


Tailored IoT & BigData Sandboxes and Testbeds for Smart,
Autonomous and Personalized Services in the European
Finance and Insurance Services Ecosystem



D5.11 – Data Management Workbench and Open APIs - II

Revision Number	3.0
Task Reference	T5.5
Lead Beneficiary	UBI
Responsible	Konstantinos Perakis
Partners	AGRO AKTIF BOI BOS CP GFT INNOV ISPRINT JSI LXS NBG PRIVE RB UBI WEA
Deliverable Type	Report (R)
Dissemination Level	Public (PU)
Due Date	2021-07-30
Delivered Date	2021-09-24
Internal Reviewers	ATOS, CTAG
Quality Assurance	CCA
Acceptance	WP Leader Accepted and Coordinator Accepted
EC Project Officer	Pierre-Paul Sondag
Programme	HORIZON 2020 - ICT-11-2018
	This project has received funding from the European Union's Horizon 2020 research and innovation programme under Grant Agreement no 856632

Contributing Partners

Partner Acronym	Role ¹	Author(s) ²
UBI	Lead Beneficiary	Konstantinos Perakis, Dimitris Miltiadou, Stamatis Pitsios
LXS	Contributor	Alejandro Ramiro, Luis Miguel Garcia, Javier Pereira
ATOS	Internal Reviewer	Jose Gatos
CTAG	Internal Reviewer	Marcos Cabeza Irida
CCA	Quality Assurance	Paul Lefrere

Revision History

Version	Date	Partner(s)	Description
0.1	2021-06-15	UBI	ToC Version
0.2	2021-06-21	UBI	Initial contributions to Section 4
0.3	2021-06-30	UBI, LXS	Contributions to Section 5
0.35	2021-07-12	UBI	Updated contribution to Section 4
0.40	2021-07-21	UBI, LXS	Updated contribution to Section 4
0.45	2021-08-09	UBI, LXS	Updated contribution to Section 4
0.50	2021-08-24	UBI	Updated contribution to Section 4
0.70	2021-09-09	UBI, LXS	Updated contributions to Sections 4 and 5
0.80	2021-09-17	UBI	Finalisation of sections 4 and 5
1.0	2021-09-21	UBI	First Version for Internal Review
2.0	2021-09-23	UBI	Version for Quality Assurance
3.0	2021-09-34	UBI	Version for Submission

¹ Lead Beneficiary, Contributor, Internal Reviewer, Quality Assurance

² Can be left void

Executive Summary

The goal of Task 5.5 “OpenAPI for Analytics and Integrated BigData/AI WorkBench” is to enable access to the added-value analytics functionalities of INFINITECH, that are offered via microservices implementations, and their respective Open APIs in an integrated way through a single entry-point. The document at hand, entitled “Data Management Workbench and Open APIs - II”, constitutes the second report of the work performed and the produced outcomes of Task 5.5. The main purpose of this deliverable is to deliver the **first prototype version of the INFINITECH Open API Gateway** which is implemented in accordance with the **detailed design specifications** which were documented in the first iteration of this deliverable. The first prototype version of the INFINITECH Open API Gateway aims at effectively addressing the accessibility and consumption challenges raised from the modular nature of the microservices implementations of added-value offerings in INFINITECH.

Towards this end, the current deliverable provides the updated documentation that will supplement the information documented in the deliverable D5.10, highlighting the updates from the existing documentation and introducing the technical details of the implementation of the first prototype version of INFINITECH Open API Gateway. The deliverable builds on top of the outcomes and knowledge extracted in D5.10 in order to provide the updated report of the work performed until M22.

Hence, the scope of the current report can be summarized in the following axes:

- To conduct a detailed and comprehensive **analysis of the challenges imposed by the microservices architecture**, in the consumption of the offered – by the underlying INFINITECH microservices – business functionalities from any client applications or services. The analysis takes into consideration commonly applied approaches, focusing on their advantages and disadvantages and documenting the rationale for the selection of the API Gateway pattern. Although the results remained unchanged from the previous iteration, they are included in the deliverable for coherency reasons.
- To document the **INFINITECH Open API Gateway design specifications** that formulate the solution of the single entry-point for the added-value analytics functionalities and other core offerings of INFINITECH, and a core ingredient of the INFINITECH Reference Architecture. The design specifications document the set of main functionalities of the component, the core design decisions and the design specifications that were driven by these decisions. Moreover, the high-level architecture of the component is defined and the two distinct modules that compose this architecture are documented in detail. Finally, the design specifications are complemented with the documentation of the supported use cases from each module and their respective sequence diagrams. Although the design specifications remained unchanged from the previous iteration, they are included in the deliverable for coherency reasons.
- To deliver the **first prototype version of the INFINITECH Open API Gateway** and provide the **technical documentation of the delivered version**. In this context, the deliverable provides the supplementary technical documentation of the delivered version, describing in detail the implementation aspects of the two core modules of the INFINITECH Open API Gateway, presenting the implemented functionalities and how these modules were successfully integrated to formulate the delivered version.
- To document the updated list of well-established **open-source technologies, libraries and frameworks** which are leveraged during the implementation phase of the INFINITECH Open API Gateway component.

The current deliverable presents the first prototype version of the INFINITECH Open API Gateway component which is delivered on M22. It should be stressed that the curation of both the definition of the design specifications, as well as the implementation of the presented component, is a living process that will last until M30 as per the INFINITECH Description of Action. Hence, the deliverable at hand constitutes a living document that will be updated in order to document all the enhancements and optimisations that will be introduced till M30, based on the feedback that will be collected by the INFINITECH stakeholders and the project’s evolution. The last iteration of the deliverable, namely D5.12, will be provided on M30 and will

constitute the final documentation of both the design specifications and the implementation details of the final version of the INFINITECH Open API Gateway.

Table of Contents

1	Introduction.....	8
1.1	Objective of the Deliverable	8
1.2	Insights from other Tasks and Deliverables.....	9
1.3	Structure	10
2	Motivation and Challenges	11
3	The INFINITECH Open API Gateway	14
3.1	Design overview.....	14
3.2	Design Specifications	16
3.3	Use Cases and Sequence Diagrams	18
3.3.1	Gateway.....	18
3.3.2	Service Registry	22
4	Implementation of the INFINITECH Open API Gateway	26
4.1	Gateway	26
4.1.1	Gateway Backend	28
4.1.2	Gateway Frontend.....	30
4.2	Service Registry.....	31
5	Baseline technologies and tools.....	34
6	Conclusions.....	36
	Appendix A: Literature	38

List of Figures

Figure 2-1:	Direct Client-to-Microservices communication.....	12
Figure 2-2:	API Gateway pattern.....	13
Figure 3-1:	INFINITECH Open API Gateway high-level architecture	17
Figure 3-2:	Open APIs discovery sequence diagram	19
Figure 3-3:	Request Handling (1-to-1).....	20
Figure 3-4:	Request Handling (1-to-many).....	22
Figure 3-5:	Self-registration sequence diagram	23
Figure 3-6:	Self-deregistration sequence diagram.....	24
Figure 3-7:	Periodic Health Check sequence diagram	25
Figure 4-1:	Gateway Module Architecture and interactions	26
Figure 4-2:	Gateway Request Processing Handling.....	27
Figure 4-3:	Discovery of microservices’ Open API documentation.....	31

List of Tables

Table 1:	INFINITECH Open API Gateway design aspects.....	16
Table 2:	Gateway – Open APIs discovery	18
Table 3:	Gateway – Request Handling (1-to-1)	19
Table 4:	Gateway – Request Handling (1-to-many).....	21

Table 5: Service Registry – Self-registration.....	22
Table 6: Service Registry – Self-deregistration.....	23
Table 7: Service Registry – Periodic Health Check	24
Table 8: INFINITECH Open API Gateway list of technologies.....	35
Table 9: Conclusions (TASK Objectives with Deliverable achievements)	37
Table 10: (map TASK KPI with Deliverable achievements)	37

Abbreviations/Acronyms

Abbreviation	Definition
AI	Artificial Intelligence
API	Application Programming Interface
DL	Deep Learning
IP	Internet Protocol
HTTP	Hypertext Transfer Protocol
JSON	JavaScript Object Notation
JWT	JSON Web Token
KPI	Key Performance Indicator
ML	Machine Language
RA	Reference Architecture
REST	Representational State Transfer
URL	Uniform Resource Locator

1 Introduction

D5.11 is released in the scope of WP5 “Data Analytics Enablers for Financial and Insurance Services”, which enumerates the associated activities and documents the updated outcomes of Task 5.5 “OpenAPI for Analytics and Integrated BigData/AI WorkBench”. The specific deliverable is prepared in accordance with the INFINITECH Description of Action and reports the updates to the activities performed and the results of the work performed as the second iteration within the context of Task 5.5. It aims at providing the updated detailed documentation of the designed solution for the effective and efficient access to the added-value analytics functionalities and other core offerings of INFINITECH through a sophisticated, integrated and single-entry point manner.

The deliverable at hand is building on top of the results documented in the first iteration, namely deliverable D5.10, providing the technical details of the first prototype version of the INFINITECH Open API Gateway. As documented also in the first iteration of the deliverable, the modular nature of the microservices architecture, together with the dynamic nature of the virtualised infrastructures where these microservices are usually deployed, dictates the need for a full solution that will address the problem of how the clients of a microservices-based platform can effectively and efficiently consume their offered business capabilities. To this end, the consortium embraced a practical and well-defined pattern to overcome the problem of clients directly accessing the underlying microservices once they have managed to discover their latest network location. This pattern has driven the design of a full solution that takes into consideration that the microservices constitute the main ingredient of the INFINITECH platform, and the various analytics and other main functionalities of INFINITECH will be based on microservices. The goal of the designed solution is to eliminate the barriers of accessing these offered functionalities from different clients. Following this approach leads to the design specifications for a solution that will be easily integrated into the INFINITECH platform, and will constitute the single-point-of-entry for those functionalities defined during the first iteration.

In this second iteration, the documentation of the first prototype version of the INFINITECH Open API Gateway is provided. The implementation of the prototype was driven by the design specifications of the first iteration that remained unchanged. The delivered solution enables the INFINITECH project’s pilots, as well as the stakeholders of the financial sector to consume the offered functionalities in a straight-forward and easy manner, and without requiring knowledge of the underlying microservices architecture. It contains a core subset of designed features and sets the path for the next release of the INFINITECH Open API Gateway that is scheduled for M30, in accordance with the INFINITECH Description of Action, which will be documented in deliverable D5.12.

1.1 Objective of the Deliverable

The purpose of this deliverable is to report the outcomes of the work performed within the context of Task 5.5 at this phase of the project (M22). The deliverable constitutes the second iteration with the main goal to document the advancements in comparison with the first iteration and in particular to provide the documentation of the first prototype version of the INFINITECH Open API Gateway whose design specifications were documented in the first iteration.

During this second period (from M14 till M22), the activities mainly focused on the implementation of the designed modules of the INFINITECH Open API Gateway and their functionalities, as well as their successful integration in order to formulate the first prototype version of the INFINITECH Open API Gateway. To this end, the implementation details for each module are documented in detail.

The deliverable aims at providing the updated documentation of the information that has been documented in the previous iteration and, for coherency reasons, it contains the information included in the previous iteration, highlighting the updates and optimisations that were introduced where needed. The revised information is presented utilising the approach that was followed in the previous iteration.

Hence, the first objective of the deliverable is to present the results, which remained unchanged from the previous iteration, of the analysis performed towards addressing the challenges imposed in the microservices architecture for the clients that need to consume the business functionalities offered by the microservices. To this end, the two most common approaches, namely the direct client-to microservices and the API Gateway pattern, are presented in detail by documenting the benefits and drawbacks of each approach.

The second object of the deliverable is to present the design specifications of the INFINITECH Open API Gateway component. The design specifications remained also unchanged from the previous iteration and they document: a) the rationale for the adoption of the API Gateway pattern which eliminates the barriers of accessing the underlying microservices implementations of the INFINITECH platform, b) the design overview of the component along with the main functionalities that the component offers, c) the documentation of the core design decisions that were taken during the design phase of the component, d) the detailed specifications of the core modules of the modular architecture of the INFINITECH Open API Gateway component, namely the Gateway and the Service Registry, describing their offered functionalities in accordance with the design decisions and e) to document the list of supported use cases and their respective sequence diagram that depicts the interactions of the involved clients and modules.

The third objective is the introduction of the detailed documentation of the implementation of the first prototype version of the INFINITECH Open API Gateway. The specific deliverable documents for each module the functionalities which were implemented in this first prototype version, the overview of the implemented version of each module, accompanied by the details of how these modules are interacting in order to provide the end-to-end functionalities of the INFINITECH Open API Gateway.

The fourth objective of the deliverable is to present the updated list of baseline technologies and tools that will be leveraged in the implementation phase of the INFINITECH Open API Gateway component. The initial list of technologies, libraries and frameworks that was documented in the first iteration received several updates in accordance with the implementation activities and decisions taken during the implementation of the first prototype version of the INFINITECH Open API Gateway.

It should be noted that according to the INFINITECH Description of Action, Task 5.5 lasts until M30. Hence, the final version of the deliverable will be released on M30 with deliverable D5.12. Thus, the upcoming version of the deliverable will constitute the final iteration and will provide the final and complete implementation details of the INFINITECH Open API Gateway component, as well as the optimisations and enhancements that will be introduced to the design specifications taking into consideration the evolvement of the project, as well as the feedback that will be collected by the pilots and the stakeholders of INFINITECH.

1.2 Insights from other Tasks and Deliverables

Deliverable D5.11 is released in the scope of WP5 “Data Analytics Enablers for Financial and Insurance Services” and documents the updated outcomes of the work performed within the context of Task 5.5 “OpenAPI for Analytics and Integrated BigData/AI WorkBench”. The task is directly related to the rest of the Tasks of WP5, enabling access to the added-value analytics functionalities of INFINITECH which are formulated by the library of ML/DL algorithms for Financial and Insurance Services (as produced by T5.4), that incorporates the incremental and parallel data analytics (produced by T5.2) as well as the declarative real-time analytics (produced by T5.3) and leverages the data collection process for the algorithm’s training and evaluation (produced by T5.1).

In addition, the task builds on top of the outcomes of WP2 “Vision and Specifications for Autonomous, Intelligent and Personalized Services” in which the overall requirements of the INFINITECH platform are defined. Specifically, Task 5.5 received as input the collected user stories of the pilots of the project and the extracted user requirements, as documented in deliverables D2.1 and D2.2 that report the work performed in Task 2.1. Furthermore, the inputs of the task included the elicited technical requirements and the fundamental building blocks of the INFINITECH platform and their specifications, as reported in deliverables D2.5 and D2.6. Finally, the task received as input the outcomes of T2.7, in which the INFINITECH Reference Architecture (INFINITECH RA) was formulated, as documented in deliverables D2.13 and D2.14, during the

design specifications definition of the INFINITECH Open API Gateway component that constitutes a part of the INFINITECH platform.

1.3 Structure

This document is structured as follows:

- Section 1 introduces the document, describing the context of the outcomes of the work performed within the task and highlights its relation to the rest of the tasks and deliverables of the project.
- Section 2 documents the motivation and challenges in the access of the deployed microservices of a microservices-based platform by the clients.
- Section 3 presents the design overview and design specifications of the INFINITECH Open API Gateway component, as well as the use cases addressed, along with the corresponding sequence diagrams.
- Section 4 presents the implementation details of the delivered version of the INFINITECH Open API Gateway component.
- Section 5 presents the list of baseline technologies and tools that are utilised in the implementation of the INFINITECH Open API Gateway component.
- Section 6 concludes the document.

2 Motivation and Challenges

Updates from D5.10:

This particular section remained unchanged from the previous version. It presents an analysis of the challenges imposed in the consumption of business functionalities in the microservices architectures and commonly applied approaches to overcome them.

In traditional monolithic applications, the application is designed and implemented with the aim of executing a specific, commonly domain-specific, business logic and interacts with the client-side or any third-party client (applications) via an exposed API. In this sense, the clients will usually retrieve the required data or invoke any specific operation by executing a single REST call, that is received and handled by the underlying server-side (backend) application. While the monolithic application approach has a number of benefits, such as simplifying development and deployment, as well as horizontal scalability by multiple running instances behind a load balancer, it also introduces a large number of drawbacks that are very critical and create severe barriers to application execution and evolution. The monolithic architecture has a limitation when it comes to application size and complexity. When the application grows in size, the number of functionalities, complexity and sustainability issues arise. The source code becomes difficult to understand and maintain, while also changes in the functionalities offered (and/or the introduction of new functionalities) require great development effort, and add multiple dependencies that should be met, hence the complexity grows. Additionally, the scaling and continuous deployment becomes problematic, since a small update might require a redeployment of all the applications once extensive testing has been successfully performed.

The introduction of the microservices architecture approach solves the drawbacks of the traditional monolithic architecture approach by structuring the underlying application as a collection of microservices that are loosely-coupled, highly maintainable and testable, independently-deployable, organised by business capabilities and can be owned by different development teams [1]. Nevertheless, while the microservices architecture enables the effective and efficient development of large and complex applications in a rapid, frequent and reliable manner, it also introduces a number of challenges that need to be addressed.

In the microservices architecture, the application is partitioned in multiple microservices, where each one of them undertakes a specific business capability of the application and a set of responsibilities. However, this business-capabilities decomposition imposes several requirements on the underlying microservices. Each microservice should provide an API that enables the needed intercommunication between the various services. As the context of each microservice is clearly defined, the provided API is usually fine-grained and restricted to the assigned business capability. Hence, the granularity of the APIs of the microservice is usually different than the one required by any client of the application, and there are cases where the client is required to invoke multiple microservices in order to obtain all the required data or execution results. Furthermore, the evolution of the application usually requires the introduction of new microservices or even the refinement of existing ones to support a new business capability. The dynamic nature of the virtualised infrastructures that are hosting the applications, and the microservices that compose these architectures, enables the deployment, upgrade, scaling and restart of each microservice independently of the rest of the microservices of the application. Hence, both the number of microservices, as well as their connection details, such as the hostname and their IP addresses might change and vary dynamically. Thus, the transition from a traditional monolithic application to a microservice application, imposes the problem of how the clients of this application can effectively and efficiently consume the offered business capabilities by either a specific microservice or a combination of microservices.

To this end, the problem can be narrowed down to the client-to-microservices communication. Two approaches can be followed to solve this problem. In the first approach, usually referenced as direct client-to-microservice communication [2], the client is able to directly make requests to the microservices of the application. In this case, each microservice provides an exposed API endpoint and can be reached through a specific URL that maps to a load balancer which in turn distributes the incoming requests to the requested microservice instances (see Figure 2-1).

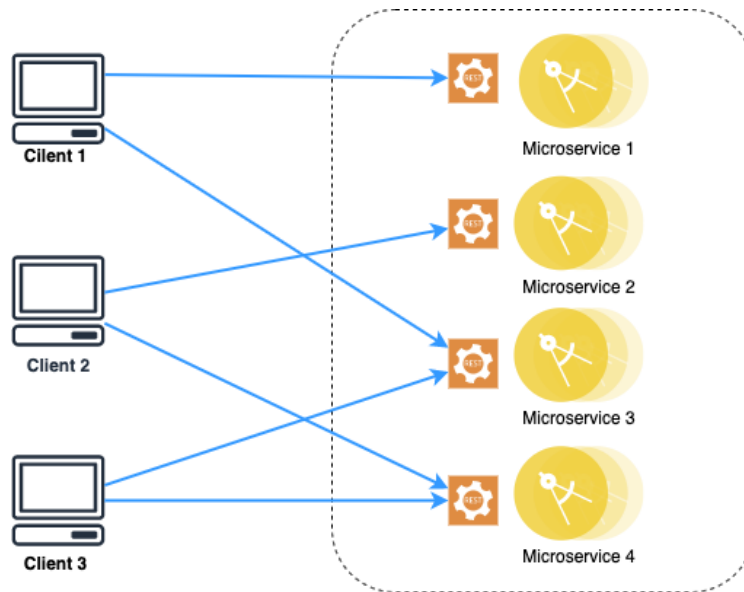


Figure 2-1: Direct Client-to-Microservices communication

However, this approach has several limitations and challenges that should be taken into consideration. The first challenge arises due to the fact that, by design, each microservice offers a fine-grained API with a specific context. Thus, there are cases where the client is forced to make multiple requests to different microservices, hence multiple server round trips, to collect the information needed to complete a single operation. This is inefficient and increases the latency in some cases, while it also increases the complexity in the source code of the client. An additional challenge is that microservices might use different and diverse communication protocols (e.g. HTTP, AMQP, Thrift) and as a consequence the client should be capable of interacting with all these diverse communication protocols. Moreover, common functionalities such as authorisation, authentication and logging, need to be taken care of at the microservices level, rather than at a common cross-cutting level on top of the microservices. Finally, as the implementation of the client is directly connected with the underlying microservices, any changes or updates in the existing microservices, such as the merging or splitting of an existing microservice, or the introduction of new ones, usually propagates changes on the clients that need to adapt to the evolution of the microservices.

The second approach that addresses all the above-mentioned challenges is the API Gateway pattern [3]. The API Gateway is introduced as an intermediate layer between the clients and the underlying microservices, and acts as a single entry-point for all clients that consume these microservices. All incoming requests are generated towards the API Gateway which serves them in two core ways. The API Gateway usually routes the request to the appropriate microservice acting as a reverse proxy, however in some cases it handles the request by invoking multiple microservices and aggregating the results, which are provided back to the requestor (see Figure 2-2). In addition to this, the API Gateway undertakes several cross-cutting functionalities such as the authorisation, authentication, monitoring and load balancing.

The API Gateway encapsulates the underlying system architecture, hiding the details of how the application is partitioned into microservices, providing different APIs for diverse clients if needed. It eliminates the need for clients to discover or keep track of the network locations of the microservices and their exposed APIs that can change dynamically. Furthermore, it reduces the problem of multiple server round trips, as it handles the invocation of multiple microservices and aggregation of the results. The API Gateway is able to handle the diverse communication protocols problem, since it exposes a well-defined and web-friendly API, undertaking the protocol conversion internally. Thus, the API Gateway is able to simplify both the communication with the clients, as well as the complexity of the client's source code. The client is only interacting with the API Gateway, hence the update or evolution of the application's microservice is hidden from the client. Furthermore, it handles cross-cutting functionalities, while the development of the underlying microservices is also simplified.

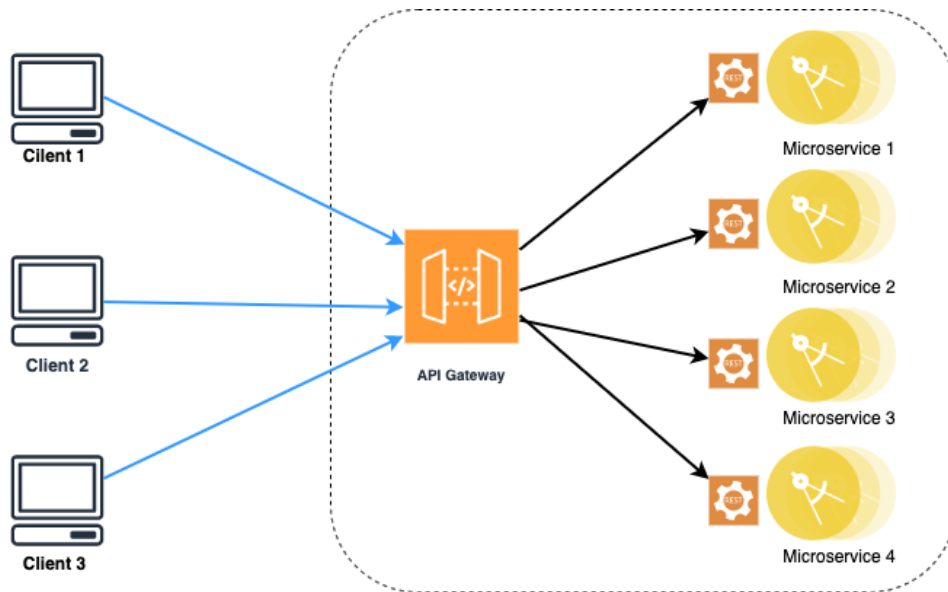


Figure 2-2: API Gateway pattern

However, the API Gateway pattern also has some drawbacks. As it constitutes a component with high availability, it should be designed, developed, deployed and maintained in a proper way in order to avoid becoming a bottleneck for the application. Hence, it is crucial that several aspects such as performance, adaptability and fault tolerance are taken into consideration to benefit from the API Gateway pattern solution. Nevertheless, despite these drawbacks, the API Gateway pattern is considered to be the proper approach in the case of INFINITECH, as it effectively solves all the fundamental issues and challenges described in the direct client-to-microservice communication.

3 The INFINITECH Open API Gateway

Updates from D5.10:

This particular section remained unchanged from the previous version. It presents the design overview of the INFINITECH Open API Gateway, documenting the design specifications of the proposed solution, as well as the use cases addressed supplemented by the relevant sequence diagrams.

3.1 Design overview

The **INFINITECH Open API Gateway** is a sophisticated API Gateway that encompasses the Open API specification in order to provide a single point of entry for the added-value functionalities of INFINITECH which are based on microservices. The work that was performed during this phase was mainly focused on the Machine Learning (ML) / Deep Learning (DL) analytics functionalities of INFINITECH, which are implemented as microservices. Following the API Gateway pattern, as presented in Section 2, it is capable of effectively handling and determining the network location of dynamically-deployed microservice instances, while at the same time hiding the internal application's architecture from the clients of the application.

Hence, the main functionalities of the INFINITECH Open API Gateway are as follows:

- To act as a single point of entry for the ML/DL analytics functionalities of INFINITECH, handling the incoming requests towards the dynamically-deployed microservice instances;
- To facilitate the communication between the underlying microservice instances;
- To enable the self-registration of the microservice instances;
- To facilitate the discovery of the microservice instances;
- To enable the discovery of the microservice instance Open APIs;
- To effectively handle the unavailability or unreachability of the underlying microservice instances;
- To provide authentication, monitoring and logging functionalities to the microservice instances.

During the design phase of the INFINITECH Open API Gateway, several design decisions were taken based on the specifications set by the INFINITECH Reference Architecture, as documented in deliverable D2.13; the elicited technical requirements are reported in deliverable D2.5, as well as the requirements of the pilots and stakeholders of INFINITECH, as documented in deliverable D2.1. Within the INFINITECH Reference Architecture, the INFINITECH Open API Gateway is positioned in the Interface layer of the architecture undertaking the task of providing the required information from the Analytics layer to the Presentation layer where the clients reside.

With regards to the *Request Handling* processes of the INFINITECH Open API Gateway, it was decided that it should support both core ways of the API Gateway pattern. Hence, the INFINITECH Open API Gateway should be able to handle one to one (1-to-1) microservices invocation, meaning it can handle requests that are simply translated by routing the incoming request to the appropriate microservice. On the other hand, the INFINITECH Open API Gateway should handle one to many (1-to-many) microservices invocation, thus a request could be translated into the invocation of multiple microservices, whose results are aggregated and returned to the requestor. In the latter case, a sophisticated logic will be employed to effectively handle the multiple invocations, dependencies between the invocations, and failure handling through modern approaches, such as reactive programming with the promises pattern instead of the traditional asynchronous call-back pattern.

As the underlying microservices are operating in a distributed manner, it is required to establish a set of *Communication Mechanisms* across the different processes that enable the invocation of the microservices in both an asynchronous and synchronous manner. To this end, both approaches will be supported with the proper mechanisms for asynchronous message-based communication and synchronous communication.

One of the core aspects of the INFINITECH Open API Gateway is the effective and dynamic *Service Registry* and *Service Discovery* processes. The INFINITECH Open API Gateway is designed to be able to find and keep track

of the network location of each microservice that it handles the requests for. As explained in Section 2, microservices are usually operating in a cloud infrastructure, within virtual machines or containers, that change dynamically in terms of IP addresses and ports. Hence, it is imperative that the INFINITECH Open API Gateway incorporates the appropriate Service Registry and Service Discovery processes that effectively handle the changes of the microservice instances in a dynamic manner.

Service Registry is a database that maintains the updated network locations of the microservices. Hence, the microservice instances are registered with the Service Registry on start-up and deregistered on shutdown. In the meanwhile, the Service Registry performs periodic health checks on the registered microservices to ensure their availability. For the service registration, there are two approaches, the self-registration pattern and the third-party registration pattern. In the Self-Registration pattern [4], the microservice instance is responsible for the registration and deregistration of itself with the service registry, while in the third-party registration pattern [5] the registration and deregistration is handled by a third party application called registrar, by performing polling operations or event subscriptions. In the INFINITECH Open API Gateway, the self-registration pattern is followed as it keeps the complexity of the component at the lower levels, while the coupling of the registration process with the microservices requires less effort, as there is a large number of available libraries performing these operations with out-of-the-box integration.

With regards to Service Discovery there are two approaches, namely the client-side service discovery and the server-side discovery. In the client-side service discovery [6], the client retrieves the latest network location of a microservice by querying the Service Registry. While this approach is considered more efficient in terms of network requests, it bounds the client implementation with the Service Registry. On the other hand, in the server-side discovery [6] the client makes requests to the API Gateway and the API Gateway queries the Service Registry before it finally routes the request to the latest network location of the requested microservice. In the INFINITECH Open API Gateway, the server-side discovery is followed as it eliminates the need to couple the client's implementation with the Service Registry, and additionally it hides the implementation details of the API Gateway and the registered microservices from the client.

A crucial aspect of the INFINITECH Open API Gateway is the embracement of **Open APIs** to enhance the accessibility of the underlying ML / DL microservices. The Open API specification is proposed by the Open API initiative, which is an open-source collaboration project of the Linux foundation, and defines a standard, language-agnostic interface to RESTful APIs which allows both humans and computers to discover and understand the capabilities of the service, without access to the source code, documentation, or through network traffic inspection [7]. The INFINITECH Open API Gateway will embrace the Open API specification which all underlying ML / DL microservices are required to follow in order to help the clients understand and consume their exposed services without the need for knowledge of the implementation details or access to the source code of the microservices. Hence, the INFINITECH Open API Gateway will facilitate publishing the offered Open APIs of the microservices in order to enable their easy and effortless discovery by the clients of the INFINITECH Open API Gateway.

Another aspect of the INFINITECH Open API Gateway is *Fault Tolerance*. For either 1-to-1 or 1-to-many microservices invocation, a mechanism performing efficient failure handling should be employed. Towards this end, the Circuit Breaker pattern [8] will be followed in the INFINITECH Open API Gateway. The specific pattern is utilised to handle the case where one of the invoked microservices is unavailable or unreachable and the whole process is not stalled, while the allocated resources are properly managed. In this pattern, if a microservice reaches a threshold of consecutive failures, then the circuit breaker that acts as a proxy will immediately reject all upcoming requests to this microservice for a predefined testing period. When this testing period ends, a limited number of test requests are made; if they succeed, the circuit breaker allows the invocation of the microservice again, else a new testing period is started again.

One final aspect of the INFINITECH Open API Gateway is the offering of cross-cutting functionalities such as authentication, monitoring and logging. With regards to the authentication, the INFINITECH Open API Gateway will embrace the Access Token pattern and specifically the JSON Web Token (JWT). The Open API Gateway that acts as the single point of entry for the ML/DL functionalities of INFINITECH will authenticate the incoming requests from the clients and will provide a JWT that securely authenticates the requestor back

to the microservices. Furthermore, the INFINITECH Open API Gateway performs continuous monitoring and detailed logging of all the received requests and executed operations for security, malicious activity and performance analysis purposes.

The following table summarizes the main design decisions that were taken during the design phase of the INFINITECH Open API Gateway.

Table 1: INFINITECH Open API Gateway design aspects

Design Aspect	INFINITECH Open API Gateway
Request Handling	Both 1-to-1 and 1-to-many microservices invocation
Communication Mechanisms	Both asynchronous message-based and synchronous communication
Service Registry	Self-registration pattern
Service Discovery	Server-side discovery pattern
API Specifications	Open API specification
Fault Tolerance	Circuit Breaker pattern
Cross-cutting functionalities	Authentication, Monitoring, Logging

It should be noted that the presented design aspects of the INFINITECH Open API Gateway do not limit the support for additional added-value functionalities which are based on microservices. On the contrary, the described approach facilitates the expansion of the list of supported microservices implementations in an effortless and effective manner as per the project's needs

3.2 Design Specifications

During the analysis of the client-to-microservices communication problem, it was clear why the API Gateway pattern is considered the dominant solution for this problem, as it effectively and efficiently resolves the problem in a solid and robust manner. Hence, the consortium decided to design and implement the INFINITECH Open API Gateway that is based on the API Gateway pattern, aiming at providing the required access to the underlying ML/DL microservices. However, as presented in the overview of the INFINITECH Open API Gateway in the previous section, there are several aspects that were carefully addressed in the design phase towards the implementation of the optimal (most effective and efficient) solution that will address the requirements of INFINITECH stakeholders.

The INFINITECH Open API Gateway has a modular architecture composed of two core modules, namely the *Gateway* and the *Service Registry*. Each module has a clear scope and a distinct context undertaking a specific set of responsibilities. Additionally, the modules are interacting through well-defined REST APIs in order to perform the main operations of the INFINITECH Open API Gateway. Their role in the INFINITECH Open API Gateway is as follows:

- The Gateway provides the single-entry point for all incoming requests from the pilots and undertakes the responsibility of invoking the appropriate microservices while also ensuring the required inter-communication between the microservices where needed.
- The Service Registry provides the database that maintains the updated network locations of the microservices. Furthermore, it provides the mechanism for self-registration and deregistration, while also performing the health check operations.

Figure 3-1 depicts the high-level architecture of the INFINITECH Open API Gateway. As illustrated, the Gateway can receive new requests from the clients for the invocation of the underlying ML/DL microservices. Upon receiving a new request, the Gateway consults the Service Registry in order to discover the updated network location of the requested microservice. Once the information is retrieved by the Service Registry,

the Gateway, acting as a reverse proxy, routes the request to the appropriate microservice by invoking the respective endpoint of the microservice. The selected microservice handles the request and provides the results back to the Gateway. At this final step, the Gateway replies to the client with the results as provided by the respective ML/DL microservice.

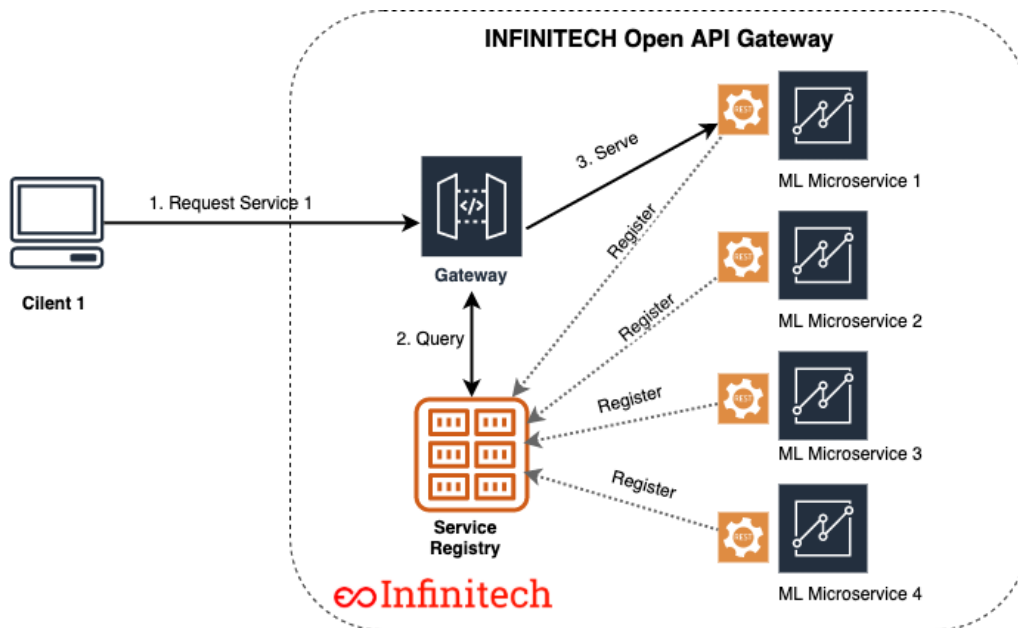


Figure 3-1: INFINITECH Open API Gateway high-level architecture

The execution is slightly different in the case where a request requires the invocation of multiple microservices. In this case, the Gateway upon receiving the request, consults the service registry and retrieves the list of microservices. Then, multiple requests are initiated to the respective microservices. Each microservice provides the results back to the Gateway and when all results are received the Gateway aggregates the results and sends them back to the client.

The scope of the Gateway is to facilitate the invocation of ML/DL microservices, hiding the details of the deployed microservices from the client. It undertakes all the core functionalities of the INFINITECH Open API Gateway utilising the services of the Service Registry to fulfil its purpose. Furthermore, the Gateway provides cross-cutting operations to the respective ML/DL microservices decoupling their implementation from these operations.

In accordance with the INFINITECH Open API Gateway overview presented in section 3.1, the main functionalities of the Gateway module are as follows:

- It performs the handling of incoming requests, acting as a reverse proxy and routing the incoming requests to the respective ML/DL microservices either for 1-to-1 microservices invocation or 1-to-many microservices invocation.
- It provides the communication mechanism for the microservices invocation in both asynchronous and synchronous ways.
- It performs the service discovery operations by interacting with the Service Registry during the request handling process.
- It enables the publishing of the Open APIs of the ML/DL microservices.
- It implements the fault tolerance mechanism for the effective handling of failures from the microservices side implementing the Circuit Breaker pattern.
- It undertakes the authentication of the requestor operations utilising JWT which are passed to the respective microservice.
- It performs the continuous monitoring and logging of all operations performed in the INFINITECH Open API Gateway

The scope of the Service Registry is to support the operations performed by the Gateway for the service discovery. In particular, the Service Registry is the complementary module that is capable of maintaining the updated network location of each microservice, and providing the appropriate information to the Gateway. Additionally, it decouples the service registration and deregistration from the INFINITECH Open API Gateway.

In accordance with the INFINITECH Open API Gateway overview presented in section 3.1, the main functionalities of the Service Registry module are as follows:

- It provides the mechanism for the self-registration and self-deregistration of the ML/DL microservices by employing the mechanism that handles these requests, and the registration client that should be integrated in the ML/DL microservices.
- It supports the service discovery operations by providing the required information to the Gateway.
- It performs periodic health checks upon the registered ML/DL microservices to ensure the availability of the registered microservices.

3.3 Use Cases and Sequence Diagrams

As explained in Section 3.2, the INFINITECH Open API Gateway is composed of two core modules, namely the Gateway and the Service Registry. In this section, the detailed use cases that each specific module addresses are documented, presenting in detail the relevant information. Additionally, for each use case, the corresponding sequence diagram that depicts the interactions of the modules and the involved clients is presented.

3.3.1 Gateway

3.3.1.1 Open APIs discovery

The specific functionality enables the discovery of the Open APIs of the registered microservices by the client. To achieve this, the Gateway maintains and displays a dynamic list where all registered microservices are listed. For each microservice, the list of Open APIs can be retrieved by the user. The list of Open APIs is made available to the Gateway during the self-registration process of each microservice.

Table 2: Gateway – Open APIs discovery

Stakeholders involved:	Client's User
Pre-conditions:	1. The existence of at least one registered microservice with Open APIs
Post-conditions:	1. The client's user is able to retrieve the list of Open APIs of a registered microservice
Data Attributes	None
Normal Flow	<ol style="list-style-type: none"> 1. The client's user navigates to the page where the dynamic list of registered microservices is displayed by the Gateway 2. The client's user selects one of the microservices from the list 3. The list of Open APIs is displayed to the user following the Open API specification
Pass Metrics	1. The client's user is able to retrieve the list of Open APIs for the selected microservice

Fail Metrics	1. The client’s user cannot retrieve the list of Open APIs for the selected microservice
---------------------	------------------------------------------------------------------------------------------

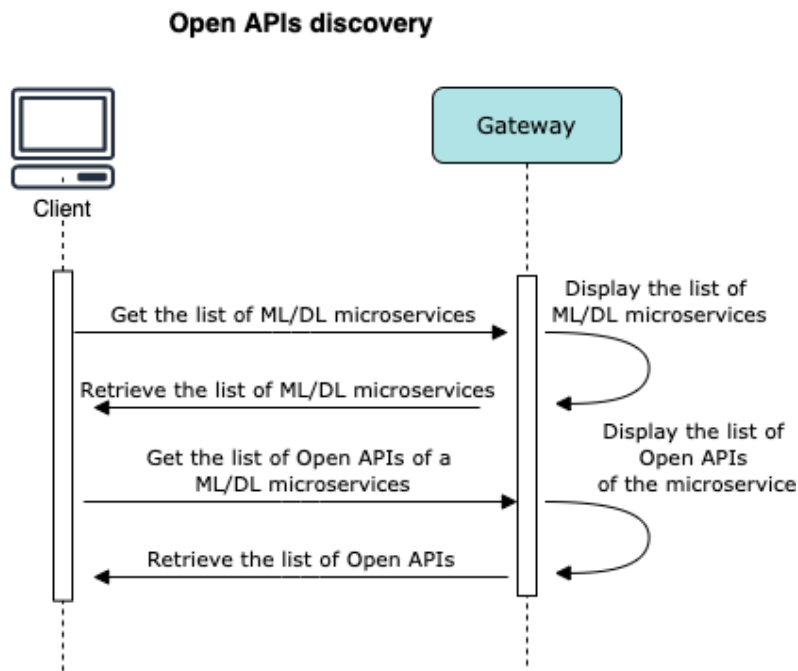


Figure 3-2: Open APIs discovery sequence diagram

3.3.1.2 Request Handling (1-to-1)

The specific functionality handles the incoming request for the invocation of the ML/DL microservice. The Gateway receives the request from the client and consults the Service Registry to retrieve the updated network location of the requested microservice. Upon retrieving the network location, the Gateway routes the request to the requested ML/DL microservice. The microservice receives and handles the request. The results are provided back to the Gateway which in turn returns them to the client. The specific use case also handles the case where the requested ML/DL microservice is unreachable or irresponsible.

Table 3: Gateway – Request Handling (1-to-1)

Stakeholders involved:	Client, Gateway
Pre-conditions:	1. The existence of at least one registered microservice with Open APIs
Post-conditions:	1. The client retrieves the results of the executed ML/DL microservice
Data Attributes	<ul style="list-style-type: none"> 1. The endpoint of the desired microservice, prefixed with the name that the corresponding microservice used to register with the registry. 2. The data that needs to be forwarded to the desired microservice 3. The JWT that will be used to authenticate the client
Normal Flow	1. The client initiates a request for a specific ML/DL microservice to the Gateway

	<ol style="list-style-type: none"> 2. The Gateway consults the Service Registry to retrieve the updated network location of the ML/DL microservice 3. Upon the retrieval of the updated network location, the Gateway routes the request to the specific ML/DL microservice <p><i>In the case where the ML/DL microservice is available and operational:</i></p> <ol style="list-style-type: none"> 4. The ML/DL microservice receives and handles the request. The results are provided to the Gateway 5. The Gateway propagates the results to the client <p><i>In the case where the ML/DL microservice is unavailable and irresponsible:</i></p> <ol style="list-style-type: none"> 4. The request is not properly received by the requested ML/DL microservice 5. The Gateway is informed and reports the failure to the requestor
Pass Metrics	<ol style="list-style-type: none"> 1. The requested ML/DL microservice is invoked successfully and the results of the execution are returned to the client 2. In the case where the requested ML/DL microservice is unavailable and irresponsible, the client is informed of the failure of the request
Fail Metrics	<ol style="list-style-type: none"> 1. The requested ML/DL microservice cannot be invoked and the request fails 2. The client is not informed when a failure occurs

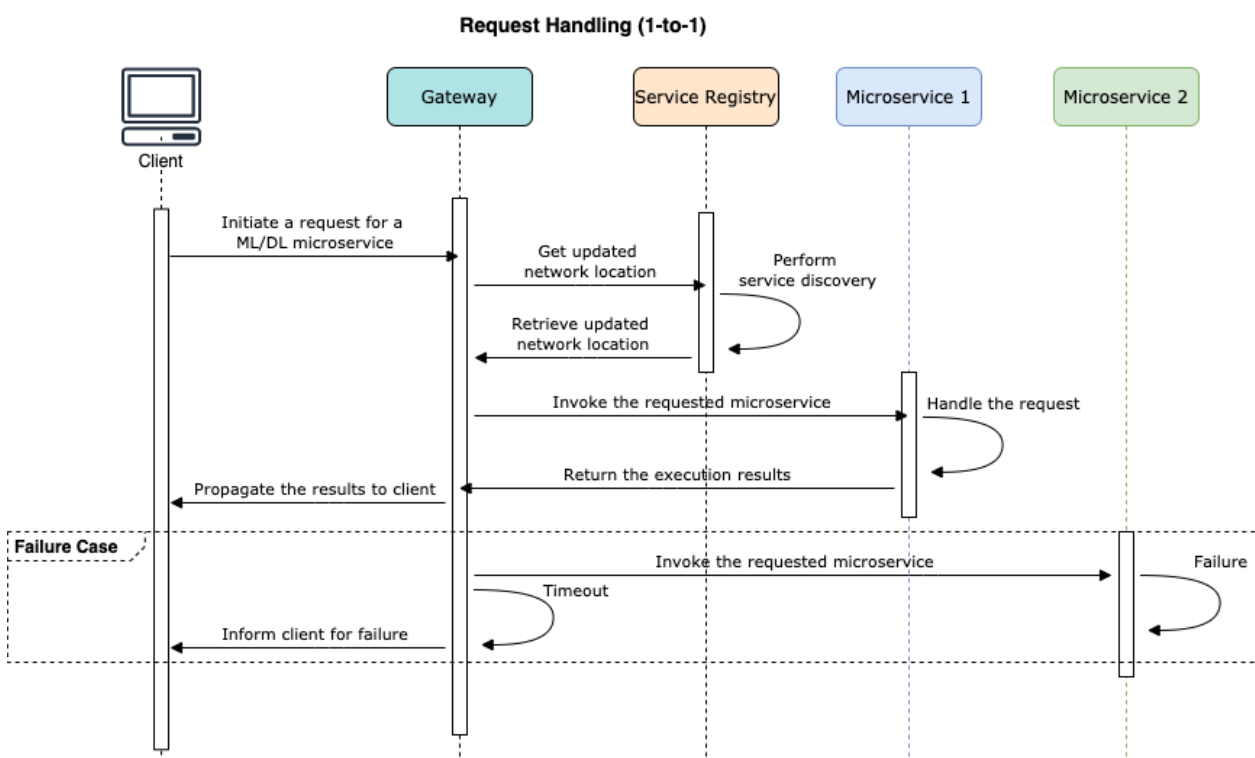


Figure 3-3: Request Handling (1-to-1)

3.3.1.3 Request Handling (1-to-many)

The specific functionality handles the incoming request for the invocation of a set of ML/DL microservices. The Gateway receives the request from the client and consults the Service Registry to retrieve the updated network location of the requested microservices. Upon retrieving the network location of the microservices,

the Gateway initiates one request to each of the ML/DL microservices involved in the request. The respective microservices receive and handle the requests. The results from each microservice execution are provided back to the Gateway. Upon receiving all the results, the Gateway aggregates the results and returns them to the client. The specific use case also handles the case where one of the requested ML/DL microservices is unreachable or irresponsive.

Table 4: Gateway – Request Handling (1-to-many)

Stakeholders involved:	Client, Gateway
Pre-conditions:	1. The existence of at least two pre-registered microservices with Open APIs
Post-conditions:	1. The client retrieves the results of the executed ML/DL microservices
Data Attributes	<ol style="list-style-type: none"> 1. A list of the endpoints of the desired microservices, prefixed with the names that the corresponding microservices used to register with the registry. 2. The data that needs to be forwarded to each of the desired microservices 3. The JWT that will be used to authenticate the client
Normal Flow	<ol style="list-style-type: none"> 1. The client initiates a request to the Gateway. 2. The Gateway retrieves the request which involves the execution of the multiple microservices. It consults the Service Registry to retrieve the updated network location of the ML/DL microservices 3. Upon the retrieval of the updated network locations of the microservices, the Gateway initiates a request to each specific ML/DL microservice <p><i>In the case where all ML/DL microservices are available and operational:</i></p> <ol style="list-style-type: none"> 4. Each ML/DL microservice receives and handles the request. The results are provided to the Gateway 5. Upon receiving all the results from all invoked microservices, the Gateway aggregates the results 6. The Gateway propagates the aggregated results to the client <p><i>In the case where one of the ML/DL microservices is unavailable and irresponsive:</i></p> <ol style="list-style-type: none"> 4. One of the requests is not properly received by the requested ML/DL microservice 5. The Gateway is informed, stops the processing and reports the failure to the requestor
Pass Metrics	<ol style="list-style-type: none"> 1. The requested ML/DL microservices are invoked successfully and the aggregate results are returned to the client 2. In the case where one of the requested ML/DL microservice is unavailable and irresponsive, the client is informed of the failure of the request
Fail Metrics	1. The requested ML/DL microservices cannot be invoked and the request fails

2. The client is not informed when a failure occurs

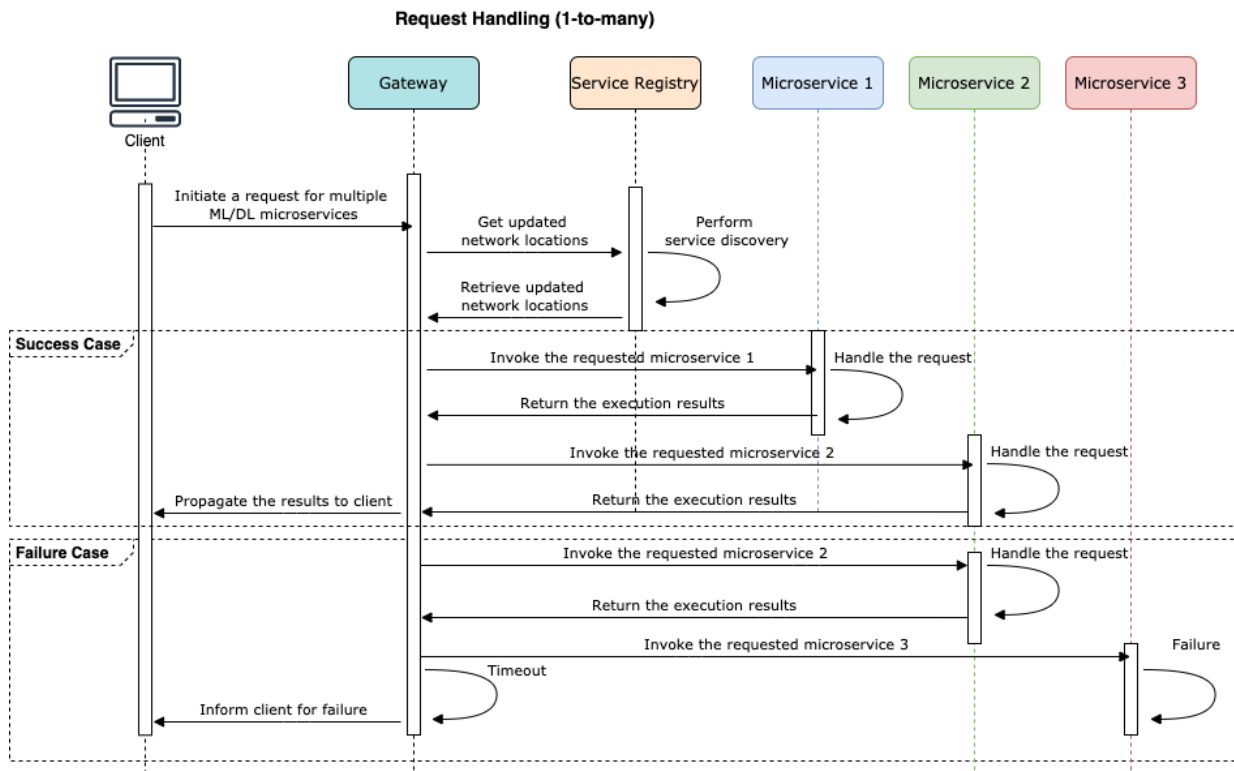


Figure 3-4: Request Handling (1-to-many)

3.3.2 Service Registry

3.3.2.1 Self-registration

The specific functionality handles the self-registration of the ML/DL microservice in the Service Registry of the INFINITECH Open API Gateway. The prerequisite for the registration is the integration of the registry client in the microservice implementation. During the microservice start-up, the microservice is self-registered to the Service Registry.

Table 5: Service Registry – Self-registration

Stakeholders involved:	ML/DL microservice, Service Registry
Pre-conditions:	1. The existence of a ML/DL microservice that has integrated the registry client and its implementation
Post-conditions:	1. The ML/DL microservice is registered in the Service Registry
Data Attributes	Upon self-registration, the microservice needs to send the following information to the service registry: <ol style="list-style-type: none"> 1. The name that the microservice will use as a unique identifier to register to registry 2. The host of the microservice

	<ol style="list-style-type: none"> 3. The port of the microservice 4. The health-check endpoint 5. The health-check interval
Normal Flow	<ol style="list-style-type: none"> 1. The ML/DL microservice is started. The integrated registry client is invoked. 2. The registry client communicates with the Service Registry and self-registers the microservice instance.
Pass Metrics	<ol style="list-style-type: none"> 1. The ML/DL microservice is successfully registered in the Service Registry and it is ready to receive new requests
Fail Metrics	<ol style="list-style-type: none"> 1. The ML/DL microservice failed to register in the Service Registry

Self-registration

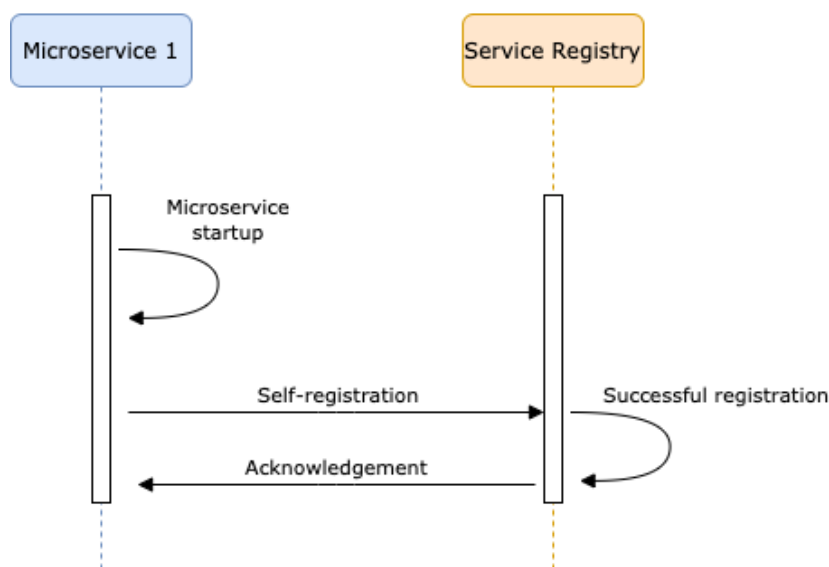


Figure 3-5: Self-registration sequence diagram

3.3.2.2 Self-deregistration

The specific functionality handles the self-deregistration of the ML/DL microservice in the Service Registry of the INFINITECH Open API Gateway. As a prerequisite, the registry client should be integrated with the microservice implementation. During the microservice shutdown, the microservice is self-deregistered from the Service Registry.

Table 6: Service Registry – Self-deregistration

Stakeholders involved:	ML/DL microservice, Service Registry
Pre-conditions:	<ol style="list-style-type: none"> 1. The existence of a ML/DL microservice that has integrated the registry client on its implementation and has been successfully registered in the Service Registry.
Post-conditions:	<ol style="list-style-type: none"> 1. The ML/DL microservice is successfully deregistered from the Service Registry

Data Attributes	1. The identifier of the microservice that needs to be deregistered
Normal Flow	<ol style="list-style-type: none"> 1. The ML/DL microservice initiates shutdown. The integrated registry client is invoked. 2. The registry client communicates with the Service Registry and self-deregisters.
Pass Metrics	1. The ML/DL microservice is successfully deregistered from the Service Registry
Fail Metrics	1. The ML/DL microservice failed to deregister from the Service Registry

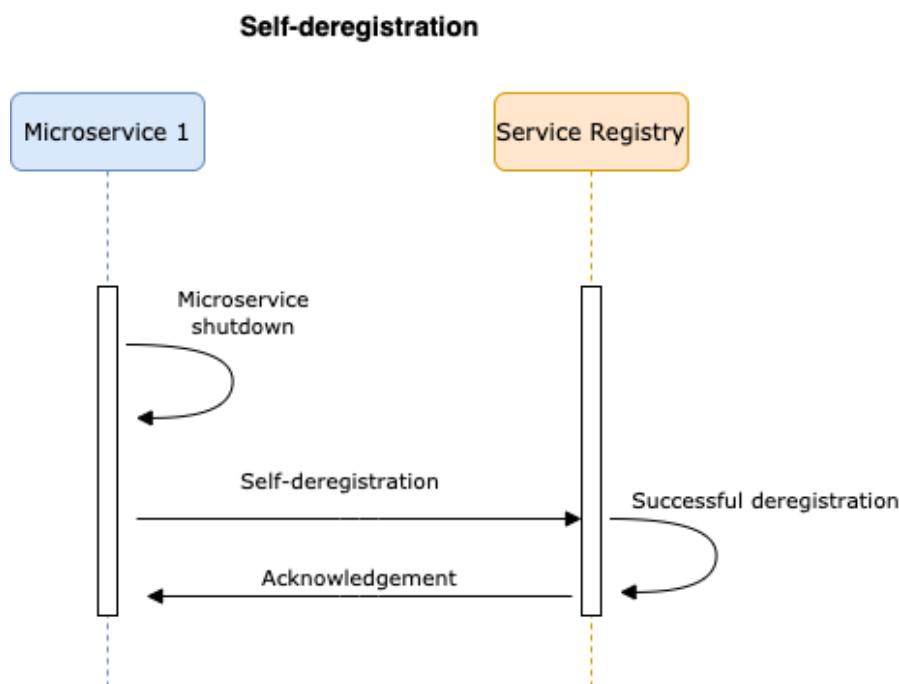


Figure 3-6: Self-deregistration sequence diagram

3.3.2.3 Periodic Health Check

The specific functionality handles the need for periodic health check execution on the registered microservices to ensure their availability. In this context, the Service Registry in every period checks if each registered service is reachable and operating as expected. In the case where a microservice is not responding, then the microservice is marked as in an “Unhealthy” state. At this point, the microservice can properly self-deregister when it comes back into service or it can be removed by the administrator of the INFINITECH Open API Gateway at any time.

Table 7: Service Registry – Periodic Health Check

Stakeholders involved:	ML/DL microservice, Service Registry
Pre-conditions:	<ol style="list-style-type: none"> 1. The existence of a ML/DL microservice that has been successfully registered in the Service Registry and is operating as expected. 2. The existence of a ML/DL microservice that has been successfully registered in the Service Registry and is not responding or is not reachable anymore.

Post-conditions:	1. The operational ML/DL microservice remains in the Service Registry in a “Healthy” state and the unresponsive ML/DL microservice is flagged as “Unhealthy”
Data Attributes	N/A
Normal Flow	<ol style="list-style-type: none"> 1. The Service Registry initiates the health check operation for a specific microservice that is not reachable or irresponsive. 2. The specific microservice is flagged as “Unhealthy” 3. The Service Registry initiates the health check operation for a specific microservice that is operational. 4. The specific microservice remains in the Service Registry in a “Healthy” state
Pass Metrics	1. The operational microservice remains in the Service Registry in a “Healthy” state, while if unresponsive is in an “Unhealthy” state.
Fail Metrics	1. The operational microservice is flagged as “Unhealthy” and / or the unresponsive is flagged as “Healthy”

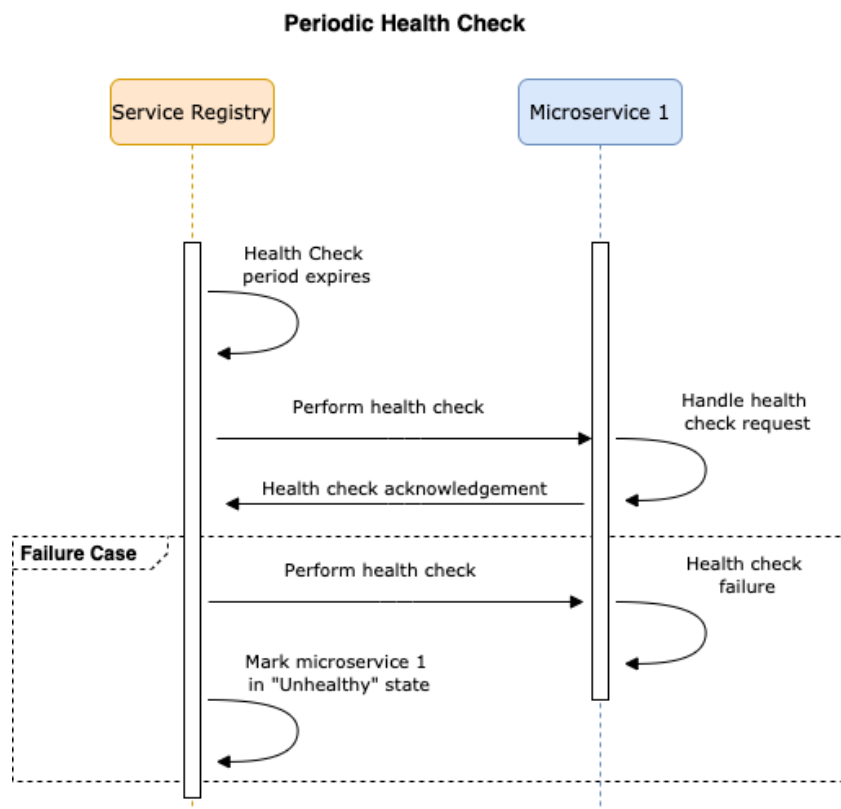


Figure 3-7: Periodic Health Check sequence diagram

4 Implementation of the INFINITECH Open API Gateway

As described in section 3.2, the INFINITECH Open API Gateway adopts a modular architecture and is composed of two core modules, namely the Gateway and the Service Registry. The integration of those two modules formulates the overall solution that provides the single-entry point for the underlying added-value analytics functionalities which are made available in the form of microservices. The implementation of both modules is based on the formulated design specifications which are also documented in section 3.2 of the current deliverable.

In the following subsections, the implementation details of the first prototype version of the INFINITECH Open API Gateway are documented in detail. In particular, for each module the functionalities which were implemented in this first prototype version are documented providing an overview of the delivered module along with the details on how these modules are interacting in order to provide the end-to-end functionalities of the INFINITECH Open API Gateway.

4.1 Gateway

The main role of the Gateway module is to provide the single-entry point for all incoming HTTP requests and to properly handle them by invoking the appropriate microservices as defined in the parameters of the request, ensuring the discovery and intercommunication of the underlying microservices where needed.

To meet its goals, the Gateway is divided into two main components namely the *Gateway Backend* and the *Gateway Frontend* components. The Gateway Backend component undertakes all the core functionalities of the Gateway module by performing all the required request-handling operations. As the Gateway module acts as an advanced reverse proxy which offers additional functionalities than just forwarding the request to the appropriate microservice, multiple operations are performed and orchestrated in the background in order to effectively handle and process all incoming HTTP requests. The Gateway Frontend is providing the single user interface of the INFINITECH Open API Gateway through which the clients can discover and explore the details of the registered microservices, providing them with access to the corresponding Open API documentation of each registered microservice. Both components interact through well-established APIs in order to realise the workflows designed in Section 3 of the current deliverable.

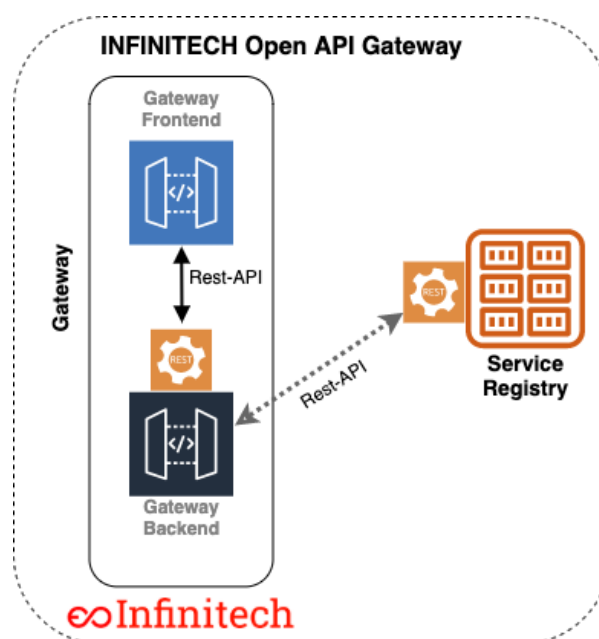


Figure 4-1: Gateway Module Architecture and interactions

The core part of the Gateway Backend component is based on the Spring Cloud Gateway library³ that is offered as part of the Java-based Spring Cloud Ecosystem⁴. The basic concepts of this library are routes and filters. Route constitutes a basic concept of the Spring Cloud Gateway as it defines the routing of a specific request to the underlying microservice upon a successful match. However, before this request is proxied to the appropriate microservice a set of filters are applied which are able to apply modifications in the incoming HTTP request or in the outgoing HTTP response. Filters are divided into the ones applied prior the routing of the request to the requested microservice, hence applying “pre” filter logic, and the ones applied on the microservice’s response, which are applying “post” filter logic.

For the purposes of the INFINITECH Open API Gateway, the specific implementation has been extended taking the implementation of the standard API Gateway offered by Spring Cloud Gateway one step further to being more dynamic with the introduction of:

- a) **Dynamic routing** which is not relying on “hardcoded” routes in the configuration but leverages the dynamic service registry that is offered by the Service Registry module. Dynamic routing is implemented taking into consideration the changes in the availability and network access information of microservices in the cloud and containerised environments.
- b) **Dynamic filtering** which is applied on all incoming HTTP requests with the dynamic configuration of the relevant filters. This dynamic configuration is compiled based on the updated and accurate metadata of each microservice which are dynamically updated and maintained during their self-registration in the Service Registry.

These new features are effectively integrated into the Gateway module’s request handling process, which constitutes the core offering of the Gateway module, extending its capabilities and facilitating the implementation of the designed features of the INFINITECH Open API Gateway. Figure 4-2 displays the implemented request handling process and how routes and filters are leveraged in this process.

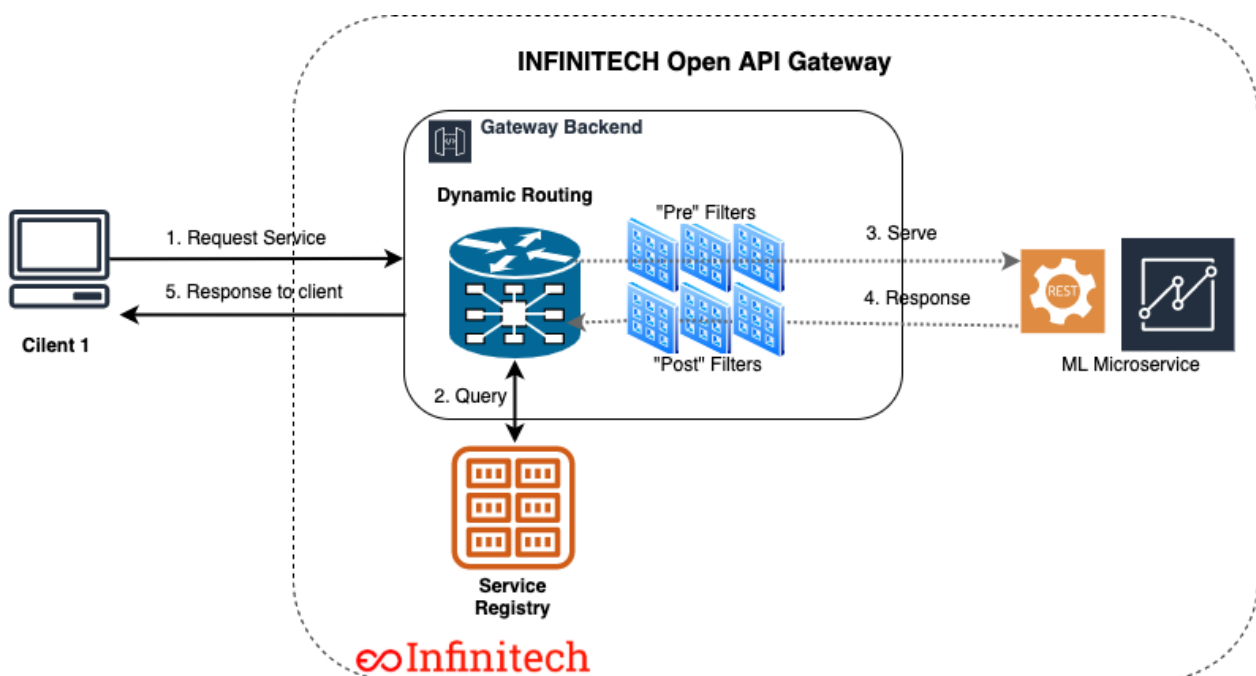


Figure 4-2: Gateway Request Processing Handling

³ Spring Cloud Gateway, <https://spring.io/projects/spring-cloud-gateway>

⁴ Spring Cloud Ecosystem, <https://spring.io/projects/spring-cloud>

4.1.1 Gateway Backend

The Gateway Backend interacts with the Service Registry module via the well-defined APIs which are exposed by the Service Registry in order to automatically and dynamically retrieve the list of registered microservices as well as their appropriate metadata in order to compile the list of dynamic routes for which the Gateway serves as an advanced reverse proxy. This dynamic routing is supplemented with the dynamic filtering with a set of filters that are applied on all incoming HTTP requests. In detail, the Gateway Backend component is leveraging a series of configurable, built-in filters, which are offered by the Spring Cloud Gateway in order to implement the design workflow of the Gateway module. In addition to this, it exploits the offer by the Spring Cloud Gateway library functionality of developing and employing custom filters where needed.

Towards this end, the list of filters currently implemented and utilised in the Gateway Backend component are as follows:

Logging Filter:

The Logging filter is a custom filter which intercepts all incoming HTTP requests and logs the following information: the request type (GET, POST, PUT, DELETE, etc.), the original incoming URI and the transformed URI/service name that will eventually handle the request. It should be noted that for privacy preservation reasons and compliance with GDPR, sensitive information items such as headers, parameters and body are not being logged, since they may contain sensitive information i.e. tokens.

Auth Filter:

The Auth Filter constitutes a custom filter which intercepts all incoming requests and checks whether the request should be authenticated or authorized before reaching the appropriate microservice. During service self-registration, each microservice can define two parameters, namely the authentication and the authorization parameters as described in section 4.2, that control this behaviour. It should be noted that the specific filter is implemented but currently not leveraged by the underlying ML/DL microservices of INFINITECH, since it requires an additional authentication/authorization which is not provided yet.

Rewrite Path Filter:

The Rewrite Path Filter⁵ is a built-in filter which is configured and utilised within the Gateway Backend to appropriately modify the URI of the incoming HTTP request in order to be proxied to the underlying microservice.

Prefix Path Filter:

The Prefix Path Filter⁶ is built-in filter which is configured and utilised within the Gateway Backend to automatically prepend the context path of the underlying microservice in the request URI of the incoming HTTP request, so that the client does not have to include it explicitly when triggering the gateway.

Add Request Header Filter:

The Add Request Header Filter⁷ is a built-in filter which is dynamically configured and utilised within the Gateway Backend to properly set the X-Forwarded-Prefix header, so that when the request arrives to the appropriate service, it knows that it was being reverse proxied before arriving. The specific approach is utilised in the case where the Swagger user interface is utilised to trigger an endpoint of a microservice.

Retry Filter:

⁵ Rewrite Path Filter, <https://docs.spring.io/spring-cloud-gateway/docs/current/reference/html/#the-rewritepath-gatewayfilter-factory>

⁶ Prefix Path Filter, <https://docs.spring.io/spring-cloud-gateway/docs/current/reference/html/#the-prefixpath-gatewayfilter-factory>

⁷ Add Request Header Filter, <https://docs.spring.io/spring-cloud-gateway/docs/current/reference/html/#the-addrequestheader-gatewayfilter-factory>

The Retry Filter⁸ is a built-in filter which is dynamically configured and utilised within the Gateway Backend in order for the gateway to automatically retry upon a failed request. The specific filter is dynamically configured for each separate microservice independently based on the appropriate metadata values of each microservice propagated during self-registration, as described in section 4.2.

Circuit Break Filter:

The Circuit Break Filter⁹ is a built-in filter which is dynamically configured and utilised within the Gateway to apply the Circuit Breaker pattern, a mechanism for properly handling requests to “slow” or “failed” services, by immediately returning an *HTTP 503 - Service Unavailable* code, preventing a further burden to that service. The specific filter constitutes a core filter of the implementation as it provides the required fault-tolerance feature. In order to be effective and efficient, the following parameters are set:

- *slidingWindowType*: The specific parameter configures the type of the sliding window which is used to record the outcome of calls when the CircuitBreaker is closed. Sliding window can either be count-based or time-based.
- *slidingWindowSize*: The specific parameter configures the size of the sliding window which is used to record the outcome of calls when the CircuitBreaker is closed.
- *minimumNumberOfCalls*: The specific parameter configures the minimum number of calls which are required (per sliding window period) before the CircuitBreaker can calculate the error rate or slow call rate. For example, if *minimumNumberOfCalls* is 10, then at least 10 calls must be recorded, before the failure rate can be calculated. If only 9 calls have been recorded the CircuitBreaker will not transition to open even if all 9 calls have failed.
- *permittedNumberOfCallsInHalfOpenState*: The specific parameter configures the number of permitted calls when the CircuitBreaker is half open.
- *failureRateThreshold*: The specific parameter configures the failure rate threshold in percentage. When the failure rate is equal or greater than the threshold, the CircuitBreaker transitions to open and starts short-circuiting calls.
- *waitDurationInOpenState*: The specific parameter configures the time that the CircuitBreaker should wait before transitioning from open to half-open.
- *slowCallDurationThreshold*: The specific parameter configures the duration threshold above which calls are considered as slow and increases the rate of slow calls.
- *slowCallRateThreshold*: The specific parameter configures a threshold in percentage. The CircuitBreaker considers a call as slow when the call duration is greater than *slowCallDurationThreshold*. When the percentage of slow calls is equal or greater than the threshold, the CircuitBreaker transitions to open and starts short-circuiting calls.

The Circuit Break Filter implementation that is utilised for the fault-tolerance feature is supplemented with a lightweight, easy-to-use fault tolerance library named Resilience4j¹⁰ from which the Time Limiter filter is used to force a timeout in the requests if they exceed a specific threshold.

Request Rate Limiter Filter:

The Request Rate Limiter Filter is a built-in filter which is dynamically configured and utilised within the Gateway Backend to optionally apply a rate limiting in the underlying microservice, in order to determine if the current request is allowed to proceed. If it is not, a status of *HTTP 429 - Too Many Requests* is returned. The algorithm behind the implementation of the rate limiting is the Token Bucket algorithm¹¹ and it uses a

⁸ Retry Filter, <https://docs.spring.io/spring-cloud-gateway/docs/current/reference/html/#the-retry-gatewayfilter-factory>

⁹ Circuit Breaker Filter, <https://docs.spring.io/spring-cloud-gateway/docs/current/reference/html/#spring-cloud-circuitbreaker-filter-factory>

¹⁰ Resilience4j, <https://resilience4j.readme.io/docs>

¹¹ Token Bucket algorithm, https://en.wikipedia.org/wiki/Token_bucket

Redis¹² instance to keep track of the number of requests performed. The specific filter is also dynamically configured for each separate microservice independently based on the appropriate metadata values of each microservice propagated during self-registration, as described in section 4.2.

4.1.2 Gateway Frontend

The Gateway Frontend is offering the single user interface of the INFINITECH Open API Gateway acting as the single point of reference of the documentation of the Open APIs of the registered microservices. The Gateway Frontend leverages the embracement of the Open API specification in order to present and publish the documentation of the offered Open APIs of the microservices enabling their easy and effortless discovery and consumption from the clients of the INFINITECH Open API Gateway. In particular, the Gateway Frontend acts as the mediator between the clients of the INFINITECH Open API Gateway and the provided by the registered microservices documentation of their Open APIs by providing the means for all registered microservices to publish their Open API documentation.

To this end, the purpose of the Gateway Frontend is two-fold: a) to display the list of registered microservices and their relevant metadata or properties and b) to display Open API documentation of each registered microservice. It should be noted at this point that the Gateway Frontend is not responsible for the generation or undertaking the compilation of the aforementioned Open API documentation for each microservice, as this is provided by each microservice during their self-registration, but only for the publishing of the provided documentation.

The Gateway Frontend constitutes a Single Page Application (SPA) which is based on the React JS library¹³ displaying the dynamic list of the registered microservices. In particular, the well-established Material-UI framework¹⁴ is leveraged which offers a novel framework for the design and implementation of the user-friendly and fast user interfaces.

For each microservice, the following information is displayed:

- *Service Name*: The name of the registered microservice as provided during its self-registration.
- *Authorization*: A Boolean value that illustrates the existence of an authorisation restriction of the specific microservice.
- *Authentication*: A Boolean value that illustrates the existence of an authentication restriction of the specific microservice.
- *Open-API-Swagger Link*: A link URL to Open API documentation of the Open APIs of the specific microservice which is based on Swagger.

By navigating to a specific microservice of the list, the user is able to access the Open API documentation of the Open APIs via two different options. The first option generates a redirection to the selected Open API documentation where the user can read all the relevant information of the provided Open APIs. The second option displays the Open API documentation in an integrated manner by expanding the list item to display the relevant documentation while collapsing the rest of the items of the list.

The Gateway Frontend is composed of a set of functions which are utilised in order to collect, aggregate and display the dynamic list of the registered microservices along with their properties and metadata. The list of implemented functions are as follows:

- `fetchMicroservicesData(endpoint: Object)`: The specific function fetches microservices' information and attributes via the `'/actuator/gateway/routes'` endpoint.

¹² Redis, <https://redis.io/>

¹³ React, <https://reactjs.org/>

¹⁴ Material UI, <https://material-ui.com/>

- `createTable(headers: Object)` : The specific function undertakes the creation of the schema, headers, dynamic pagination, searching bar and all the functionalities which are supported by the final fulfilled Table.
- `dataTransformation(res: Object)` : The specific function transforms the fetched microservices' data into usable by user interface format before ingesting it into the Table.
- `dataParsing(data: Object)` : The specific function undertakes the parsing of the data into the Table and makes them reachable by the user.
- `swaggerIntegration(swaggerPath: Object)` : The specific function utilises an iframe component and custom Material Table's actions to integrate the OpenAPI-Swagger Documentation Page of the desired microservice.

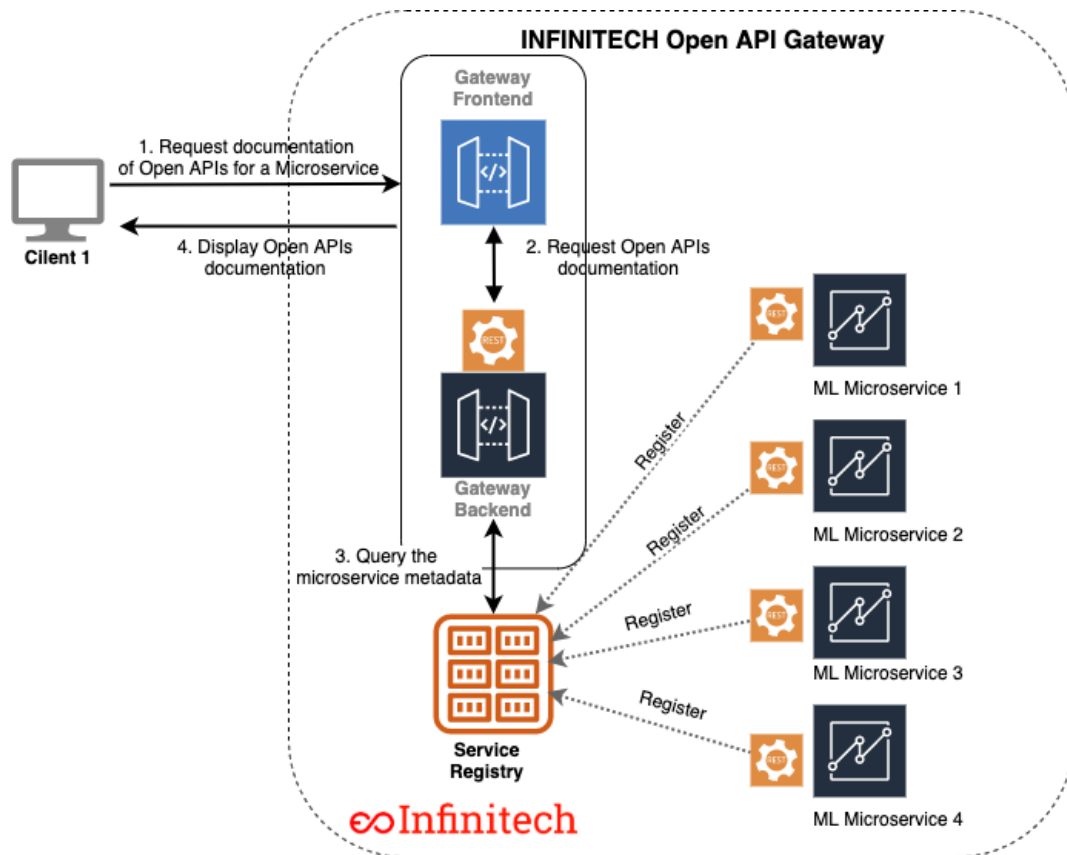


Figure 4-3: Discovery of microservices' Open API documentation

4.2 Service Registry

The main role of the Service Registry is to provide the database that effectively maintains the updated network locations of the registered microservices, offering the mechanism for self-registration and deregistration and the sophisticated mechanism that performs health check operations on these microservices.

Towards this end, the dominant open-source tool named Consul¹⁵ is leveraged that provides out-of-the-box service registry, service discovery and health checking. Service Registry is realised by properly configuring Consul and integrated it with the Gateway module via well-defined API interfaces in order to realise the workflows designed in Section 3 of the current deliverable. In particular, the Service Registry, which is based on the design specifications that are documented in Section 3 also, is complementing the Gateway module

¹⁵ Consul, <https://www.consul.io/>

by supporting the operations performed by the Gateway covering the service discovery aspects of the solution.

Consul undertakes the critical part of the integration of the microservices and provides the required service discovery operations. However, the successful registration and utilisation of any candidate microservice adheres to a specific set of rules and requirements that should be met. The integration requirements imposed by the Service Registry are as follows:

- *Self-registration and self-deregistration*: Each microservice should be able to self-register and deregister to the Service Registry utilising the service registry mechanism provided by Consul.
- *Required metadata on registration*: During service registration to consul, the following metadata should be provided:
 - *Authentication* (Boolean): It indicates if anyone accessing this service should be authenticated.
 - *Authorization* (Boolean): It indicates if the entity accessing the specific microservice should be authorized to do so.
 - *ContextPath* (String): The context path of the microservice. It can be an empty string "", or something like "/microservice1", "/linear-regression-model", etc.
 - *SwaggerPath* (String - Optional): The path holding the JSON representation of the Open API documentation.
 - *rateLimiterEnabled* (Boolean - Optional): It indicates if rate limiting should be applied. Default is true.
 - *rateLimiterScope* (String - Optional): Defines the rate limiting scope. Accepted values are "global" or "user", default value is "global". Global means that rate limiting will be applied for all requests, while user means that rate limiting will be applied per user (requires an authentication/authorization mechanism to be in place).
 - *rateLimiterReplenishRate* (Integer - Optional): The number of requests that are being replenished every second. Default is 100.
 - *rateLimiterBurstCapacity* (Integer - Optional): The total number of requests that can be accepted per second. Default is 100.
 - *rateLimiterRequestedTokens* (Integer - Optional): The cost of each request. Default is 1.
 - *retryAttempts* (Integer - Optional): The number of times the request will be retried in case of failure. Default is 3.
 - *retryFirstBackoff* (String - Optional): The amount of time after which the first retry attempt will take place. Default is 10ms.
 - *retryMaxBackoff* (String - Optional): The maximum backoff to apply. Default is 50ms.
 - *retryBackoffFactor* (Integer - Optional): The backoff retry factor. Default is 2.
 - *retryBackoffBasedOnPreviousValue* (Boolean – Optional): If false, retries are performed after a backoff interval of $\max(\text{retryFirstBackoff} * (\text{retryBackoffFactor} ^ n), \text{retryMaxBackoff})$, where n is the iteration. If true, the backoff is calculated by using $\max(\text{prevBackoff} * \text{retryBackoffFactor}, \text{retryMaxBackoff})$. Default is false.
- *Unique service name*: During registration, the name of the microservice under registration should be unique, otherwise the problem of one microservice overriding another microservice may occur or multiple microservices may be considered as multiple instances of the same microservice which will inevitably lead to an unwanted load balancing.
- *Microservices wrappers*: As per the design specifications of the INFINITECH Open API Gateway, only microservices can be registered. Hence, in the case of ML/DL algorithms each algorithm should be wrapped as a microservice using a corresponding framework, such as Spring Boot for Java, FastAPI for Python, and all the algorithm's functionalities (e.g. train, test, etc) should be exposed as endpoints.
- *Open API specification*: The INFINITECH Open API Gateway embraces Open APIs. While the publishing of the documentation of the APIs of the microservice is optional, if the specific feature is exploited each microservice should document their endpoints using Open API specification.

- *Invoking via the Swagger UI:* In the case where the triggering of the microservice's endpoints is allowed also through the Swagger User Interface then the microservice should be configured properly, as if it was to be behind a reverse proxy, in order to handle X-Forwarded-For and X-Forwarded-Prefix headers.

Upon successful registration to the service registry, each microservice can be accessed in the following path:

`{GATEWAY_IP:GATEWAY_PORT | GATEWAY_DOMAIN}/{microservice-name}/{endpoint}`

It should be noted that the context-path is not required as the Gateway module adds it automatically based on the provided metadata.

5 Baseline technologies and tools

Updates from D5.10:

This particular section has been updated to reflect the updates on the usage of technologies and tools which have been leveraged for the implementation of the first prototype of the INFINITECH Open API Gateway.

During the design phase of the INFINITECH Open API Gateway, several high-potential technologies and tools were investigated and an evaluation of multiple aspects of each candidate technology and tool was conducted. The basic criteria during the selection process were as follows:

- The offered functionalities of each technology and tool, as well as their relevance to the aspired functionalities of the INFINITECH Open API Gateway.
- The applicability of the offered functionalities to the designed use cases for the INFINITECH Open API Gateway.
- The level of maturity of each technology and tool.
- The effort for the implementation of the various components with the specific technologies and tools.
- The integration options and the level of compatibility between them.

For each of the two modules incorporated in the INFINITECH Open API Gateway, different technologies, tools and frameworks were selected as a result of this selection process. During the implementation phase of the first prototype of the INFINITECH Open API Gateway, the list of technologies, tools and frameworks has been updated and further extended with additional ones that were evaluated as the most appropriate ones for the realisation of the INFINITECH Open API Gateway's features.

In detail, the Gateway module that constitutes the core module of the component's architecture is based on the Java programming language and specifically Java version 11. To facilitate the implementation of the module, the dominant open-source Java-based Spring Boot Framework⁴ is utilised. Spring Boot is a powerful lightweight application framework that facilitates the implementation of Java-based enterprise applications in an efficient and modular manner. Spring Boot offers out-of-the-box functionalities capable of supporting the application development such as the embedded Netty Server and easy integration with a large variety of Java-based libraries and technologies. The already embedded and well-established open-source Netty Server¹⁶ is leveraged in the implementation of the module. On top of the Spring Boot Framework, the Spring Cloud Gateway³ library is leveraged. Spring Cloud Gateway provides the means to implement an API Gateway utilising the Spring Boot Framework, offering several functionalities such as flexible routing and request handling, as well as cross-cutting concerns such as security, resiliency and monitoring. With regard to the Circuit Breaker pattern implementation, the Spring Cloud Circuit Breaker⁹ is exploited in order to provide the required abstraction layer and out-of-the-box integration with multiple open-source implementations such as the Resilience4j¹⁰ which is also exploited for the fault-tolerance aspects and the forced timeout functionality. Additionally, there is a compatibility a number of complementary libraries such as the Spring Data Redis Reactive¹⁷ that enables the integration with the Redis key-value store and SpringDoc Open API¹⁸ which enables the automatic generation of an Open API documentation based on the provided YAML or JSON source. The user interface of the Gateway module is built on top of Node.js¹⁹, the dominant JavaScript framework, React¹³ which is one of the most commonly used Single Page Applications frameworks and the complementary libraries Material UI¹⁴ and Material Table²⁰.

¹⁶ Netty, <https://netty.io/>

¹⁷ Spring Data Redis Reactive, <https://spring.io/projects/spring-data-redis>

¹⁸ SpringDoc. Open API, <https://github.com/springdoc/springdoc-openapi>

¹⁹ Node.js, <https://nodejs.org/en/>

²⁰ Material Table, <https://material-table.com/#/>

As explained in previous sections, the Service Registry module is based on Consul¹⁵ which is the well-established open-source solution for service discovery and health checking. Consul provides a sophisticated service registry that enables the self-registration and self-deregistration of microservices, maintaining their dynamic network location and providing simplified service discovery via an open and extensible API and an HTTP interface. Furthermore, it provides real-time health monitoring in order to track the availability and operational status of the registered microservices. Additionally, Consul provides the registration client in multiple programming languages that is integrated within the implementation of each microservice. The registration client undertakes the responsibility of self-registering the respective microservice during the microservice start-up process and the self-deregistration of the microservice during the microservice shutdown process.

In addition to this, the integration of the API Gateway with the Service Registry is enabled with the Spring Cloud Consul²¹ that covers the integration aspects of Spring Boot applications with Consul.

The following table summarises the main technologies and tools which are utilised in the implementation of the INFINITECH Open API Gateway:

Table 8: INFINITECH Open API Gateway list of technologies

Module Name	Software Artefact Name
<i>Gateway</i>	<ul style="list-style-type: none"> • Java version 11 • Spring Boot Framework • Netty server • Spring Cloud Gateway • Spring Cloud Circuit Breaker • Spring Cloud Consul • Resilience4J • Spring Data Redis Reactive • SpringDoc Open API • Node.js • React • Material UI • Material Table
<i>Service Registry</i>	<ul style="list-style-type: none"> • Consul

²¹ Spring Cloud Consul, <https://cloud.spring.io/spring-cloud-consul/reference/html/>

6 Conclusions

The purpose of the deliverable at hand, entitled “D5.11 - “Data Management Workbench and Open APIs - II”, was to report the updated outcomes of the work performed within the context of Task 5.5 “OpenAPI for Analytics and Integrated BigData/AI WorkBench” of WP5. To this end, the deliverable included the updated documentation that supplemented the information included in the previous iteration, namely deliverable D5.10, with regard to the challenges of accessing the microservices implementations of added-value offerings in INFINITECH, the detailed design specifications of the INFINITECH Open API Gateway component and the technical details of its first prototype version that eliminates these challenges with a sophisticated (high-TRL) solution that has potential for novel extensibility.

In detail, at first the deliverable reported the results of a thorough analysis of the two most common approaches, namely the direct client-to-microservices pattern and the API Gateway pattern, that can be used to solve the problem of accessing underlying microservice-based added-value functionalities from the clients. During this analysis, the list of advantages and disadvantages of each approach were elaborated and the reasons for selecting the API Gateway pattern to build the appropriate solution within the context of INFINITECH was clarified. It should be noted that the results remained unchanged from the previous iteration of the deliverable and they were reported here for coherency reasons.

In addition to the analysis of the two approaches, the deliverable documented in detail the design decisions and design specifications of the INFINITECH Open API Gateway component. Within this context, the core design decisions that were taken during the design phase in order to formulate the core aspects of the INFINITECH Open API Gateway were documented. These decisions were the basis for the formulation of the design specifications and the main functionalities that are offered by the INFINITECH Open API Gateway, as well as the modular architecture of the designed solution. The modular architecture is composed of two core modules, namely the Gateway and the Service Registry. Each module has a solid role and operates under a clear context and is assigned with a subset of functionalities from the overall list of functionalities of the INFINITECH Open API Gateway. The design specifications were supplemented with the list of supported use cases and the detailed sequence diagrams that depict the interactions between the modules and the clients of the INFINITECH Open API Gateway. The design specifications remained also unchanged from the previous iteration of the deliverable and they were reported here for coherency reasons.

The deliverable introduces the detailed technical specifications of the first prototype version of the INFINITECH Open API Gateway. In particular, for each module of the component the deliverable presented the details of their implemented functionalities in accordance with their design specifications, as well as the technical details of the integration of the two modules in order to formulate the integrated solution of the INFINITECH Open API Gateway. The Gateway module is formulated by the Gateway Backend that undertakes all the backend functionalities of the INFINITECH Open API Gateway and the Gateway Frontend that constitutes the single user interface of the solution for the discovery of the documentation of the registered microservices and their respective Open APIs. The Service Registry module undertakes the implementation of the service registry, service discovery and health checking functionalities of the solution.

Finally, the deliverable presented the updated list of baseline technologies and tools which were exploited and combined for the implementation of the INFINITECH Open API Gateway component. The list is composed of mature and well-established technologies and tools which were leveraged and smoothly integrated during the implementation phase, effectively addressing the requirements elicited by the defined design specifications.

It should be noted at this point that the current deliverable constitutes the second report of the work performed within the context of Task 5.5. In accordance with the INFINITECH Description of Action, the final iteration will be delivered on M30 with deliverable D5.12. The outcomes of this deliverable will drive the implementation activities of the INFINITECH Open API Gateway component. Nevertheless, as the project evolves and the implementation of the specific component is a living process that will last until M30, updates and refinements will be introduced based on the feedback that will be collected, as well as any requirements

that may arise. The final and complete documentation of both the design specifications as well as the implementation details will be documented in the upcoming final version of this deliverable.

Table 9: Conclusions (TASK Objectives with Deliverable achievements)

Objectives	Comment
<i>Exploit the API Gateway pattern offerings.</i>	The proposed solution successfully exploits the offerings of the API Gateway pattern related to the effective and efficient service discovery and straight-forward access to the dynamically deployed microservices, the elimination of multiple server round trips and the hiding of the microservice’s architecture details and complexity from the clients.
<i>Provide a single-point-entry for the added-value offerings functionalities of INFINITECH.</i>	The INFINITECH Open API Gateway delivers the required single entry-point for all the ML/DL microservices and their exposed Open APIs that will be available in INFINITECH. Furthermore, it provides the basis for the support of additional microservice-based added-value functionalities in INFINITECH.
<i>Design and deliver the required solution that enables the discovery and consumption of the ML/DL microservices of INFINITECH and their exposed Open APIs.</i>	The detailed design specifications of the INFINITECH Open API Gateway are delivered with the current deliverable; its offered functionalities successfully cover both the discovery and the consumption of the underlying ML/DL microservices of INFINITECH as well as of their exposed Open APIs. The first prototype version of the INFINITECH Open API Gateway is delivered with the current deliverable while its final implementation will be delivered on M30.

Table 10: (map TASK KPI with Deliverable achievements)

KPI	Comment
<i>API Gateways available for enabling the access to INFINITECH ML/DL microservices.</i>	Target Value = 1 The specific KPI is successfully achieved with the presented solution, namely the INFINITECH Open API Gateway, whose detailed design specifications and first prototype implementation were documented in the current deliverable.
<i>Coverage of ML/DL microservices exposing Open APIs.</i>	Target Value = 100% Per the design specifications documented in the current deliverable, the INFINITECH Open API Gateway is capable of covering all ML/DL microservices that are exposing Open APIs within the context of INFINITECH, hence the coverage is 100%.
<i>Number of services exposed to the clients of the API Gateway.</i>	Target Value >= 4 The INFINITECH Open API Gateway exposes five main services: <ul style="list-style-type: none"> a) the request handling service, b) the discovery service of the Open APIs of the ML/DL microservices, c) the ML/DL microservices service registry, d) the self-registration service, and e) the self-deregistration service.

Appendix A: Literature

- [1] “Microservices.io,” [Online]. Available: <https://microservices.io/>. [Accessed 20 October 2020].
- [2] M. Dudjak and G. Martinović, “An API-first methodology for designing a microservice-based Backend as a Service platform.,” *Information Technology and Control*, vol. 49, no. 2, pp. 206-223, 2020.
- [3] F. Montesi and J. Weber, “Circuit breakers, discovery, and API gateways in microservices,” 2016.
- [4] “Self-Registration,” 2020. [Online]. Available: <https://microservices.io/patterns/self-registration.html>. [Accessed 10 October 2020].
- [5] “3rd-Party Registration,” 2020. [Online]. Available: <https://microservices.io/patterns/3rd-party-registration.html>. [Accessed 15 October 2020].
- [6] D. Taibi, V. Lenarduzzi and C. Pahl, “Architectural patterns for microservices: a systematic mapping study,” 2018.
- [7] “Open API specification,” 2020. [Online]. Available: [https://swagger.io/specification/#:~:text=The%20OpenAPI%20Specification%20\(OAS\)%20defines,or%20through%20network%20traffic%20inspection..](https://swagger.io/specification/#:~:text=The%20OpenAPI%20Specification%20(OAS)%20defines,or%20through%20network%20traffic%20inspection..) [Accessed 20 October 2020].
- [8] “Circuit Breaker Pattern,” 2020. [Online]. Available: <https://docs.microsoft.com/en-us/azure/architecture/patterns/circuit-breaker>. [Accessed 10 October 2020].