared to its trivial counterpart, where all server instances are running at all times.

The performance with the *Trivial* approaches defines the baseline for our test infrastructure. The proposed approach is evaluated in different scenarios. Each scenario is composed by variables of a fractional factorial design experiment, and correspond to a distinct user activity pattern. The variables are listed below.

- **Active User Ratio (AUR).** The maximum amount of users that the application is expected to have active at any given time. Since the *Standby Servers* approach allows for the dynamic serving of users through a fixed set of servers, it is necessary to define the target user ratio. The total amount of users is defined by the baseline experiment, which uses the *Trivial* (or all servers running at all times) approach. The levels for this variable are: 50% and 25%.
- **Write Requests Ratio (WR).** Defines the write/total request ratio. For instance, if WR is 50%, half of the client requests, on average, will be WRITE operations. The levels for this variable are 50% and 10%.
- **Request Frequency (RF).** Defines how often the client will send requests within a session. The levels for this variable are: 10 seconds and 60 seconds.

The output variables of this experiment are define below, and allow to evaluate the server performance from the user perspective (latency), and the resource consumption from the server perspective.

- **Latency** (client). Time spent by the client in each request, in milliseconds. Clients start the tracking of time right before sending the request, and stop immediately after a server response arrives. Therefore, request preparation and response inspection are not included in this time.
- **CPU utilization** (server). CPU time ratio for the server instances, per node, as reported by the Docker engine through cAdvisor. Each node reports a ratio up to 1 (e.g., 100%), which comprises all available cores (e.g., two cores report up to 1 instead of 2).
- **Total RAM memory usage** (server). Total memory allocated to the server instances, as reported by the Docker engine through cAdvisor. This metric is the sum of the memory allocated for each node, expressed in GB.

The duration of the user session is fixed at 600 seconds, or 10 minutes. All sessions are started with a LOGIN operation, and are closed by a LOGOUT operation.

## 5.2 Experiment environment

The servers (along with routing managers, load balancers and application wrappers) are deployed to an 8-node Kubernetes cluster. Each cluster node is a virtual machine (VM) in a private cloud environment managed by OpenStack. The cluster has 1 master node, which is used exclusively to run the cluster control plane and Kubernetes-related workloads (such as the API server). The master node VM is of the type `sgx-scalable.epc0.medium`, with 2 vCPUS, 4 GB of RAM, 60 GB of storage and no EPC memory (i.e., it is not able to run Intel SGX workloads). The rest of the cluster is composed by 7 similar worker nodes, whose VM type is `sgx-scalable.epc30.large`. Each worker node has 4 vCPUs, 8 GB of RAM, 100 GB of storage and 30 MB of EPC memory (i.e., able to run Intel SGX workloads). The 7 worker nodes are distributed among 3 physical host machines, which means that co-located VMs probably affect one another at some degree. This cloud environment does over commit CPU (i.e., the sum of virtualized allocatable resources exceeds the real amount of resources in the physical host), but does not over commit RAM nor EPC.

The clients run on a single physical machine, with a 3.1 GHz dual-core Intel Core i7 5557U, 16 GB of RAM, 512 GB of storage and no EPC memory. Each client is spawned on its own Docker container. The client machine is not in the same physical network as the server VMs, as it is expected in this use case, where clients access services hosted in a remote cloud.

## 5.3 Experiment tools

To collect the latency of the servers in different scenarios, the client application (as described in Chapter 3) is modified to track and log the time spent in every request, until the server response arrives. Request preparation and response inspection are not included in this time. The client applications use the standard Go time library (`time`) to track time. The approach used to deploy the server instances is completely transparent to the client application, which always talk to the same endpoint (i.e., the load balancer endpoint). An auxiliary *bash* script spawns parallel clients within a given scenario and collects their results and logs in a CSV file.

To evaluate the resource consumption of the cluster, we use TEEMon, a continuous performance monitoring framework for TEEs. TEEMon supports Kubernetes clusters and provide a series of dashboards where it is possible to monitor not only TEE state (e.g. active enclaves, active EPC pages, EPC page fault occurrences) but also the infrastructure (CPU, RAM, network) or specific Docker containers. Metrics are stored in Prometheus, a popular monitoring system and time-series database, and visualized in Grafana, a popular observability platform. Our metrics of interest, on the server side, come from two exporters deployed by TEEMon: cAdvisor (v0.30.4), which exposes container metrics (e.g., CPU, memory and network); and SGX-exporter, which provides EPC-related metrics (e.g., free pages, allocated pages, evicted pages). Metrics are scraped from exporters every 15 seconds.

## 5.4 Baseline performance

The baseline performance is defined by the *Trivial* approach, since we want to define the maximum amount of servers that can run comfortably in the test infrastructure, whilst still providing a reasonable tail latency for all operations. The preliminary experiment aims at defining a reasonable amount of servers for our experiment environment, which will then be considered as the total amount of users. Table 5.1 shows the baseline experiment levels, which configure the most intensive usage scenario: a write-intensive (WR is 50%) set of frequent operations (every 10 seconds).

The output variable is *Latency*, or the total time spent by each client request to be completed, in milliseconds. We are specially interested in the tail latency, here defined as the 95<sup>th</sup> and the 99<sup>th</sup> percentiles of *Latency*. The  $n^{\text{th}}$  percentile of the *Latency* distribution is the maximum latency for the fastest  $n\%$  of all requests. For instance, if the 99<sup>th</sup> percentile latency of a given operation is 500 milliseconds, it means that 99% of the requests were processed under 500 milliseconds. The median, which is the 50<sup>th</sup> percentile, is also shown.

<b>Servers per wrapper</b>	<b>Total servers</b>	<b>Write ratio</b>	<b>Request frequency</b>	<b>Session</b>
15	105	50%	10 seconds	600 seconds
20	140	50%	10 seconds	600 seconds
25	175	50%	10 seconds	600 seconds

Table 5.1: Levels for baseline experiment with *Trivial* approach.

<b>Servers</b>	<b>LOGIN p50 latency</b>	<b>LOGIN p95 latency</b>	<b>LOGIN p99 latency</b>
105	676.812 ms	1012.090 ms	1131.732 ms
140	736.101 ms	1141.018 ms	2105.009 ms
175	755.821 ms	1824.261 ms	2719.989 ms

Table 5.2: Latency percentiles for LOGIN operations.

<b>Servers</b>	<b>LOGOUT p50 latency</b>	<b>LOGOUT p95 latency</b>	<b>LOGOUT p99 latency</b>
105	365.972 ms	740.487 ms	1081.462 ms
140	443.303 ms	905.167 ms	1427.415 ms
175	591.772 ms	3079.142 ms	5307.941 ms

Table 5.3: Latency percentiles for LOGOUT operations.

<b>Servers</b>	<b>QUERY_AVG p50 latency</b>	<b>QUERY_AVG p95 latency</b>	<b>QUERY_AVG p99 latency</b>
105	173.064 ms	489.750 ms	1211.930 ms
140	185.184 ms	731.192 ms	1299.095 ms
175	186.661 ms	819.088 ms	1299.523 ms

Table 5.4: Latency percentiles for QUERY\_AVG operations.

Servers	WRITE p50 latency	WRITE p95 latency	WRITE p99 latency
105	344.935 ms	706.655 ms	1386.174 ms
140	381.992 ms	975.072 ms	1565.018 ms
175	422.520 ms	1163.695 ms	1824.909 ms

Table 5.5: Latency percentiles for WRITE operations.

Figure 5.1 shows the performance degradation caused by the increase in the number of active application instances. More active applications result in more clients being imposed prohibitive latency times. Figure 5.2 shows the *Latency* distribution, which lets us analyze the tail latency for each setting.

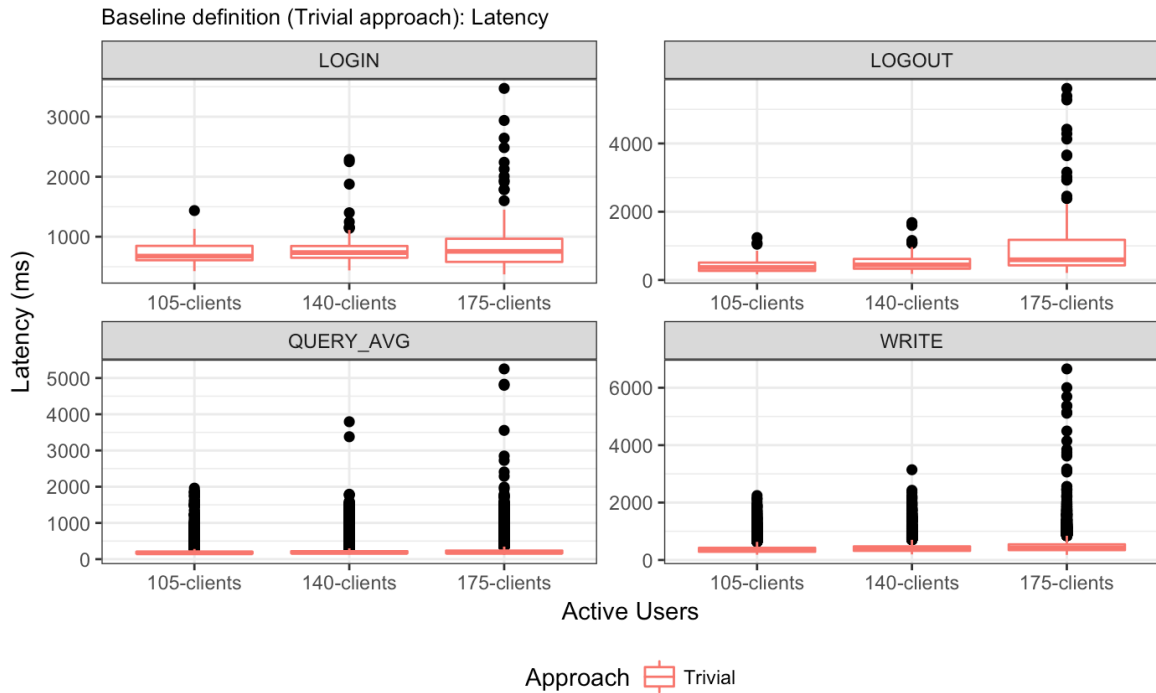


Figure 5.1: Baseline: per-operation latency (Trivial approach).

As seen in Figure 5.2, the performance of 15 servers per wrapper presents the best latency curve for all operations, being the most reasonable when it comes to the test infrastructure. Therefore, we consider 105 servers to be the total amount of servers (and users) in this experiment.

Resource-wise, we have a similar landscape. Since SCONE enclaves are statically allocated at startup time, we expect memory usage to be constant for a given amount of running servers. The

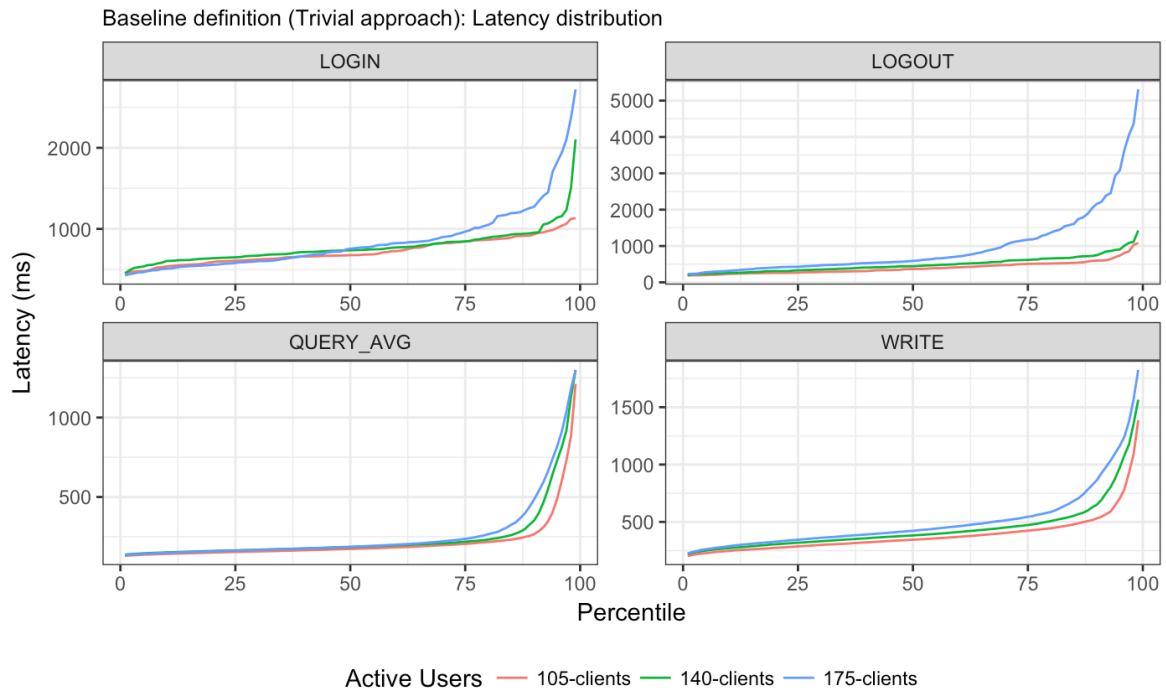


Figure 5.2: Baseline: per-operation latency distribution (Trivial approach).

application instances inside enclaves are not allowed to allocate more memory. The memory usage is shown in Table 5.6. CPU utilization graphs (Figure 5.3) show a clear increase in CPU time with more running servers, which is also expected.

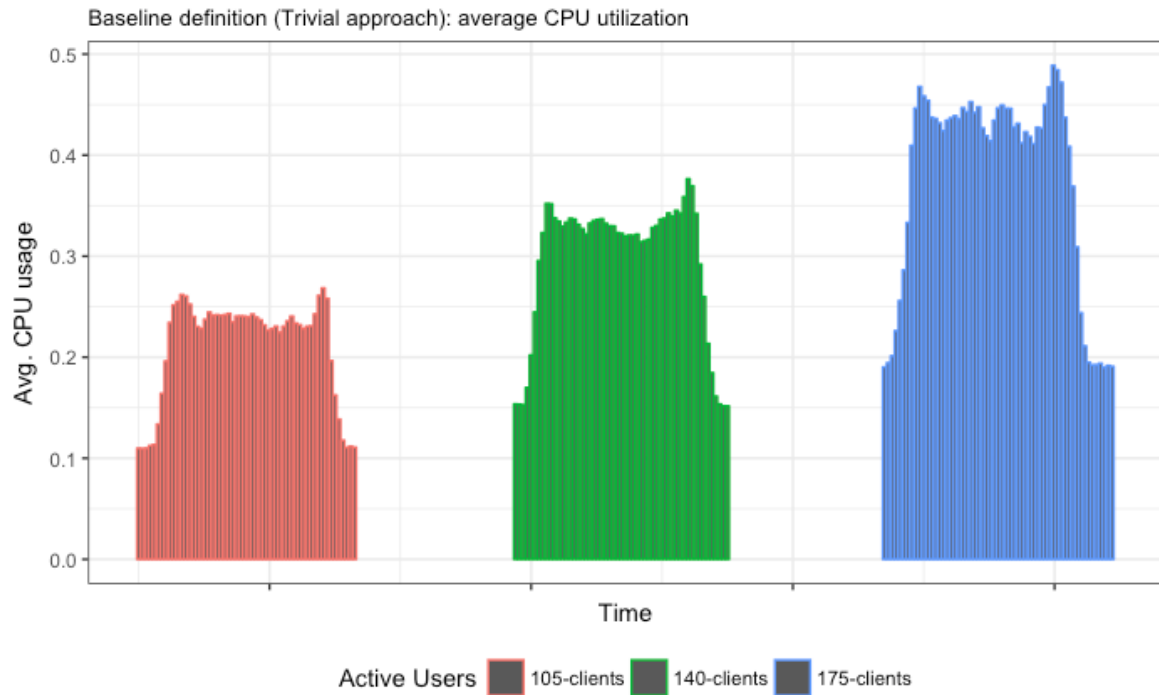


Figure 5.3: Baseline: average CPU utilization (Trivial approach).

Servers	Allocated memory
105	10.99 GB
140	14.61 GB
175	18.26 GB

Table 5.6: Baseline: total memory allocated to application instances in the cluster.

## 5.5 Discussion

We compare the per-operation latency distribution of *Trivial* and *Standby Servers* approaches for an Active User Ratio of 50%. Thus, we have 56 active users. The session duration is 600 seconds, and the request frequency is 10 seconds. For this comparison, we use a Write Request ratio of 50%. We only consider Type 2 clients.



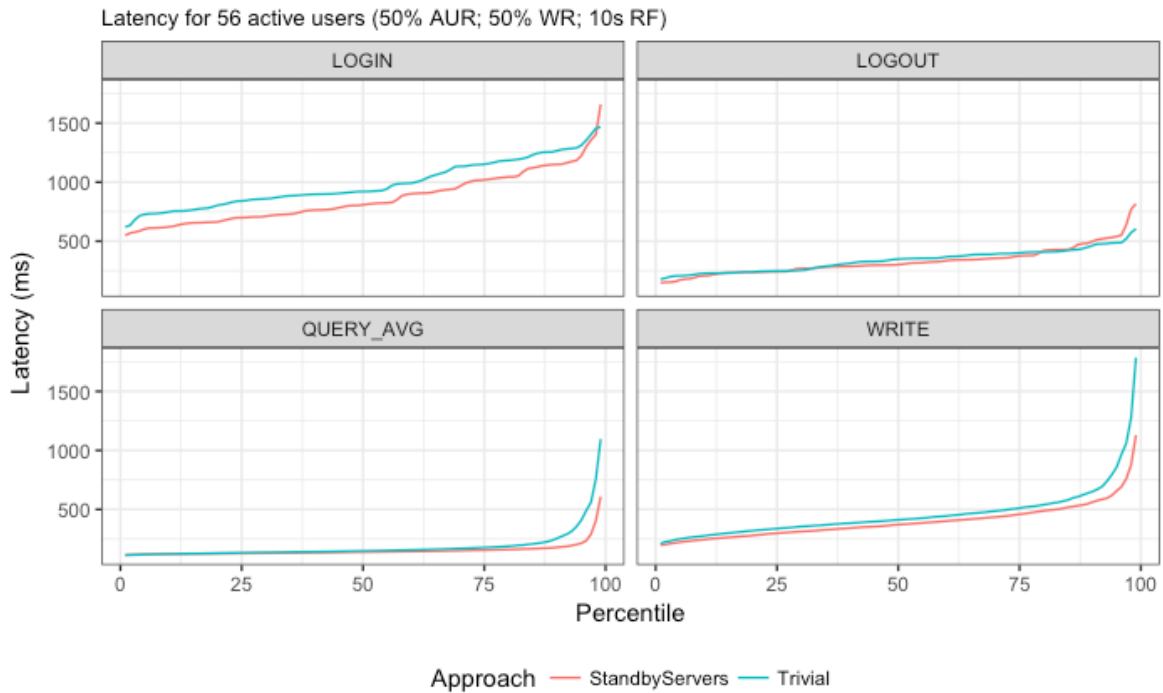


Figure 5.4: Latency distribution for 56 clients, or 50% Active User Ratio (RF 10 seconds; WR 50%).

In this scenario (Figure 5.4, the client sends requests every 10 seconds during the whole session duration of 10 minutes. Half of the requests, on average, are `WRITE` operations. There are 56 active users (out of 105), or an Active User Ratio of 53.33%. It means that, in the Trivial approach, almost half of the allocated resources are not used. This extra overhead can be seen in the latency distribution, as the Standby Servers approach is faster. If we consider an access pattern that has fewer, more sparse requests, the behavior is similar. Figure 5.5 shows the comparison when the Request Frequency is 60 seconds (i.e., client sends requests every 60 seconds).

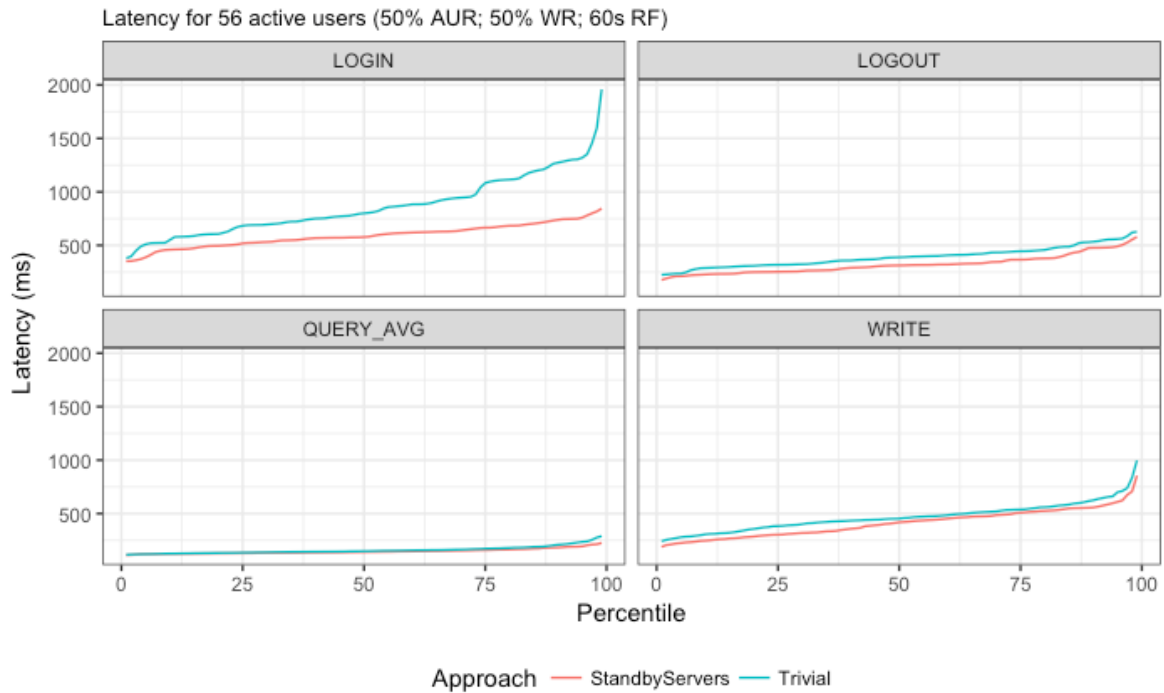


Figure 5.5: Latency distribution for 56 clients, or 50% Active User Ratio (RF 60 seconds; WR 50%).

Considering another type of scenario (frequent requests, read-intensive), the results are similar. The Request Frequency is 10 seconds, and the Write Request Ratio is 10%, indicating a read-intensive workload. Figure 5.6 shows that the *Standby Servers* approach is still more efficient than the *Trivial* counterpart.

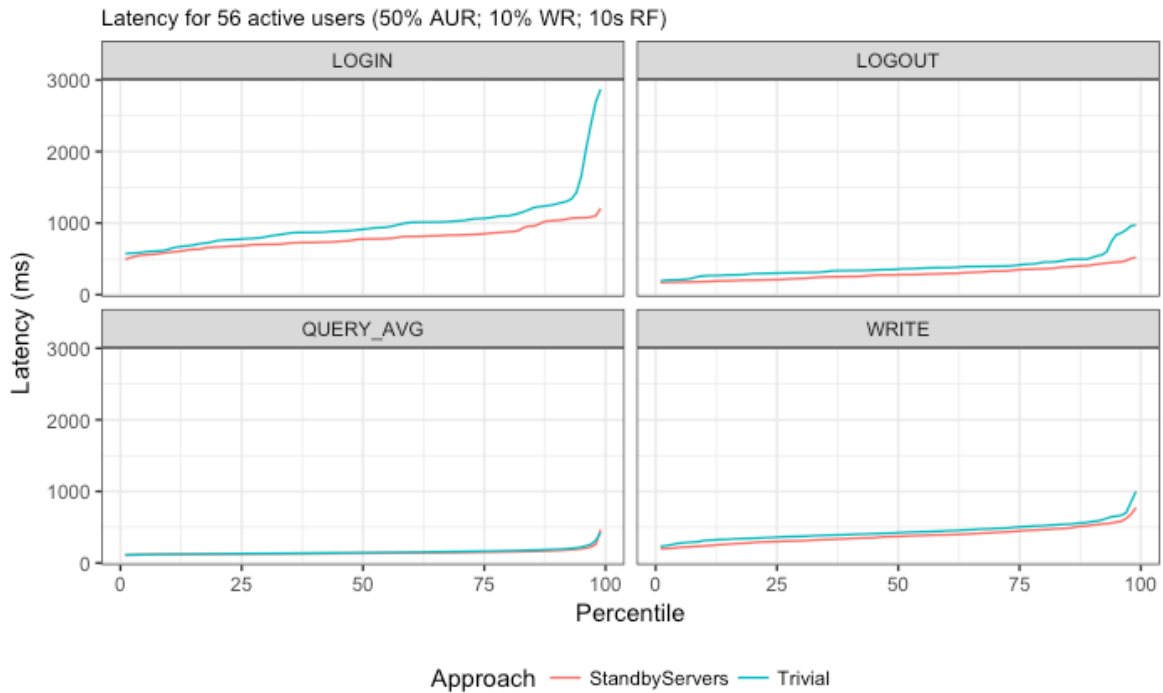


Figure 5.6: Latency distribution for 56 clients, or 50% Active User Ratio (RF 10 seconds; WR 10%).

When running the experiment for 28 active users, or 26.67%, the behavior is similar. `WRITE` and `QUERY_AVG` operations have both very similar latencies in both approaches. For `LOGIN` and `LOGOUT` operations, however, the *Standby Servers* approach shows better performance.

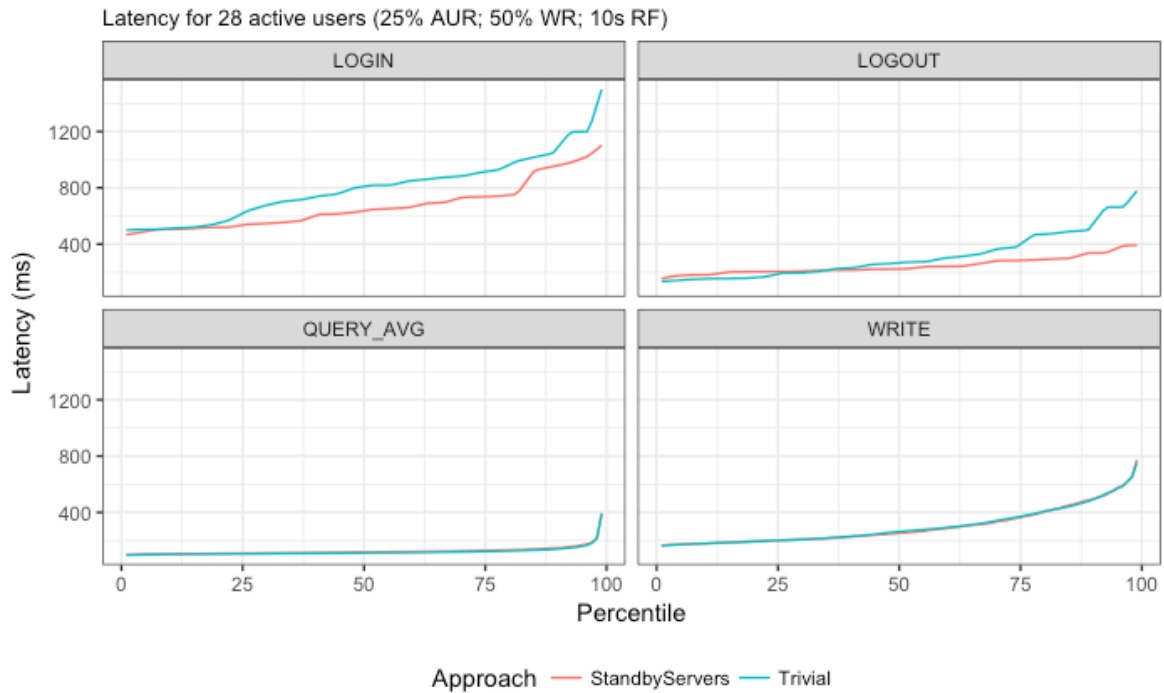


Figure 5.7: Latency distribution for 28 clients, or 25% Active User Ratio (RF 10 seconds; WR 50%).

Figure 5.8 compares the behavior of both approaches (w.r.t. to *Latency* per operation) when the Active User Ratio increases. The *Standby Servers* approaches performs equally or better than the *Trivial* approach whilst consuming a fraction of the resources. The difference is more noticeable in the most expensive operation, LOGIN.

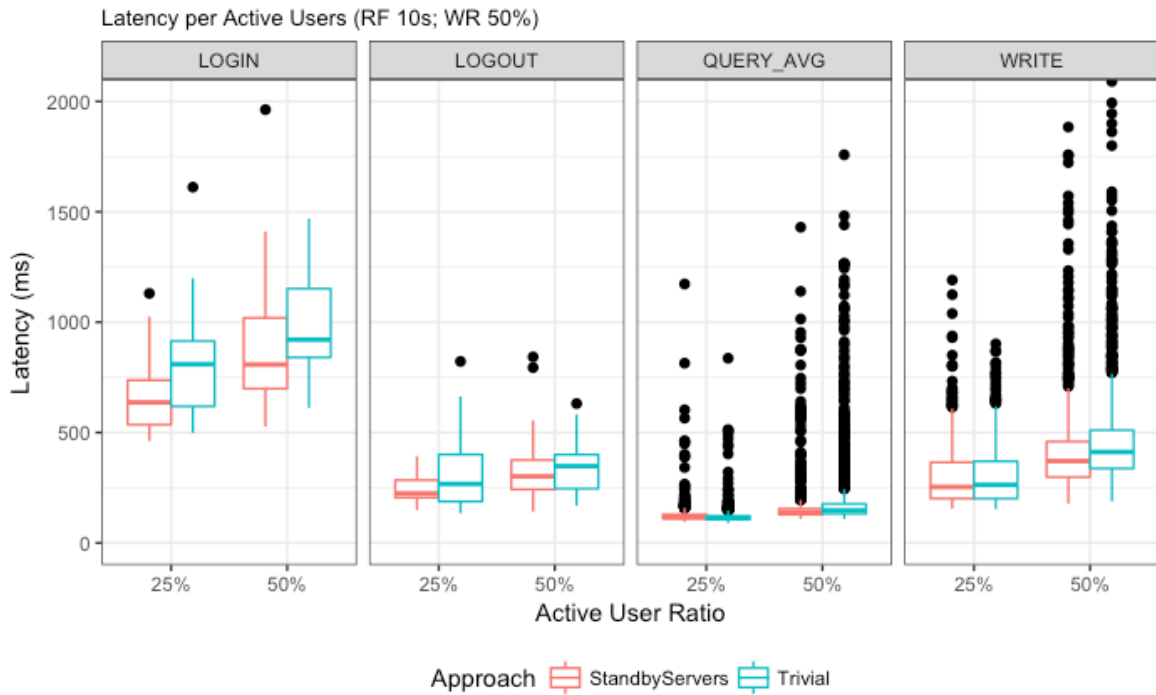


Figure 5.8: Per-operation latency versus Active User Ratio (RF 10 seconds; WR 50%).

Since the enclave sizes are fixed, having them managed dynamically to serve only active users is more resource-efficient. Table 5.7 compares the total memory utilization for the server instances in the cluster, as reported by the Docker engine through cAdvisor. Total memory of the cluster is 56 GB (8 GB x 7 nodes). The value shown under the *Memory utilization* column is the sum of the memory allocated by server instances in each node. The reduction in memory utilization provided by the *Standby Servers* approach is of 46.13% for an Active User Ratio of approximately 50%, and of 72,52% for an Active User Ratio of approximately 25%.

Active Users	Active User Ratio	Approach	Memory utilization (out of 56 GB)
56	53.33%	<i>Trivial</i>	10.99 GB
56	53.33%	<i>Standby Servers</i>	5.92 GB
28	26.67%	<i>Trivial</i>	10.99 GB
28	26.67%	<i>Standby Servers</i>	3.02 GB

Table 5.7: Total memory allocated to application instances in the cluster.

Based on the observed latency distribution and average resource consumption profile, the *Standby*

*Server* approach presents itself as viable alternative to deploy applications that fit into the architecture, and whose requirements are satisfied by it. Its performance is similar or better than the alternative *Trivial* approach in almost all tested scenarios, whilst using less computing resources. Nonetheless, all experiments were run on Kubernetes, the industry-standard container orchestration tool, which allowed a much simpler and more straightforward infrastructure management. Other popular tools, such as Prometheus and Grafana were also taken advantage of in the task of monitoring resource consumption and even TEE metrics, thanks to the TEEMon framework. The uniform API provided by Kubernetes also ensures that the architecture can be deployed to any other Kubernetes cluster with Intel SGX support.

## 5.6 Threats to validity

The work has the following threats to validity.

- The experiments do not consider important aspects of distributed systems, such as network capacity and storage performance.
- The experiments do not consider extra SCONE tuning parameters, such as number of queues, e-threads and s-threads, which can yield significant performance improvements.
- The experiments consider only one application.

# Chapter 6

## Final remarks

### 6.1 Conclusion

This work presented an approach to deploy applications with a strict set of confidentiality, integrity and isolation requirements. These requirements include end-to-end communication encryption and data-at-rest encryption. The *Standby Servers* approach allow the dynamic provisioning of such application instances, that run inside of Intel SGX enclaves, aiming at attenuating some of the overheads imposed by running inside of a TEE. The approach leverages a stateful routing mechanism powered by off-the-shelf load balancer solutions and a Routing Manager to dynamically assign incoming users to server instances that are ready and running, but do not hold any particular state. With the help of an Application Wrapper, which manages the application instances, the (encrypted) state is dynamically loaded into the applications' file system. After the client and server exchange establish a secure communication channel and the data is provisioned, the client provides the decryption key, and the server starts a session, serving all the subsequent requests from that same user within the session duration.

The architecture is generic enough so that different applications can be integrated to the Application Wrapper with little implementation effort.

This work also evaluated the approach applicability in the context of industry-standard cloud infrastructure tools, such as Kubernetes, which provides simplicity, portability and a rich and well-established ecosystem of tools and services. The approach presented itself as viable with respect to end-user latency distribution and resource footprint, including scarce resources, such as Enclave Page Cache, when compared to the alternative, static approach of having all the application instances running at all times.

Furthermore, the following publications were produced, directly or indirectly, within the context

of this research work, some of them in international collaboration.

- *"Implementing Quality of Service and Confidentiality for Batch Processing Applications"* [3], on the 2018 IEEE/ACM International Conference on Utility and Cloud Computing Companion;
- Demo *"Asperathos: Running QoS-Aware Sensitive Batch Applications with Intel SGX"* [12], on the 37th Simpósio Brasileiro de Redes de Computadores e Sistemas Distribuídos;
- *"TEEMon: A continuous performance monitoring framework for TEEs"*, on the 21st ACM/I-FIP International Middleware Conference (accepted, to be published);
- Tutorial *"Processamento confidencial de dados de sensores na nuvem"*, on the 20th Simpósio Brasileiro em Segurança da Informação e de Sistemas Computacionais (accepted, to be published).

Finally, this work also led to the presentation of the technical talk *"Secure Data Analysis with OpenStack and Asperathos"* at an international computer industry conference, the 2019 OpenInfrastructure Summit<sup>1</sup>.

## 6.2 Future work

This work evaluated a complex architecture. Naturally, there are many interesting paths that deserve being further investigated and evolved. A few of them are listed below.

- **Explore auto-scaling.** One of the limitations of the *Standby Servers* approach is that the number of application instances is fixed per wrapper. However, the Wrapper offers a REST API that allow an external entity to start and stop application instances. This way, the Wrappers could be managed by an auto-scaling controller, for example, that could add or remove Wrappers (e.g., in the face of an infrastructure or topology change). Furthermore, the Wrappers themselves can be extended to listen to events (e.g., the application could log response times for each request) and take scaling decisions based on them, in order to meet previously defined performance constraints.

---

<sup>1</sup><https://www.openstack.org/summit/denver-2019/summit-schedule/events/23579/secure-data-analysis-with-openstack-and-asperathos>



- **Support generic applications seamlessly.** Although the architecture is generic enough to accommodate different applications, some implementation effort is still needed to make the sidecar threads listen to the right events. The Application Wrapper can be extended to allow these events (and respective actions) to be described in a declarative way.
- **Compare to serverless.** Compare the advantages and disadvantages of this approach with that of serverless architectures. Dynamic loading code to enable more efficient resource usage is also a goal of the serverless paradigm. Generalizing the applications supported would also enable a thorough comparison with this alternative paradigm.

## 6.3 Acknowledgments

This work was partially funded by CAPES (Coordenação de Aperfeiçoamento de Pessoal de Nível Superior) through the Master's program at the Departamento de Sistemas e Computação at Universidade Federal de Campina Grande. The current use case and architecture has been developed within the project UFCG-EMBRAPII ScalableSGX in a research collaboration between Universidade Federal de Campina Grande (UFCG), Scontain U.G., and Technische Universität Dresden (TU Dresden). The Routing Manager implementation is a contribution by Sergei Arnautov (Scontain U.G.). Fábio Silva (UFCG) contributed to purging mechanisms in the Intel SGX driver. Gabriel Vinha (UFCG) contributed mainly to the implementation of monitoring approaches for SGX enclaves. Lucas Cavalcante (UFCG) contributed mainly to the dynamic provision of client data using Ceph RBD. Although some of these aspects were not included in the scope of this work, their participation in the project provided continuous feedback.

# Bibliography

- [1] Ittai Anati, Shay Gueron, Simon Johnson, and Vincent Scarlata. Innovative technology for cpu based attestation and sealing. In *Proceedings of the 2nd international workshop on hardware and architectural support for security and privacy*, volume 13, page 7. Citeseer, 2013.
- [2] Sergei Arnautov, Bohdan Trach, Franz Gregor, Thomas Knauth, Andre Martin, Christian Priebe, Joshua Lind, Divya Muthukumar, Dan O’keeffe, Mark L Stillwell, et al. SCONE: Secure linux containers with intel SGX. In *12th USENIX Symposium on Operating Systems Design and Implementation (OSDI 16)*, pages 689–703, 2016.
- [3] Igor Ataide, Gabriel Vinha, Clenimar Souza, and Andrey Brito. Implementing quality of service and confidentiality for batch processing applications. In *2018 IEEE/ACM International Conference on Utility and Cloud Computing Companion (UCC Companion)*, pages 258–265. IEEE, 2018.
- [4] David Bernstein. Containers and cloud: From lxc to docker to kubernetes. *IEEE Cloud Computing*, 1(3):81–84, 2014.
- [5] Stefan Brenner and Rüdiger Kapitza. Trust more, serverless. In *Proceedings of the 12th ACM International Conference on Systems and Storage*, pages 33–43, 2019.
- [6] Eric A Brewer. Kubernetes and the path to cloud native. In *Proceedings of the sixth ACM symposium on cloud computing*, pages 167–167, 2015.
- [7] Brendan Burns, Brian Grant, David Oppenheimer, Eric Brewer, and John Wilkes. Borg, omega, and kubernetes. *Queue*, 14(1):70–93, 2016.
- [8] David Kaplan, Jeremy Powell, and Tom Woller. Amd memory encryption. *White paper*, 2016.
- [9] Frank McKeen, Ilya Alexandrovich, Alex Berenzon, Carlos V Rozas, Hisham Shafi, Vedvyas Shanbhogue, and Uday R Savagaonkar. Innovative instructions and software model for isolated execution. *Hasp@ isca*, 10(1), 2013.

- 
- [10] Saeid Mofrad, Fengwei Zhang, Shiyong Lu, and Weidong Shi. A comparison study of intel sgx and amd memory encryption technology. In *Proceedings of the 7th International Workshop on Hardware and Architectural Support for Security and Privacy*, pages 1–8, 2018.
- [11] Christian Priebe, Divya Muthukumaran, Joshua Lind, Huanzhou Zhu, Shujie Cui, Vasily A Sartakov, and Peter Pietzuch. Sgx-kl: Securing the host os interface for trusted execution. *arXiv preprint arXiv:1908.11143*, 2019.
- [12] Lília Rodrigues Sampaio, Clenimar Souza, Gabriel Silva Vinha, and Andrey Brito. Asperathos: Running qos-aware sensitive batch applications with intel sgx. In *Anais Estendidos do XXXVII Simpósio Brasileiro de Redes de Computadores e Sistemas Distribuídos*, pages 89–96. SBC, 2019.
- [13] Carlos Segarra, Ricard Delgado-Gonzalo, Mathieu Lemay, Pierre-Louis Aublin, Peter Pietzuch, and Valerio Schiavoni. Using trusted execution environments for secure stream processing of medical data. In *IFIP International Conference on Distributed Applications and Interoperable Systems*, pages 91–107. Springer, 2019.
- [14] G. Teodoro, T. Tavares, B. Coutinho, W. Meira, and D. Guedes. Load balancing on stateful clustered web servers. In *Proceedings. 15th Symposium on Computer Architecture and High Performance Computing*, pages 207–215, 2003.
- [15] Chia-Che Tsai, Donald E Porter, and Mona Vij. Graphene-sgx: A practical library OS for unmodified applications on SGX. In *2017 USENIX Annual Technical Conference (USENIXATC 17)*, pages 645–658, 2017.
- [16] Nico Weichbrodt, Pierre-Louis Aublin, and Rüdiger Kapitza. sgx-perf: A performance analysis tool for intel sgx enclaves. In *Proceedings of the 19th International Middleware Conference*, pages 201–213, 2018.
- [17] Johannes Winter. Trusted computing building blocks for embedded linux-based arm trustzone platforms. In *Proceedings of the 3rd ACM workshop on Scalable trusted computing*, pages 21–30, 2008.