



HORIZON 2020 - ICT-14-2016-1

## AEGIS

Advanced Big Data Value Chains for Public Safety and Personal Security

### WP3 - System Requirements, User stories, Architecture and MicroServices



## D3.5 – Architecture and Revised Components, Microservices and APIs Designs v4.00

Version 1.0

**Due date:** 31/01/2019

**Delivery Date:** 31/01/2019

**Author(s):** Dimitrios Miltiadou, Konstantinos Perakis, Stamatis Pitsios (UBITECH), Maurizio Megliola, Ambra De Bonis (GFT), Mahmoud Ismail (KTH), Yury Glikman, Fabian Kirstein, Fritz Meiners (Fraunhofer), Giannis Tsapelas, Panagiotis Kokkinakos, Spiros Mouzakis, Christos Botsikas, Michael Kontoulis (NTUA), Sotirios Koussouris, George Bikas (SUITE5)

**Editor:** Dimitrios Miltiadou (UBITECH)

**Lead Beneficiary of Deliverable:** UBITECH

**Dissemination level:** Public

**Nature of the Deliverable:** Report

**Internal Reviewers:** Marios Phinikettos (SUITE5), Charalampos Kladouhos (Konkat), Georgios Avouris (Konkat)

VERSIONING (ONLY MAJOR VERSIONS)			
VERSION	DATE	NAME, ORGANIZATION	DESCRIPTION OF THE NEW VERSION
0.1	10/01/2019	UBITECH	Toc
0.2	11/01/2011	UBITECH	CONTRIBUTION TO SECTION 1, 2
0.3	14/01/2019	UBITECH	CONTRIBUTION TO SECTION 2
0.4	15/01/2019	UBITECH, KTH, NTUA, SUITE5, FRAUNHOFER, GFT	CONTRIBUTION TO SECTION 3,4
0.5	18/01/2019	UBITECH	CONTRIBUTION TO SECTION 3,4
0.6	23/01/2019	KTH, NTUA, SUITE5, FRAUNHOFER, GFT	CONTRIBUTION TO SECTION 3,4
0.7	25/01/2019	UBITECH	CONTRIBUTION TO SECTION 5
0.8	28/01/2019	UBITECH	REVIEW READY VERSION
0.9_KONKAT	29/01/2019	KONKAT	INTERNAL REVIEW
0.90_SUITE5	29/01/2019	SUITE5	INTERNAL REVIEW
1.0	30/01/2019	UBITECH	FINAL VERSION

**Remark:** The versioning is only for the word documents in the formation phase and should be kept internally. Please delete the versioning before creating the final pdf that goes to the commission. It can be provided to the commission on request. Please document only major versions and such versions that indicate through the versioning, who (person and which partner) has contributed/was responsible for the different chapter, if this is feasible.

## EXPLANATIONS FOR FRONTPAGE

**Author(s):** Name(s) of the person(s) having generated the Foreground respectively having written the content of the report/document. In case the report is a summary of Foreground generated by other individuals, the latter have to be indicated by name and partner whose employees he/she is. List them alphabetically.

**Editor:** Only one. As formal editorial name only one main author as responsible quality manager in case of written reports: Name the person and the name of the partner whose employee the Editor is. For the avoidance of doubt, editing only does not qualify for generating Foreground; however, an individual may be an Author – if he has generated the Foreground - as well as an Editor – if he also edits the report on its own Foreground.

**Lead Beneficiary of Deliverable:** Only one. Identifies name of the partner that is responsible for the Deliverable according to the AEGIS DOW. The lead beneficiary partner should be listed on the frontpage as Authors and Partner. If not, that would require an explanation.

**Internal Reviewers:** These should be a minimum of two persons. They should not belong to the authors. They should be any employees of the remaining partners of the consortium, not directly involved in that

deliverable, but should be competent in reviewing the content of the deliverable. Typically this review includes: Identifying typos, Identifying syntax & other grammatical errors, Altering content, Adding or deleting content.

**AEGIS KEY FACTS**

<b>Topic:</b>	ICT-14-2016 - Big Data PPP: cross-sectorial and cross-lingual data integration and experimentation
<b>Type of Action:</b>	Innovation Action
<b>Project start:</b>	1 January 2017
<b>Duration:</b>	30 months from <b>01.01.2017</b> to <b>30.06.2019</b> (Article 3 GA)
<b>Project Coordinator:</b>	Fraunhofer
<b>Consortium:</b>	10 organizations from 8 EU member states

**AEGIS PARTNERS**

<b>Fraunhofer</b>	Fraunhofer-Gesellschaft zur Förderung der angewandten Forschung e.V.
<b>GFT</b>	GFT Italia SRL
<b>KTH</b>	Kungliga Tekniska högskolan
<b>UBITECH</b>	UBITECH Limited
<b>VIF</b>	Kompetenzzentrum - Das virtuelle Fahrzeug, Forschungsgesellschaft-GmbH
<b>NTUA</b>	National Technical University of Athens – NTUA
<b>EPFL</b>	École polytechnique fédérale de Lausanne
<b>SUITE5</b>	SUITE5 Limited
<b>KONKAT</b>	ANONYMOS ETAIREIA KATASKEVON-TECHNIKON ERGON, EMPORIKON, VIOMICHANIKONKAI NAUTILIAKON EPICHEIRISEON KON'KAT
<b>HDIA</b>	HDI Assicurazioni S.P.A

**Disclaimer:** AEGIS is a project co-funded by the European Commission under the Horizon 2020 Programme (H2020-ICT-2016) under Grant Agreement No. 732189 and is contributing to the BDV-PPP of the European Commission.

The information and views set out in this publication are those of the author(s) and do not necessarily reflect the official opinion of the European Communities. Neither the European Union institutions and bodies nor any person acting on their behalf may be held responsible for the use which may be made of the information contained therein.

© Copyright in this document remains vested with the AEGIS Partners

## EXECUTIVE SUMMARY

The document at hand, entitled “Architecture and Revised Components, Microservices and APIs Design - v4.00” documents the efforts within the context of the tasks 3.1, 3.2, 3.3, 3.4 and 3.5 of WP3. Towards this end, the scope of D3.5 is to provide final detailed documentation with regard to the high-level and technical architecture of the platform, the components of the platform, as well as the workflows of the platform. The deliverable, which is the final deliverable of WP3, builds directly on top of the outcomes and knowledge extracted from D3.4 with the aim of providing a complete documentation containing all the necessary updates and modifications that have been performed during WP3.

More specifically, the objectives of the deliverable D3.5 are as follows:

- Document the final version of the high-level architecture of the AEGIS platform. Within this context, the document at hand describes all the functionalities of the final version of the AEGIS platform, as well as a comprehensive description of the components of the platform, focusing on the role within the platform and their positioning in the architecture. In this description the core functionalities of each component are also highlighted.
- Document the final version of the technical architecture in which the functional decomposition of the components, the relationship among them and the data flow are presented.
- Provide the detailed final documentation for each component with regard to their design and the specifications.
- Present the final list of functionalities of each component of the platform and the list of technologies and tools utilised towards the implementation of the aforementioned functionalities.
- Document the list of microservices that were designed in the context of each component in order to enable the realisation of the functionalities of each component of the platform.
- Document the technical interfaces and exposed outcomes offered by each component in detail that were designed in order to enable the smooth integration of the components as well as the execution of the workflows of the platform.
- Present the final BPMN diagrams that are illustrating how, from a user perspective, the AEGIS platform’s workflows are enabled and the AEGIS platform’s components interact on a high-level.

The outcomes of this deliverable will drive the final implementation activities of the project towards the implementation of the AEGIS Platform Release 4.00 which constitutes as the final release of the platform.

## Table of Contents

<b>EXPLANATIONS FOR FRONTPAGE.....</b>	<b>2</b>
<b>AEGIS KEY FACTS .....</b>	<b>4</b>
<b>AEGIS PARTNERS.....</b>	<b>4</b>
<b>EXECUTIVE SUMMARY.....</b>	<b>5</b>
<b>LIST OF FIGURES .....</b>	<b>8</b>
<b>LIST OF TABLES .....</b>	<b>9</b>
<b>ABBREVIATIONS .....</b>	<b>11</b>
<b>1. INTRODUCTION.....</b>	<b>13</b>
1.1. OBJECTIVE OF THE DELIVERABLE .....	13
1.2. INSIGHTS FROM OTHER TASKS AND DELIVERABLES .....	13
1.3. STRUCTURE .....	14
<b>2. AEGIS ARCHITECTURE.....</b>	<b>15</b>
2.1. HIGH LEVEL ARCHITECTURE.....	15
2.2. TECHNICAL ARCHITECTURE.....	17
2.3. AEGIS INTEGRATED NOTEBOOKS .....	18
<b>3. AEGIS COMPONENTS AND APIS SPECIFICATIONS .....</b>	<b>20</b>
3.1. DATA HARVESTER .....	20
3.1.1. Overview.....	20
3.1.2. Supported Data Sources .....	22
3.1.3. List of microservices.....	23
3.1.4. Technologies to be used.....	24
3.1.5. APIs and exposed outcomes.....	25
3.2. CLEANSING TOOL.....	26
3.2.1. Overview.....	26
3.2.2. List of microservices.....	28
3.2.3. Technologies to be used.....	30
3.2.4. APIs and exposed outcomes.....	30
3.3. ANONYMISATION TOOL.....	40
3.3.1. Overview.....	40
3.3.2. List of microservices.....	41
3.3.3. Technologies to be used.....	41
3.3.4. APIs and exposed outcomes.....	42
3.4. BROKERAGE ENGINE .....	43
3.4.1. Overview.....	43
3.4.2. List of microservices.....	44
3.4.3. Technologies to be used.....	45
3.4.4. APIs and exposed outcomes.....	45
3.5. AEGIS DATA STORE.....	49
3.5.1. Overview.....	49
3.5.2. HopsFS filesystem.....	49
3.5.3. AEGIS Metadata Service .....	52
3.6. AEGIS INTEGRATED SERVICES .....	56
3.6.1. Overview.....	56
3.6.2. List of microservices.....	57
3.6.3. Technologies to be used.....	58
3.6.4. APIs and exposed outcomes.....	59
3.7. QUERY BUILDER .....	60
3.7.1. Overview.....	60
3.7.2. List of microservices.....	61

3.7.3. <i>Technologies to be used</i> .....	63
3.7.4. <i>APIs and exposed outcomes</i> .....	63
3.8. VISUALISER.....	64
3.8.1. <i>Overview</i> .....	64
3.8.2. <i>List of microservices</i> .....	66
3.8.3. <i>Technologies to be used</i> .....	67
3.8.4. <i>APIs and exposed outcomes</i> .....	68
3.9. ALGORITHM EXECUTION CONTAINER.....	68
3.9.1. <i>Overview</i> .....	68
3.9.2. <i>List of microservices</i> .....	70
3.9.3. <i>Technologies to be used</i> .....	70
3.9.4. <i>APIs and exposed outcomes</i> .....	71
3.10. AEGIS FRONT-END.....	71
3.10.1. <i>Overview</i> .....	71
3.10.2. <i>List of microservices</i> .....	74
3.10.3. <i>Technologies to be used</i> .....	74
3.10.4. <i>APIs and exposed outcomes</i> .....	75
3.11. MULTILINGUALISM SUPPORT.....	75
3.11.1. <i>Overview</i> .....	75
3.11.2. <i>Approach</i> .....	75
3.11.3. <i>Technologies to be used</i> .....	76
3.11.4. <i>List of microservices</i> .....	77
3.11.5. <i>APIs and exposed outcomes</i> .....	77
3.12. HOLISTIC SECURITY APPROACH.....	78
3.12.1. <i>Overview</i> .....	78
3.12.2. <i>Technologies to be used</i> .....	80
3.12.3. <i>API</i> .....	81
<b>4. USER INTERACTION WORKFLOWS.....</b>	<b>82</b>
4.1. SIGN-UP AND LOGIN.....	82
4.2. DATA IMPORT.....	82
4.2.1. <i>Importing data for a new dataset</i> .....	82
4.2.2. <i>Anonymisation workflow</i> .....	83
4.2.3. <i>Data cleansing workflow</i> .....	83
4.3. DATA AND SERVICE EXPLORATION (SEARCH).....	84
4.3.1. <i>From the main AEGIS platform</i> .....	84
4.3.2. <i>Using query builder</i> .....	84
4.4. DATA EXPORT FROM AEGIS.....	86
4.5. ARTEFACT SHARING/REUSE.....	86
4.6. SERVICE CREATION.....	87
4.7. SERVICE CONSUMPTION.....	88
<b>5. CONCLUSION.....</b>	<b>89</b>
<b>APPENDIX A: LITERATURE.....</b>	<b>90</b>

**LIST OF FIGURES**

Figure 2-1: AEGIS high-level architecture .....	16
Figure 2-2: AEGIS Technical Architecture .....	19
Figure 3-1: Sequence diagram of the Harvester component .....	21
Figure 3-2: The harvester interface .....	24
Figure 3-3: Harvester Orchestration Concept .....	25
Figure 3-4: Offline data cleansing sequence diagram .....	28
Figure 3-5: Data anonymisation sequence diagram .....	40
Figure 3-6: Brokerage Engine sequence diagram .....	44
Figure 3-7: Query building and execution workflow .....	61
Figure 3-8: Sequence diagram of the visualiser component .....	65
Figure 3-9: Algorithm Execution Container sequence diagram .....	69
Figure 3-10: Mock-up of the AEGIS platform landing page .....	72
Figure 3-11: Mock-up of the main menu of the AEGIS platform .....	73
Figure 4-1: Sign-up and Login workflow .....	82
Figure 4-2: Importing data and metadata and registering them as a part of a new dataset .....	83
Figure 4-3: Data anonymisation workflow .....	83
Figure 4-4: Data cleansing workflow .....	84
Figure 4-5: Data and service exploration workflow .....	84
Figure 4-6: Dataset exploration through query builder workflow .....	85
Figure 4-7: Data acquisition sub-process workflow .....	86
Figure 4-8: Data export workflow .....	86
Figure 4-9: Artefact Sharing Workflow .....	87
Figure 4-10: Service creation workflow .....	88
Figure 4-11: AEGIS Service consumption workflow .....	88



**LIST OF TABLES**

Table 3-1: Harvester list of microservices .....	23
Table 3-2: Data Harvester technical interface .....	26
Table 3-3: Cleansing Tool list of microservices .....	30
Table 3-4: Offline Cleansing tool technical interface .....	31
Table 3-5 Offline Cleansing tool add new provider .....	32
Table 3-6: Offline Cleansing tool update provider .....	32
Table 3-7: Offline Cleansing tool delete provider .....	33
Table 3-8: Offline Cleansing tool new dataset .....	34
Table 3-9: Offline Cleansing tool update an existing dataset .....	34
Table 3-10: Offline Cleansing tool delete dataset .....	35
Table 3-11: Offline Cleansing tool add new variable .....	35
Table 3-12: Offline Cleansing tool update an existing variable .....	36
Table 3-13: Offline Cleansing tool delete variable .....	37
Table 3-14: Offline Cleansing tool add validation rule .....	37
Table 3-15: Offline Cleansing tool delete validation rule .....	38
Table 3-16: Offline Cleansing tool add cleaning rule .....	38
Table 3-17: Offline Cleansing tool delete cleaning rule .....	39
Table 3-18: Offline Cleansing tool update missing value rule .....	39
Table 3-19: Anonymisation Tool list of microservices .....	41
Table 3-20: Anonymisation tool technical interface .....	43
Table 3-21: Brokerage engine list of microservices .....	45
Table 3-22: Brokerage Engine technical interface 1 .....	46
Table 3-23: Brokerage Engine technical interface 2 .....	47
Table 3-24: Brokerage Engine technical interface 3 .....	48
Table 3-25: Brokerage Engine technical interface 4 .....	49

Table 3-26: HopsFS list of microservices .....	50
Table 3-27: AEGIS Data Store technical interface 1 .....	51
Table 3-28: AEGIS Data Store technical interface 2 .....	52
Table 3-29: AEGIS Metadata service list of microservices .....	54
Table 3-30: AEGIS Metadata Service technical interface .....	56
Table 3-31: AEGIS Metadata Service technical interface 2 .....	56
Table 3-32: AEGIS Integrated services list of microservices .....	58
Table 3-33: Query Builder list of microservices .....	62
Table 3-34: Visualiser list of microservices.....	67
Table 3-35: Algorithm Execution Container list of microservices .....	70
Table 3-36: AEGIS Front-End list of microservices.....	74
Table 3-37 - Translation Middleware microservice .....	77
Table 3-38: Translation Middleware technical interface .....	78
Table 3-39: Holistic Security Approach summary.....	80

**ABBREVIATIONS**

API	Application programming interface
BPMN	Business Process Model and Notation
CO	Confidential, only for members of the Consortium (including the Commission Services)
CSS	Cascading Style Sheets
CSV	Comma Separated Value files
D	Deliverable
DLT	Distributed ledger technology
DoW	Description of Work
DPF	Data Policy Framework
FLOSS	Free/Libre Open Source Software
HTML	Hypertext Markup Language
H2020	Horizon 2020 Programme
JSON	JavaScript Object Notation
JWT	JSON Web Token
NLP	Natural language processing
OSS	Open Source Software
PSPS	Public Safety and Personal Security
PU	Public
PM	Person Month
R	Report
RDF	Resource Description Framework
REST	Representational State Transfer
RTD	Research and Development
SQL	Structured Query Language
SSL	Secure Sockets Layer
T	Task

TLS	Transport Layer Security
UI	User Interface
URL	Uniform Resource Locator
WP	Work Package
XML	Extensible Markup Language

## 1. INTRODUCTION

### 1.1. Objective of the deliverable

The scope of D3.5 is to document the efforts within the context of the tasks 3.1, 3.2, 3.3, 3.4 and 3.5 of WP3. Towards this end, the main objective of the current deliverable is to document the AEGIS platform final architecture, as well as to provide the final documentation of the platform's components in terms of design and functionalities, the developed microservices within the context of each component, the technologies that are used and the APIs and exposed outcomes of each component. The document at hand concludes all activities performed in this work package.

The current deliverable aims at providing the updated documentation that supplements the information documented in the deliverable D3.4. Hence, the purpose of the current deliverable is threefold. Firstly, deliverable D3.5 documents the final version of AEGIS platform architecture incorporating all the updates and enhancements that were introduced in the course of the development of the AEGIS platform. Secondly, deliverable D3.5 provides the updated description of the AEGIS platform's components, highlighting all the refinements and improvements in the design and specifications. Thirdly, deliverable D3.5 presents the final AEGIS platform workflows that illustrate the provided functionalities of AEGIS platform. For coherency reasons, the current document builds on top of the information included in the deliverable D3.4, indicating the necessary updates at the end of each section.

The revised information is presented using the approach followed in the previous versions. At first, the final high-level architecture of the AEGIS platform is presented, focusing on the detailed description of the role of each component within the platform, as well as the functionalities undertaken by each component. Additionally, the description is supplemented with the relevant information with regards to the positioning of each component within the architecture. Following the high-level architecture, the AEGIS platform's technical architecture is documented towards the aim of presenting the functional decomposition of the components, the relationship among them and the data flow.

Following the AEGIS platform architecture description, for each component of the platform a detailed description is provided with the complete documentation of the design and functionalities of the component, while also documenting the component's interaction with the rest of the components. Furthermore, for each component the list of designed microservices is documented, the technologies that are utilised in the component implementation are presented and finally the technical interfaces and the exposed outcomes of each component are documented.

Finally, the current deliverable presents the final AEGIS platform workflows in the form of BPMN diagrams in order to provide an overview of the functionalities of the platform focusing on the user perspective.

### 1.2. Insights from other tasks and deliverables

The deliverable builds on top of the work reported in WP3 and WP4. In particular, the work performed in WP3, as reported in D3.1, D3.2, D3.3 and D3.4, provided valuable information concerning the functional and technical requirements collected, the design of the platform's

components, the high-level architecture of the platform, the platform's workflows and all the updates and refinements that were introduced on all these in the course of the development of the platform.

Another useful insight is the work performed within the context of WP4, as reported in D4.1 D4.2 and D4.3, where the first three versions of the platform were delivered. The evaluation and feedback received from the project's demonstrators on all released versions served as the basis upon which the updates and refinements on the platform and the platform's components were built.

### 1.3. Structure

Deliverable D3.5 is organised in five main sections as indicated in the table of contents.

- The first section introduces the deliverable. It documents the objectives of the deliverable and the relation of the current deliverable with the other deliverables by describing how the outcomes of other deliverables and work-packages serves as input to the current deliverable. Finally, a brief description is provided on how the document is structured.
- The second section presents the final high-level architecture of the AEGIS platform, focusing on describing the role and the positioning of each component within the architecture. In addition to this, the final technical architecture of the AEGIS platform is presented, in which the functional decomposition of the components is illustrated, along with their relationships and the respective data flow. Finally, in this section the decision to provide an integrated notebook containing three major components of the AEGIS platform is documented.
- The third section presents the details documentation of the components of the AEGIS platform. Within the context of this section, for each component the final design and functionalities are presented, as well as the list of designed microservices incorporated within each component. Furthermore, for each component the technologies that are utilised for their implementation are presented, as well as the technical interfaces and exposed outcomes are documented. Within each subsection, the updates from deliverable D3.4 are highlighted.
- The fourth section is presenting the BPMN diagrams that correspond to the provided functionalities of the AEGIS platform focusing on the user perspective and on summarising the component interactions in a high-level without the technical details.
- The fifth section concludes the deliverable. It outlines the main findings of the deliverable which will guide the future research and technological efforts of the consortium.

## 2. AEGIS ARCHITECTURE

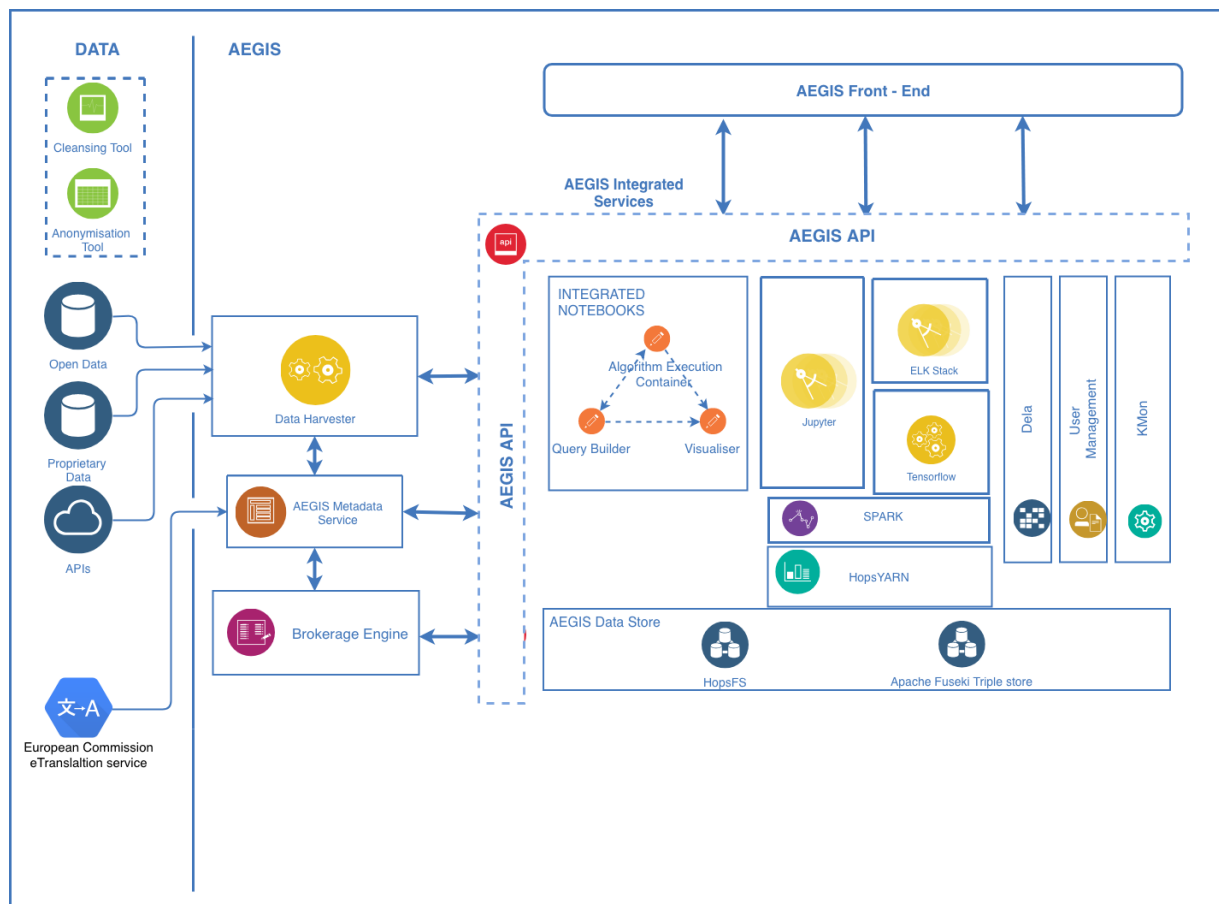
### 2.1. High level architecture

The design of the AEGIS high-level architecture was driven by the results of the thorough analysis of all the technical requirements that was conducted with the aim of addressing the goals and the expectations of the AEGIS stakeholders. Moreover, the design of the AEGIS high-level architecture facilitated the realisation of the designed workflows that enable the data-driven innovation in the PSPS domain as envisioned by the consortium.

The high-level architecture describes the complete lifecycle of the AEGIS supported processes, starting from the transformation and harmonisation of the valuable datasets from the identified data sources, to their semantic annotation and the rich metadata generation, to their exploration with dynamic queries and view creation on top of them, and continuing with advanced analytics including state of the art big data focused algorithm execution and sophisticated processing, complemented by advanced visualizations of the analysis results.

The AEGIS high-level architecture is a modular architecture composed of multiple key components, where each component was designed with a clear business context, scope and set of functionalities. As the project matured after the initial version of the high-level architecture, additional functionalities were designed and introduced in the platform. Moreover, as a result of the comprehensive analysis of the feedback received by the end-users of the platform from the released versions of the platform, a series of adjustments and refinements were introduced in the components of the platform in order to better address the identified requirements, but also to facilitate the implementation of the functionalities of the platform.

Figure 2-1 illustrates the final AEGIS high-level architecture, which incorporates all the adjustments and refinements that were introduced in the course of development of the AEGIS platform. This architecture had driven the implementation and the release of the AEGIS Platform Release 3.00 and will also drive the implementation and the release of the final version of the AEGIS platform, namely the AEGIS Platform Release 4.00.



**Figure 2-1: AEGIS high-level architecture**

Residing at the location of the data, two optional components are offered by the AEGIS ecosystem, namely the **Anonymisation tool** and the **Cleansing tool**. The Anonymisation tool is an offline tool ensuring that sensitive or personal data are not uploaded in the platform and will address the privacy and anonymity requirements by applying a set of anonymisation techniques on the initial dataset. The Cleansing tool provides the necessary cleansing processes with a variety of techniques that will be offered in both offline and online mode (through custom processes incorporated inside the integrated notebooks of the platform) depending on the context of the processes and required corrective actions.

The **Data Harvester** is providing the data entry point to the AEGIS platform offering the transformation, harmonisation and annotation functionalities required within the context of the platform as well as the rich metadata generation for the imported data. In the core of the AEGIS platform lays the **AEGIS Data Store** component, composed by the HopsFS and the Apache Fuseki Triple store. HopsFS is a fast, reliable and scalable distributed file system that undertakes the responsibility for storing the imported datasets, while the Apache Fuseki Triple store is utilised for storing the metadata associated with the imported datasets. The **AEGIS Metadata Service** is responsible for storing the metadata generated using the AEGIS ontology and vocabulary for each dataset, as provided by the Data Harvester, in the Apache Fuseki Triple store.

The **AEGIS Integrated Services** consists of a list of services responsible for the data management and processing within the platform. In addition to the multi-tenant data



management, data exploration, data parallel processing and resource management, these services implement as well the user management and service monitoring aspects of the AEGIS platform. The list of services in AEGIS Integrated Services includes the Jupyter service offering interactive notebooks, the Elasticsearch Logstash Kibana (ELK) stack, the Tensorflow, the Apache Spark and the HopsYARN, as well as the Dela, the User Management and the KMon services.

In addition to the AEGIS Integrated Services, the AEGIS platform incorporates three more components in the form of integrated notebooks using Jupyter, namely the **Query Builder**, the **Algorithm Execution Container** and the **Visualiser** that are supplementing the delivered functionalities of the AEGIS platform. More specifically, the Query Builder is simplifying and empowering the querying capabilities of the platform by providing an intuitive graphical interface for powerful data pre-processing capabilities, data retrieval and view creation on the data in order to generate a new dataset or provide an input to Algorithm Execution Container and Visualiser. The Algorithm Execution Container is enabling the execution of the data analysis algorithms over selected datasets in order to provide the data analysis results in the Visualiser. The Visualiser is the component facilitating the visualisation functionalities of the platform for either the querying and filtering results as produced by the Query Builder or the analysis results as produced by the Algorithm Execution Container.

The **Brokerage Engine** is responsible for access control and recording of actions performed over the artefacts (e.g. datasets) which are placed on the platform and have a price tag (e.g. they are not made available for free). More specifically, the Brokerage Engine is ensuring conformance with the Data Policy Framework of AEGIS while also utilising a distributed ledger supported by a blockchain implementation in order to record all transactions over these artefacts. Finally, the AEGIS Front-End is the component implementing the presentation layer of the platform using an innovative user-friendly interface to enable the easy navigation and exploitation of the platform services to the AEGIS stakeholders.

The AEGIS platform offers two additional cross-platform functionalities that are incorporated within the components of the platform with a set of technologies and tools. The first functionality is the cross-platform security that follows a holistic security approach. Within this approach the security aspects of data in storage and data in transit, as well as the applied security on the platform operations and the technical interfaces, are covered. The second functionality is the multilingualism support that is supported in three different aspects of the platform, namely the static content, the metadata and the data. For the multilingualism support the eTranslation service of the European Commission (EC) is exploited.

For each component, a detailed description documenting the functionalities and the technical details is elaborated in Section 3 of the current deliverable.

## 2.2. Technical Architecture

In addition to the updated high-level architecture presented in section 2.1, Figure 2-2 illustrates the functional decomposition of the components of the AEGIS platform, as well as the relationship of the components and the corresponding data flow during run-time. The details for the design and specification of each component are described in Section 3.

### 2.3. AEGIS Integrated Notebooks

In the course of the development of the Query Builder, the Algorithm Execution Container and the Visualiser the technical partners decided to leverage the capabilities and features provided by the Jupyter notebook service that is already integrated within the AEGIS platform. Jupyter is providing functionalities for data ingestion, data discovery, data analytics and data visualisation to the data scientists, support for various languages such as Python and Scala, integration with data processing frameworks like Spark and support for user interface implementation in JavaScript.

Although the aforementioned components were initially developed as separate predefined notebooks containing several paragraphs, the consortium decided to integrate them also into one complete notebook towards the aim of offering a holistic toolset for data query and retrieval, data pre-processing, data analysis execution and advanced visualisations. Within this holistic notebook, all the described functionalities and features of the aforementioned three components are integrated in such a way that enables the end users of the platform to perform all the desired tasks from this complete notebook, providing intuitive and advanced user experience. The parallel existence of the separate predefined notebooks and the holistic notebook is facilitated by the advanced functionalities of the Jupyter environment that is utilised in all the aforementioned notebooks.

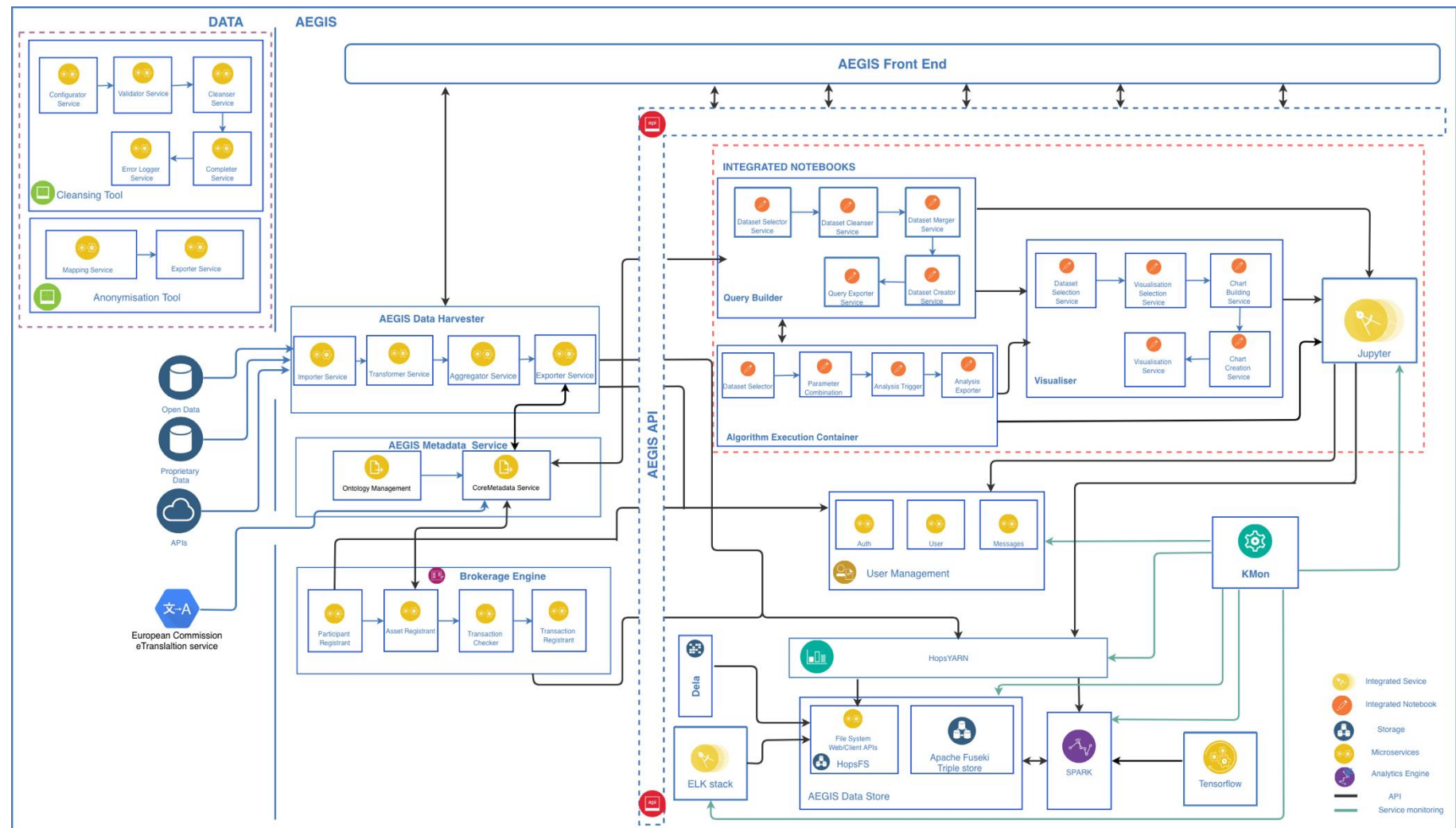


Figure 2-2: AEGIS Technical Architecture

### 3. AEGIS COMPONENTS AND APIs SPECIFICATIONS

#### 3.1. Data Harvester

The Data Harvester is an additional component, which has two main purposes. Firstly, it acts as an administrative tool for the operators of the AEGIS platform for providing a rich selection of already available datasets. Secondly, it can be used by the users of the AEGIS platform for easily retrieving data from well-known sources and interfaces. It manages both the retrieval of the actual datasets and the creation of the corresponding metadata.

##### 3.1.1. Overview

The Data Harvester is orchestrated out of several sub-components, including microservices and front-end modules. In connection they represent the process of harvesting, transforming, harmonising, annotating and providing the required data and metadata for the AEGIS platform. Therefore, they will be described as one component and from here on simply denominated as *Harvester*. The Harvester interacts tightly with the AEGIS Metadata Service and the AEGIS Data Store and is based on several basic concepts. In the following paragraphs these concepts are described in detail:

##### **Repository**

A repository represents a specific data source and handles the respective connection to it. Each repository represents descriptive and required data about the data source, where the address (in most cases a URL) is the most significant one.

##### **Annotation**

An annotation constitutes metadata of a project, dataset or file within the AEGIS platform. Hence it uses the AEGIS vocabularies and ontologies (see deliverable D2.1 Ch.3).

##### **Transformation**

A transformation describes all processing rules for converting the source data to the suitable target format. This may include mapping of fields, harmonisation and any converting.

##### **Harvester**

A harvester describes a concrete instance of retrieving data from a data source. It holds metadata about the harvesting process itself, like the execution schedule. One harvester is linked to a repository, the corresponding annotation and responsible transformation.

##### **Run**

A run depicts the single execution of a harvester, where the data and the metadata are generated and harvested respectively. It stores metadata about the performance and success of a harvesting process, which includes detailed logging information. A run can be scheduled and executed periodically.

These concepts are represented in the four microservices and the front-end of the Harvester. It is important to notice, that each microservice may have multiple instances or rather specialised implementations. E.g. they may be one importer for CSV data and one importer for JSON data.

##### **Importer**

An importer implements all functionality for retrieving data from a specific data source. It needs

to specifically support the characteristic of that data source, including protocol, serialisation format, security etc. It has to export the harvested data as JSON to the next stage.

### Transformer

A transformer converts the retrieved data from an importer into the target format of the AEGIS platform. Hence, a tabular format is specified for each data source.

### Aggregator

An aggregator collects converted data from a transformer over a configurable time interval. It allows to adjust the granularity of the available data in one file within the AEGIS platform.

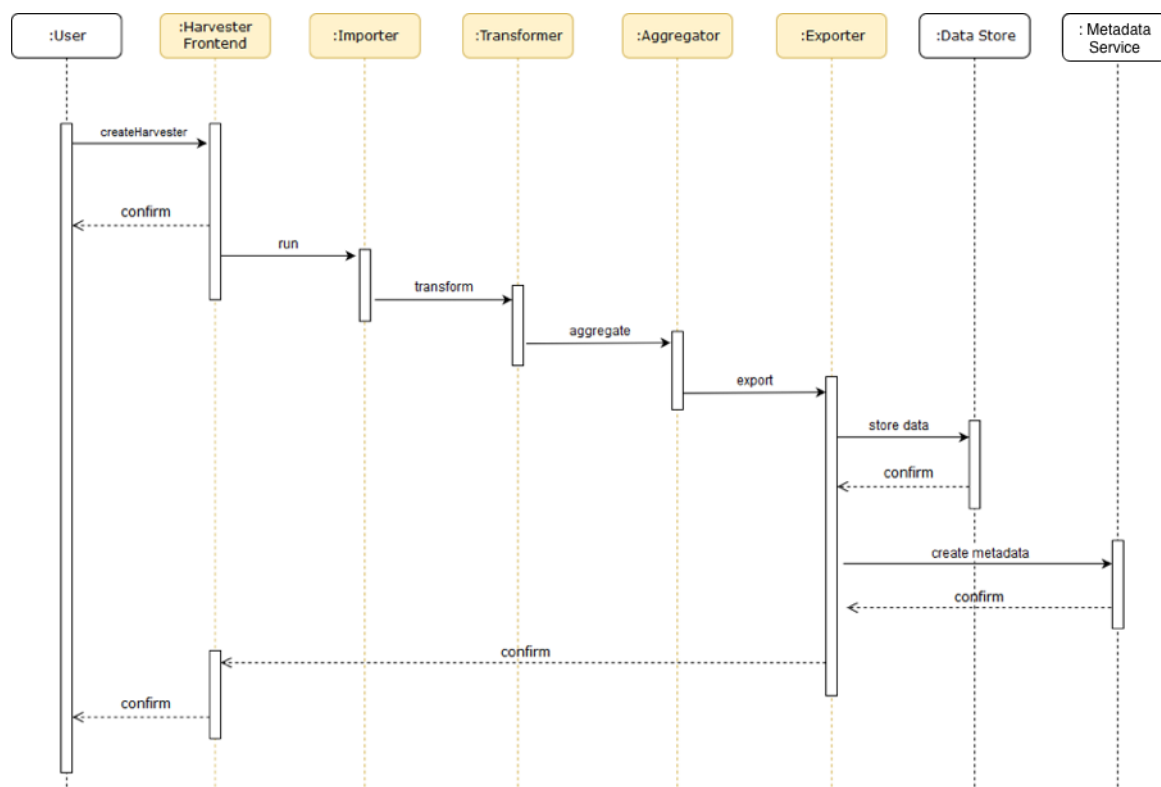
### Exporter

The exporter uploads transformed and/or aggregated data to the AEGIS platform. In addition, it creates the corresponding metadata in the AEGIS Metadata Service. There will be only one implementation for the exporter.

### Front-end

The front-end orchestrates the microservices and offers the visual interface for creating, editing and existing specific harvesting processes.

Figure 3-1 shows the process of harvesting data from a data source to the AEGIS platform.



**Figure 3-1: Sequence diagram of the Harvester component**

#### Updates from V3.0:

- No updates were introduced in terms of design and functionalities.

### 3.1.2. Supported Data Sources

The Data Harvester importer supports both concrete well-known data sources and generic standardised interfaces. The following data providers are facilitated:

Data Provider	Description	Interfaces
OpenWeatherMap ( <a href="https://openweathermap.org/">https://openweathermap.org/</a> )	Continuous harvesting of European Weather Data	RESTful JSON
European Data Portal - Single Dataset ( <a href="https://www.europeandataportal.eu/data">https://www.europeandataportal.eu/data</a> )	Harvest the resources and metadata of a single dataset	JSON Action API for metadata, Basic HTTP for files
European Data Portal – Multiple Datasets ( <a href="https://www.europeandataportal.eu/data">https://www.europeandataportal.eu/data</a> )	Harvest resources and metadata from multiple datasets by providing search parameters.	JSON Action API for metadata, Basic HTTP for files
AEGIS Event Detector	Push collected event data to the AEGIS platform	RESTful JSON (Push Method)
EM-DAT The International Disaster Database ( <a href="https://www.emdat.be/">https://www.emdat.be/</a> )	Collect disaster data	HTML Crawling Basic HTTP for files
Generic Dataset Creation	Creates CSV datasets and metadata in the AEGIS platform based on tabular data provided as JSON.	RESTful JSON
Generic Dataset Upload	Comprehensive interface for uploading file and metadata to the AEGIS platform. It includes support for bulk upload.	RESTful JSON

#### Updates from V3.0:

- No updates were introduced.

### 3.1.3. List of microservices

For the Harvester component four microservices are developed. Each service depicts one distinct task within the harvesting process. The orchestration of the services is done via a single-page-application front-end, which will be tightly integrated into the AEGIS platform. All services expose their functionality via a RESTful-API.

Component Name	Microservice Name	Functionalities
AEGIS Data Harvester	HarvesterImportService	<ul style="list-style-type: none"> <li>• Handling of repositories, e.g. creation and modification</li> <li>• Management of specific repository connectors</li> <li>• Execution of the importing process</li> <li>• Logging of importing process</li> <li>• Transfer of the imported data to the transformer service</li> </ul>
	HarvesterTransformerService	<ul style="list-style-type: none"> <li>• Management of transformations rules and scripts</li> <li>• Execution of the transformation from source data to the AEGIS data format</li> <li>• Logging of transformation process</li> <li>• Transfer of the transformed data to the aggregator or exporter service</li> </ul>
	HarvesterAggregatorService	<ul style="list-style-type: none"> <li>• Optional service for aggregating imported data for a specified time interval before exporting it to the AEGIS platform</li> <li>• Transfer of aggregated data to the exporter service</li> </ul>
	HarvesterExporterService	<ul style="list-style-type: none"> <li>• Handling of the export of the data to the AEGIS platform</li> <li>• Direct communication with the RESTful-API of AEGIS</li> <li>• Creation of the metadata in the Metadata-Service based on the given annotations</li> </ul>

**Table 3-1: Harvester list of microservices**

**Updates from V3.0:**

- No updates were introduced.

*3.1.4. Technologies to be used*

The foundation of the AEGIS Harvester is the EDP Metadata Transformer Service (EMTS)<sup>1</sup>, an open source solution for harvesting metadata from diverse Open Data sources.

For the purpose of the AEGIS platform, the EMTS is refined and updated to fit the needs of the project. This includes restructuring the application into small and scalable microservices. This is done by extracting the respective functionalities into new standalone services. These services are developed with the event-driven Java-framework Eclipse Vert.x<sup>2</sup>. It allows a much tighter integration into the AEGIS platform and the straightforward extension with additional functionalities. Correspondently, the web front-end is modified to single-page-application in order to act as an orchestrator of the various microservices. It is implemented based on the JavaScript Vue.js<sup>3</sup> framework that is allowing a better integration into the existing front-end of AEGIS platform.

Figure 3-2 shows the latest version of the front-end of the Harvester.

The screenshot shows the AEGIS Harvester web interface. At the top, there is a navigation bar with links for HOME, PIPES, NEW PIPE, and ACCOUNT. Below the navigation bar, there are two buttons: SAVE and RUN NOW. The main content area is divided into three sections: Modules, Data Flow, and Settings.

**Modules:** This section contains four buttons: OWM, CKAN, UPLOAD, and EVENT.

**Data Flow:** This section contains a form for configuring a data flow. The form has the following fields: pipeline ID (111), hopsProjectId (3434), hopsDataset (test), fetchType (ID), url (https://www.europeandataportal.eu/data/api/3/action), resourceId (liste-des-espaces-publics-numeriques-de-la-gironde2), durationInHours (0), and frequencyInMinutes (2).

**Settings:** This section contains a form for configuring the harvester's settings. The form has the following sections: Schedule (Frequency: Daily, Select date: 2018-10-17), Environment (Key-Value pairs), and Metadata (Name: Test, Tags).

**Figure 3-2: The harvester interface**

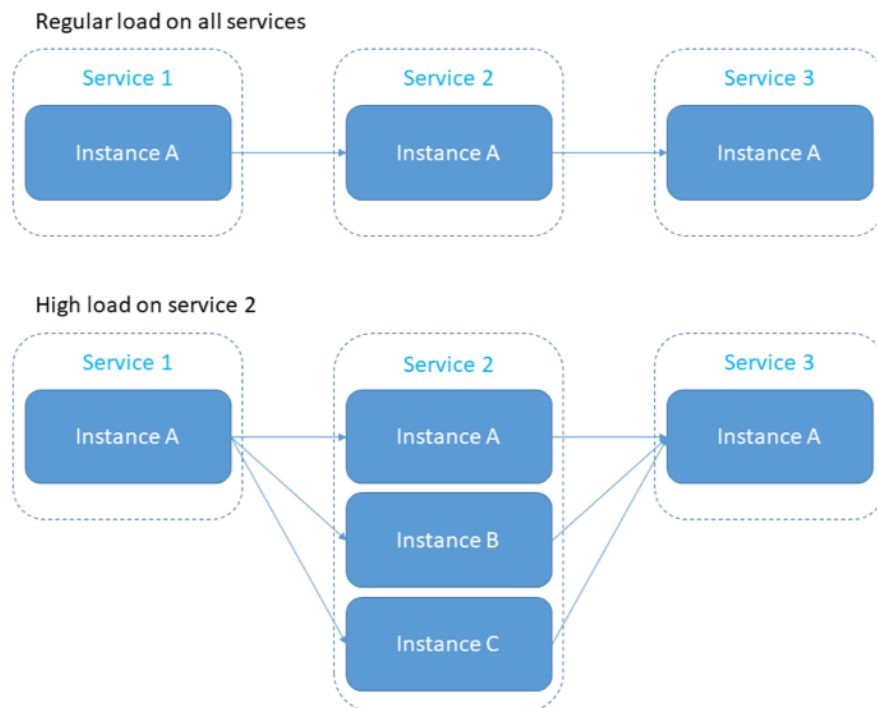
<sup>1</sup> The original source code can be found here: (<https://gitlab.com/european-data-portal/MetadataTransformerService>)

<sup>2</sup> <https://vertx.io/>

<sup>3</sup> <https://vuejs.org/>



The orchestration architecture used for the AEGIS Harvester follows a “pipeline” pattern, in which data is passed through several services, with each service manipulating the data in some sort. Each service is responsible for exactly one task. This permits a rather generic implementation of each service. The aim is to encourage a separation of concerns in order to enhance reusability, as well as allowing the dynamic scaling in times of high load. The latter is achieved by deploying additional instances of the services demanded most. Once new instances are spawned, request may dynamically be routed to the instance of a service with the least load. This is shown in Figure 3-3.



**Figure 3-3: Harvester Orchestration Concept**

The order and type of services participating in handling a certain use case is initially defined for later utilisation by the pipe implementation. The framework then builds the suitable requests (as well as the handling the concrete routing between instances) and provides the applicable configurations. This makes each service agnostic of its surroundings, aiding in the generic design mentioned earlier.

**Updates from V3.0:**

- No updates were introduced.

**3.1.5. APIs and exposed outcomes**

**Updates from V3.0:**

- No updates were introduced.

The following table documents the API of the Data Harvester component.

<b>Technical Interface</b>	
<b>Reference Code</b>	AH#01
<b>Function</b>	Orchestrate, schedule and manage harvesting processed
<b>Subsystems</b>	EDP Metadata Transformer Service
<b>Type, State</b>	
RESTful API, Web Front-end	
<b>Endpoint URI</b>	
<a href="http://aegis-harvester.fokus.fraunhofer.de">http://aegis-harvester.fokus.fraunhofer.de</a>	
<b>Input Data</b>	
Harvester endpoints	
<b>Output Data</b>	
Status logs	

**Table 3-2: Data Harvester technical interface**

## 3.2. Cleansing Tool

### 3.2.1. Overview

Data cleansing is an umbrella term for tasks that span from simple data pre-processing, like restructuring, predefined value substitutions and reformatting of fields (e.g. dates) to more advanced processes, such as outliers' detection and elimination from a dataset. Particularly in the AEGIS context of big data processing and analysis, cleansing may, by itself, be a process requiring big data technologies to be applied.

Within the context of the AEGIS platform, the consortium has decided for the data cleansing tasks to support a two-fold approach: (a) to offer an offline cleansing tool that is residing where the data are located and is enabling the execution of a variety of cleansing processes and functionalities utilised for the dataset preparation prior to importing them in the AEGIS platform (b) to provide an online cleansing tool that is facilitating the execution of simple but computationally intense data cleansing and manipulation tasks during the data query processing leveraging the computational power of the AEGIS platform.

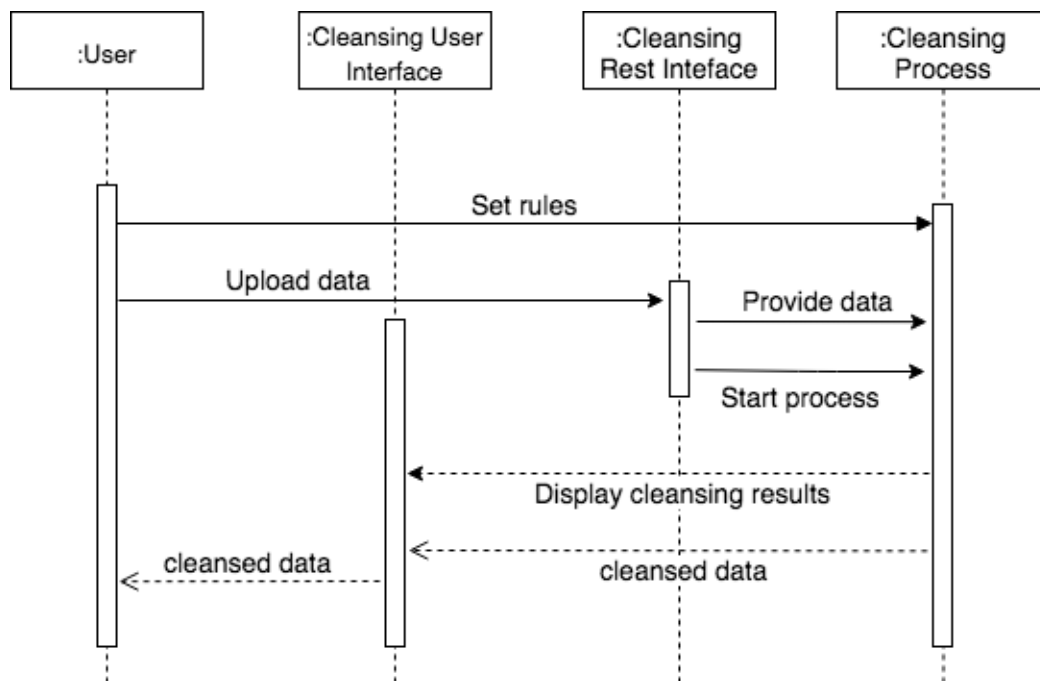
For simple data cleansing, it has become obvious that cleansing tasks that are both easy/straightforward and at the same time computationally intense, may emerge as steps of the analysis to be performed, i.e. they may be dependent on the specific application and not on intrinsic characteristics of the original dataset. These tasks will need to be performed as part of the online data manipulation. Hence, in order to provide a more intuitive user experience and also leverage the computational power of the AEGIS system, it was decided to make certain simple data cleansing functionalities available to the user during the data query creation process, i.e. when he/she should be more confident about the desired data manipulation needed to perform in order to use the data for further analysis. Hence, some simple custom cleansing processes are incorporated inside Query Builder as part of the data selection process and will be described in the corresponding section.

Additionally, the consortium identified the added value of providing an offline cleansing tool that will offer a level of customisation to the users and will be easily adaptable to the user's needs depending on the nature of the data source. Thus, it was decided to implement the offline cleansing tool that enables data validation, data cleansing and data completion processes towards the aim of increasing the reliability, accuracy and completeness of the data that are imported in the AEGIS platform. The tool is customisable, in terms of rules definitions for validation, cleansing and missing data handling, by the user and provides web-based user interface to display the cleansing process results. The tool supports simple cleansing tasks such as the predefined value substitutions and range constraints, but also complex cleansing task such as the outlier detection and removal.

The main functionalities of the offline cleansing tool are as follows:

- Define an extended list of rules for cleansing process (data validation, data cleansing, data completion).
- Provide a RESTful interface to facilitate the uploading of the dataset that will be used in the cleansing process and provide the cleaned data. Additionally, the corresponding documentation for the RESTful interface is provided to the users.
- Report the cleansing process results through an intuitive user interface.
- Offer an easy installation process in the premises of the user via Docker image.

The following figure shows the sequence diagram for the offline data cleansing. The sequence diagrams for data cleansing performed through other tools will be provided in the corresponding sections.



**Figure 3-4: Offline data cleansing sequence diagram**

#### Updates from V3.0:

- Support for easy installation via Docker image.
- The list of rules for cleansing process has been expanded with new rules.
- The REST-API documentation is provided using the Swagger framework.
- A series of optimisations were introduced for performance improvement especially for large datasets.

#### 3.2.2. List of microservices

For the offline data cleansing a list of microservices are developed and are orchestrated towards the execution of the cleansing tasks and the successful handling of the incoming requests for data cleaning transformations and corrective actions. In particular, the Cleansing Process, as shown in Figure 3-4, is composed of four microservices. The first microservice, the ConfiguratorService is undertaking the management of the constraints/rules for validation and data completion, as well as the corrective actions/rules. Additionally, three microservices, the ValidatorService, the CleanserService and the CompleterService, are responsible for the data validation, the data cleansing and the data completion respectively. The ErrorLoggerService is the microservice responsible for the collection and management of the log records that contain the identified errors and the corrective actions from the execution of the microservices of the Cleansing Process. Moreover, the ErrorLoggerService is providing the input for the Cleansing User Interface that reports the execution results to the user.

In total five microservices are developed and are described in the following table:

Component Name	Microservice Name	Functionalities
Cleansing Tool (Offline)	ConfiguratorService	<ul style="list-style-type: none"> <li>• Maintain and manage the constraints/rules for validation (e.g. specific data types, value representation, uniformity, range, regular expression patterns, cross-field validity)</li> <li>• Maintain and manage the corrective actions/rules (e.g. rejection of values, logical error identification)</li> <li>• Maintain and manage the data completion rules</li> </ul>
	ValidatorService	<ul style="list-style-type: none"> <li>• Perform data validation in accordance to the constraints/rules</li> <li>• Compile the list of identified errors identified in the validation</li> <li>• Log the errors in the appropriate log file</li> <li>• Provide interface for remote execution</li> </ul>
	CleanserService	<ul style="list-style-type: none"> <li>• Perform data cleaning based on the defined rules</li> <li>• Log the corrective actions in the appropriate log file</li> <li>• Provide interface for remote execution</li> </ul>
	CompleterService	<ul style="list-style-type: none"> <li>• Implement a list of methods / algorithms for data completion (e.x. Last Observation Carried Forward, Last Non-Zero, moving average, Linear regression, mean, median, k-NN).</li> <li>• Perform data completion based on the defined rules, algorithms and methods</li> <li>• Log the corrective actions in the appropriate log</li> <li>• Provide interface for remote execution</li> </ul>
	ErrorLoggerService	<ul style="list-style-type: none"> <li>• Create and display log records containing the errors and corrective actions</li> <li>• Manage log files generated by the rest of the microservices</li> </ul>

		<ul style="list-style-type: none"> <li>• Provide an interface for the rest of the microservices for log record creation</li> </ul>
--	--	--

**Table 3-3: Cleansing Tool list of microservices**

For the cleansing functionalities that are offered through the notebooks and notebook-based components, the relevant microservices are described in the corresponding sections.

**Updates from V3.0:**

- No updates were introduced.

*3.2.3. Technologies to be used*

For the online data cleansing, either through the dedicated data cleansing UI in the Query Builder or through custom processes -implemented with the help of the Jupyter notebook (which are part of the AEGIS Integrated Services) will be used. More details are provided in the corresponding tools' sections.

For the offline cleansing processes, which will be applied before importing data in AEGIS with the aim of making the data more easily processable by subsequent components in the data flow, the microservices architecture is followed and the corresponding microservices, as described in section 3.2.2, are written in Python, using Flask microframework<sup>4</sup> and a set of libraries such as Pandas<sup>5</sup> and NumPy<sup>6</sup>. For the documentation of the RESTful interface the Swagger framework<sup>7</sup> was utilised. In order to enable the easy installation of the offline cleansing tool, a Docker image is provided.

**Updates from V3.0:**

- The documentation of the RESTful interface is available through the Swagger framework.
- The tool is provided as a Docker image in order to facilitate the easy installation.

*3.2.4. APIs and exposed outcomes***Updates from V3.0:**

<sup>4</sup> <http://flask.pocoo.org/>

<sup>5</sup> <https://pandas.pydata.org/>

<sup>6</sup> <http://www.numpy.org/>

<sup>7</sup> <https://swagger.io/>

- No updates were introduced.

For the offline data cleansing a REST API interface is provided in order to enable the uploading of the dataset that will be cleansed and provide the cleaned dataset once the cleansing process is completed. Moreover, several additional interfaces are provided in order to facilitate the cleansing process execution. The details of these interfaces are documented in the tables below.

Since the online data cleansing is performed with the help of other components, more information is available in the corresponding sections.

Technical Interface	
Reference Code	CT#01
Function	Upload the dataset that will be used in the cleansing process
Subsystems	Offline Cleansing Tool
Type, State	
RESTful-API	
Endpoint URI	
<server url:5000>/cleaner/api/clean	
Input Data	
The data that will be used in the cleansing process in JSON format	
Output Data	
The cleansed data in JSON format	

**Table 3-4: Offline Cleansing tool technical interface**

Technical Interface	
Reference Code	CT#02
Function	Add new data provider
Subsystems	Offline Cleansing Tool

<b>Type, State</b>
RESTful-API
<b>Endpoint URI</b>
<server url:5000>/datasets/provider
<b>Input Data</b>
The data provider name
<b>Output Data</b>
-

**Table 3-5 Offline Cleansing tool add new provider**

<b>Technical Interface</b>	
<b>Reference Code</b>	CT#03
<b>Function</b>	Update existing data provider
<b>Subsystems</b>	Offline Cleansing Tool
<b>Type, State</b>	
RESTful-API	
<b>Endpoint URI</b>	
<server url:5000>/datasets/provider	
<b>Input Data</b>	
The old data provider name and the new data provider name	
<b>Output Data</b>	
-	

**Table 3-6: Offline Cleansing tool update provider**



Technical Interface	
Reference Code	CT#04
Function	Delete data provider
Subsystems	Offline Cleansing Tool
Type, State	
RESTful-API	
Endpoint URI	
<server url:5000>/datasets/provider	
Input Data	
The deleted data provider name	
Output Data	
-	

Table 3-7: Offline Cleansing tool delete provider

Technical Interface	
Reference Code	CT#05
Function	Add new dataset
Subsystems	Offline Cleansing Tool
Type, State	
RESTful-API	
Endpoint URI	
<server url:5000>/datasets/dataset	
Input Data	
The data provider name and the dataset name	

<b>Output Data</b>
-

**Table 3-8: Offline Cleansing tool new dataset**

<b>Technical Interface</b>	
<b>Reference Code</b>	CT#06
<b>Function</b>	Update an existing dataset
<b>Subsystems</b>	Offline Cleansing Tool
<b>Type, State</b>	
RESTful-API	
<b>Endpoint URI</b>	
<server url:5000>/datasets/dataset	
<b>Input Data</b>	
The data provider name, the old dataset name and the new dataset name	
<b>Output Data</b>	
-	

**Table 3-9: Offline Cleansing tool update an existing dataset**

<b>Technical Interface</b>	
<b>Reference Code</b>	CT#06
<b>Function</b>	Delete a dataset
<b>Subsystems</b>	Offline Cleansing Tool
<b>Type, State</b>	
RESTful-API	

<b>Endpoint URI</b>
<server url:5000>/datasets/dataset
<b>Input Data</b>
The data provider name and the dataset name
<b>Output Data</b>
-

**Table 3-10: Offline Cleansing tool delete dataset**

<b>Technical Interface</b>	
<b>Reference Code</b>	CT#07
<b>Function</b>	Add new variable
<b>Subsystems</b>	Offline Cleansing Tool
<b>Type, State</b>	
RESTful-API	
<b>Endpoint URI</b>	
<server url:5000>/datasets/variable	
<b>Input Data</b>	
The data provider name, the dataset name and the variable name	
<b>Output Data</b>	
-	

**Table 3-11: Offline Cleansing tool add new variable**

<b>Technical Interface</b>	
<b>Reference Code</b>	CT#08

<b>Function</b>	Update an existing variable
<b>Subsystems</b>	Offline Cleansing Tool
<b>Type, State</b>	
RESTful-API	
<b>Endpoint URI</b>	
<server url:5000>/datasets/variable	
<b>Input Data</b>	
The data provider name, the dataset name, the old variable name and the new variable name	
<b>Output Data</b>	
-	

**Table 3-12: Offline Cleansing tool update an existing variable**

<b>Technical Interface</b>	
<b>Reference Code</b>	CT#09
<b>Function</b>	Delete variable
<b>Subsystems</b>	Offline Cleansing Tool
<b>Type, State</b>	
RESTful-API	
<b>Endpoint URI</b>	
<server url:5000>/datasets/variable	
<b>Input Data</b>	
The data provider name, the dataset name and the variable name	
<b>Output Data</b>	

-

**Table 3-13: Offline Cleansing tool delete variable**

Technical Interface	
<b>Reference Code</b>	CT#10
<b>Function</b>	Add or update a validation rule
<b>Subsystems</b>	Offline Cleansing Tool
<b>Type, State</b>	
RESTful-API	
<b>Endpoint URI</b>	
<server url:5000>/rules/validation	
<b>Input Data</b>	
The data provider name, the dataset name, the variable name and the validation rule	
<b>Output Data</b>	
-	

**Table 3-14: Offline Cleansing tool add validation rule**

Technical Interface	
<b>Reference Code</b>	CT#11
<b>Function</b>	Delete a validation rule
<b>Subsystems</b>	Offline Cleansing Tool
<b>Type, State</b>	
RESTful-API	
<b>Endpoint URI</b>	

<server url:5000>/rules/validation	
<b>Input Data</b>	
The data provider name, the dataset name, the variable name and the validation rule	
<b>Output Data</b>	
-	

**Table 3-15: Offline Cleansing tool delete validation rule**

<b>Technical Interface</b>	
<b>Reference Code</b>	CT#12
<b>Function</b>	Add or update a cleaning rule
<b>Subsystems</b>	Offline Cleansing Tool
<b>Type, State</b>	
RESTful-API	
<b>Endpoint URI</b>	
<server url:5000>/rules/cleaning	
<b>Input Data</b>	
The data provider name, the dataset name, the variable name and the cleaning rule	
<b>Output Data</b>	
-	

**Table 3-16: Offline Cleansing tool add cleaning rule**

<b>Technical Interface</b>	
<b>Reference Code</b>	CT#13
<b>Function</b>	Delete a cleaning rule

<b>Subsystems</b>	Offline Cleansing Tool
<b>Type, State</b>	
RESTful-API	
<b>Endpoint URI</b>	
<server url:5000>/rules/cleaning	
<b>Input Data</b>	
The data provider name, the dataset name, the variable name and the cleaning rule	
<b>Output Data</b>	
-	

**Table 3-17: Offline Cleansing tool delete cleaning rule**

<b>Technical Interface</b>	
<b>Reference Code</b>	CT#14
<b>Function</b>	Update missing value rule
<b>Subsystems</b>	Offline Cleansing Tool
<b>Type, State</b>	
RESTful-API	
<b>Endpoint URI</b>	
<server url:5000>/rules/missing	
<b>Input Data</b>	
The data provider name, the dataset name, the variable name and the missing value rule	
<b>Output Data</b>	
-	

**Table 3-18: Offline Cleansing tool update missing value rule**

### 3.3. Anonymisation Tool

#### 3.3.1. Overview

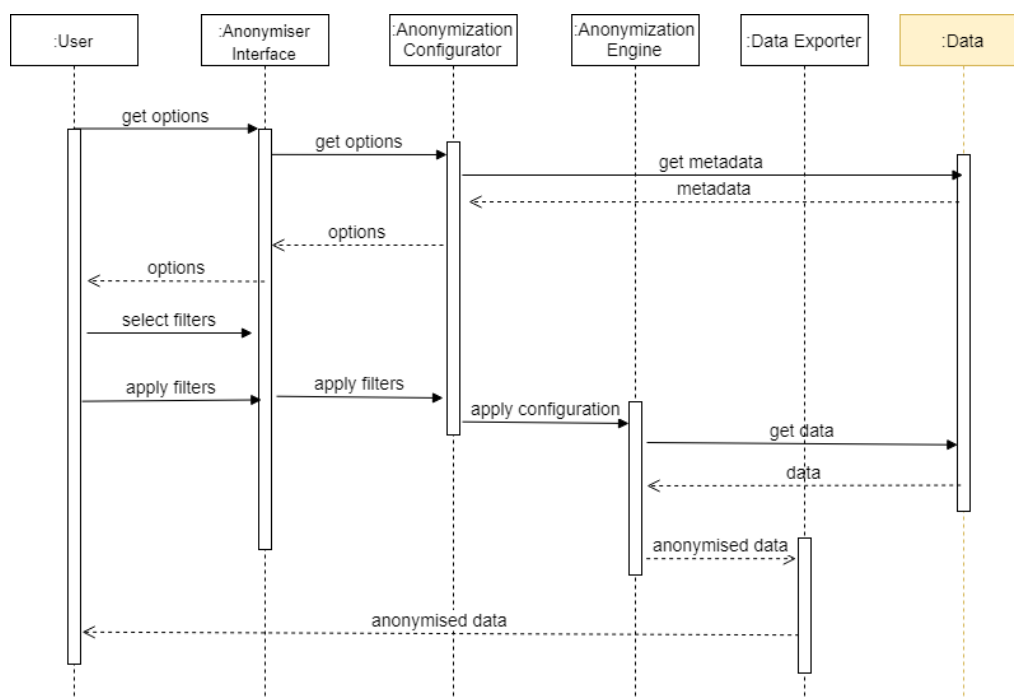
The anonymisation tool is an extensible, schema-agnostic plugin that allows real-time efficient data anonymisation. The anonymisation tool has been utilised for offline, private usage but offers the ability to output the anonymised data through a secured, web API. With emphasis on performance, the anonymisation tool syncs with private database servers and executes anonymisation functions on datasets of various sizes with little or no overhead. The purpose of the anonymisation is to enable the potential value of raw data in the system by accounting for privacy concerns and legal limitations.

The anonymisation process is optional to the AEGIS data flows and the tool is external to the core AEGIS platform, residing where the data to be anonymised are located. This decision ensures that no potentially sensitive data leave company premises, i.e. by-design eliminates any vulnerability risks entailed in uploading the initial eponymised, thus sensitive, data to the platform. Therefore, the AEGIS anonymisation solution will be used offline.

The main functionalities of the anonymisation tool are as follows:

- Connection to various data sources, including PostgreSQL, MySQL and csv files.
- Provision of anonymisation alternatives (generalisation, k-anonymity, pseudonymity), depending on the data schemas, the data values and the user's intended usage of the anonymised dataset.
- Export of anonymised data in files and as RESTful services, if desired.

Overall, the tool will help the user generate an anonymised dataset as an output, making sure that the individual sensitive records or subjects of the data cannot be re-identified.



**Figure 3-5: Data anonymisation sequence diagram**



**Updates from V3.0:**

- No updates were introduced in terms of design and functionalities.

*3.3.2. List of microservices*

The anonymisation tool, i.e. the AEGIS Anonymiser, comprises two microservices which are orchestrated towards the execution of the anonymisation workflow. The first microservice (Mapping Service) includes the functionalities provided by the Anonymisation Configurator and Anonymisation Engine shown in Figure 3-5. In the same figure, the Data Exporter corresponds to the second microservice, i.e. the Exporter Service. The Anonymiser Interface orchestrates the two microservices towards applying the anonymisation process and constitutes the interaction point with the user where required.

Component Name	Microservice Name	Functionalities
Anonymiser	Mapping Service	<ul style="list-style-type: none"> <li>• Connect to a database as data source</li> <li>• Read a csv file as data source</li> <li>• Provide anonymisation alternatives (e.g. generalisation, pseudonymity) per field</li> <li>• Apply the selected anonymisation action</li> </ul>
	Exporter Service	<ul style="list-style-type: none"> <li>• Provide the anonymised dataset through a REST API</li> <li>• Save the anonymised dataset as csv file</li> </ul>

**Table 3-19: Anonymisation Tool list of microservices****Updates from V3.0:**

- No updates were introduced.

*3.3.3. Technologies to be used*

The anonymisation tool is based on the Anonymiser, an anonymisation and persona-building tool, developed in the context of the European project CloudTeams<sup>8</sup>.

The tool performs a type of generalisation, which can be used to achieve k-anonymity. It allows users to customise the level of anonymisation per data field, i.e. sensitive data fields can be

<sup>8</sup> <https://github.com/cloudteams>

completely stripped out or suppressed from the output with asterisks or can be generalised. With the generalisation mapping, individual values of input data fields are replaced by a broader category. For example, the value '15' of the attribute 'Age' may be replaced by ' $\leq 18$ ', the value '23' by ' $20 < \text{Age} \leq 30$ '. The user may then apply a threshold (k) on the minimum number of entries with the same value, thus ensuring k-anonymity. A pseudonymity functionality is also available to hide personal information and all data fields can be masked with ranged data.

The original tool is written in Python, using the Django web framework. These technologies are also used to deliver the necessary updates and extensions in order to support the AEGIS anonymisation requirements. Specifically, the tool is extended to support csv files as data sources.

As the anonymisation tool is not integrated with the other AEGIS components but only offers limited interaction points, there is a flexibility in diversifying the provided anonymisation solution. Hence, in the course of the project, the tool is extended and adapted to the project's requirements.

#### Updates from V3.0:

- No updates were introduced.

#### 3.3.4. APIs and exposed outcomes

#### Updates from V3.0:

- No updates were introduced.

The outcome of the Anonymisation tool is available through a REST API, documented in the following table.

Technical Interface	
Reference Code	AZ#01
Function	Retrieve the anonymised data
Subsystems	None / Standalone API
Type, State	
RESTful-API	
Endpoint URI	

<server url>: anonymizer/api/<secret key>/<parameters>
<b>Input Data</b>
<b>Secret key:</b> The secret access key generated for the user through the Anonymiser’s user interface  <b>Parameters:</b> Parameters for the data to be returned, including <b>limit</b> , <b>offset</b> , <b>filters on properties</b> and <b>count</b>
<b>Output Data</b>
The anonymised data in JSON format

**Table 3-20: Anonymisation tool technical interface**

### 3.4. Brokerage Engine

#### 3.4.1. Overview

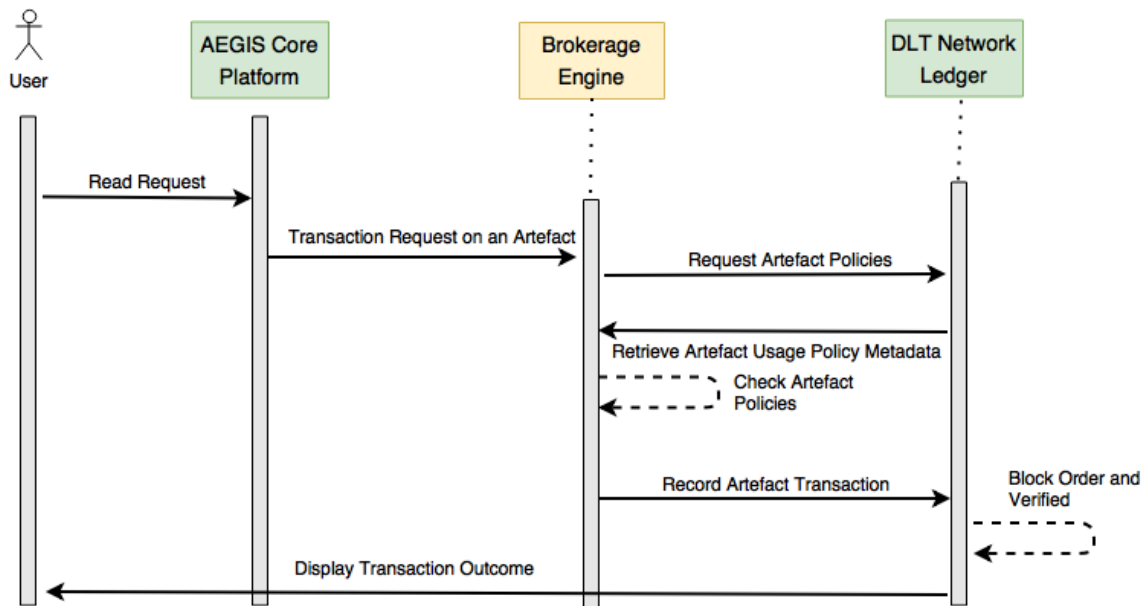
The AEGIS Brokerage Engine is responsible for instantiating parts of the Data Policy and the Business Brokerage frameworks, that have been finalised in D2.3 and constitute the models of the above-mentioned frameworks that provide the different attributes that are either necessary or desirable to enable the option of running transactions on top of either an AEGIS cluster, or a network of AEGIS clusters.

The final brokerage engine will only be used for transactions of assets (datasets) that are not free on the AEGIS platform, but are provided for a fee, depending on each owners’ wills. As such, the overall blockchain will only include transactions that refer to “for a fee” datasets, leaving out operations that are performed over freely and publicly available datasets.

In this context, the engine listens to activities that are to be performed on the AEGIS Cluster and the Data Store to prepare the Distributed Ledger Network records. These refer to adding new users on the cluster, which are also added to the ledger using the “Participant Registrant” microservice and to registering Assets on the Data Store, which are added to the ledger using the “Asset Registrant” microservice.

Upon a transaction request, and in case this does not refer to a public/free dataset, the “Transaction Checker” microservice checks each artefact’s metadata to conclude if a certain operation is possible. The first checks take place on the AEGIS Data Store, checking metadata stored there, increasing in this way the response rate as no extra API calls to the core Brokerage Engine are needed. In case a transaction is deemed as permissible, then the Brokerage Engine is engaged, where it checks against its ledger to see whether a condition applies that does not permit the operation for the data artefact under observation (for example not enough AEGIS “coins”, or having a dataset blocked by allowing its explicit use by a user for a certain period that is not over, etc.).

Once, a successfully concluded transaction is marked on the ledger using the “Transaction Registrant” microservice.



**Figure 3-6: Brokerage Engine sequence diagram**

#### Updates from V3.0:

- Revised and simpler Blockchain Model to accommodate the shared data brokerage framework that was devised in D2.3
- Improved centralised design for the core blockchain ledger, based on a cluster that connects to the AEGIS Core Platform via API calls

#### 3.4.2. List of microservices

The microservices of the Brokerage Engine are tasked with the storage, checking and updating of data in the AEGIS Distributed Ledger network.

Component Name	Microservice Name	Functionalities
Brokerage Engine	Participant Registrant	<ul style="list-style-type: none"> <li>• Listens to the registration facility of AEGIS</li> <li>• Registers the AEGIS users as participants of the Brokerage Engine of the AEGISDL network</li> </ul>
	Asset Registrant	<ul style="list-style-type: none"> <li>• Registers assets in the AEGIS Distributed Ledger network once a user selects an asset to be available over the ledger via:               <ul style="list-style-type: none"> <li>○ Communication with the AEGIS metadata storing</li> </ul> </li> </ul>

		methods listening to datasets storing ○ Communication with the Harvester and listening to datasets storing
	Transaction Checker	• Checks transaction details against details stored on the AEGIS Distributed Ledger
	Transaction Registrant	• Registers transactions in the Distributed Ledger • Updates the “wallets” of the transaction participants • Exposes executed transactions through a REST API

**Table 3-21: Brokerage engine list of microservices****Updates from V3.0:**

- The Asset registrant service refers only to datasets in this version, based on the implementation decision to utilise the ledger only for datasets transactions.

*3.4.3. Technologies to be used*

The AEGIS Brokerage engine is built on top of the of Hyperledger Fabric<sup>9</sup> framework and provides an API that is consumed by the AEGIS platform for providing the interconnection between the core platform and the Brokerage Engine. The models of the Blockchain engine have been constructed based on the AEGIS DPF presented in deliverable D2.1, while Hyperledger Composer<sup>10</sup> is being used for testing and further optimising the overall engine, and for providing an interface to easily manage the overall network that has been deployed.

**Updates from V3.0:**

- No updates were introduced.

*3.4.4. APIs and exposed outcomes***Updates from V3.0:**

<sup>9</sup> <https://www.hyperledger.org/projects/fabric>

<sup>10</sup> <https://hyperledger.github.io/composer/>

- Removal of API methods that change costs and status of assets.

The following tables present the most crucial interfaces used by the Brokerage Engine which are necessary for the interconnection with the AEGIS core platform.

Technical Interface	
Reference Code	BE#01
Function	User
Subsystems	Brokerage Engine
Type, State	
RESTful-API	
Indicative Endpoints	
GET /api/User	Get a list of all users registered with the brokerage engine
POST /api/User	Add a user to the brokerage engine
GET /api/User/{id}	Get user's details
Input Data (for POST)	
<pre>{   "\$class": "eu.aegis.User",   "uid": "string",   "balance": "0.0",   "externalAssets": [     {}   ] }</pre>	

**Table 3-22: Brokerage Engine technical interface 1**

Technical Interface	
Reference Code	BE#02
Function	AEGISAsset

<b>Subsystems</b>	Brokerage Engine
<b>Type, State</b>	
RESTful-API	
<b>Indicative Endpoints</b>	
<b>GET</b> /api/AEGISAsset	Get a list of all AEGIS assets
<b>POST</b> /api/AEGISAsset	Add an asset to the brokerage engine
<b>GET</b> /api/AEGISAsset/{id}	Get asset's details
<b>Input Data (for POST)</b>	
<pre>{   "\$class": "eu.aegis.AEGISAsset",   "aid": "string",   "type": "Dataset",   "cost": "0.0",   "status": "Free",   "exclusivity": "None",   "contractText": "string",   "owner": {} }</pre>	

**Table 3-23: Brokerage Engine technical interface 2**

<b>Technical Interface</b>	
<b>Reference Code</b>	BE#03
<b>Function</b>	BuyAsset
<b>Subsystems</b>	Brokerage Engine
<b>Type, State</b>	
RESTful-API	
<b>Indicative Endpoint</b>	

<b>POST</b> /api/BuyAsset	Buy an asset
<b>Input Data</b>	
<pre>{   "\$class": "eu.aegis.BuyAsset",   "buyer": {},   "relatedAsset": {},   "transactionId": "string",   "timestamp": "2018-03-21T10:10:29.343Z" }</pre>	

**Table 3-24: Brokerage Engine technical interface 3**

<b>Technical Interface</b>	
<b>Reference Code</b>	BE#04
<b>Function</b>	LoadBalance
<b>Subsystems</b>	Brokerage Engine
<b>Type, State</b>	
RESTful-API	
<b>Indicative Endpoints</b>	
<b>POST</b> /api/LoadBalance	Loads currency to the balance of a user
<b>Input Data</b>	
<pre>{   "\$class": "eu.aegis.LoadBalance",   "amount": 0,   "user": {},   "transactionId": "string",   "timestamp": "2018-03-21T10:10:29.384Z" }</pre>	



**Table 3-25: Brokerage Engine technical interface 4**

### 3.5. AEGIS Data Store

#### 3.5.1. Overview

The AEGIS Data Store has two distinct components; the HopsFS, which is the distributed file system mainly used for storing large amounts of data, as well as, the AEGIS Metadata Service, which is responsible for storing the metadata about the datasets.

#### 3.5.2. HopsFS filesystem

The AEGIS Data Store component is responsible for storing data that were collected and curated by the Harvester. A distributed file system approach was chosen for flexibility, reliability, and scalability. The distributed file system will allow storing large amounts of data while enabling access to the data from other AEGIS supported services such as the Query Builder and the Visualiser. In particular, the distributed file system is primarily responsible for storing large files, that is, files ranging from megabytes to terabytes in size. However, as seen in many production Big Data clusters such as the ones at Yahoo and Spotify [1], it has been observed that almost 20% of the files in the cluster are less than 4 KB in size and as much as 42% of all the file system operations are performed on files less than 16 KB in size.

Under the hood, AEGIS uses HopsFS as the main file system to store the data. HopsFS is a reliable, highly scalable, and fault tolerant distributed file system. A file is stored as list of blocks that is triple replicated for fault tolerance. Unlike HDFS that stores the file system metadata in memory, HopsFS keeps all the file system metadata in an in-memory distributed database providing bigger clusters with higher throughput.

In addition to the traditional POSIX permissions model, HopsFS supports extended Access Control Lists (ACLs). ACLs are useful for implementing permission requirements beyond the usage of only users and groups.

Moreover, HopsFS supports heterogeneous storage model that allows Datanodes to annotate each of their datanode directory with a storage type. The storage type identifies the underlying storage media (HDD, SSD, etc.). A storage policy can then be specified per directory in HopsFS to dictate which storage types to be used when adding new files and/or directories. For example, “All\_SSD” policy enforces all replicas to be stored on SSD.

**Updates from V3.0:**

- No updates were introduced in terms of design and functionalities.
- Several bug fixes and performance improvements were introduced.

### 3.5.2.1. List of microservices

HopsFS runs as a service in the AEGIS cluster where users can interact with it using the AEGIS user interface and the REST API provided by Hopsworks. Under the hood, the AEGIS user interface communicates with HopsFS using the client APIs.

Component Name	Microservice Name	Functionalities
HopsFS	File System (Client/Web APIs)	<ul style="list-style-type: none"> <li>Perform file system operations such as create, mkdir, delete, append, etc</li> </ul>

**Table 3-26: HopsFS list of microservices**

#### Updates from V3.0:

- No updates were introduced.

### 3.5.2.2. Technologies to be used

The AEGIS platform uses a file system, HopsFS, as the main store for Big Data. HopsFS is a drop-in replacement for Hadoop Distributed File System (HDFS). HopsFS is designed primarily to store large files, however, as reported most of production clusters contains a large number of small files (< 64KB). Therefore, we have extended HopsFS to efficiently manage large number of small files using a tiered storage solution. The tiers range from the highest tier where an in-memory database stores very small files (<1 KB), to the next tier where small files (<64 KB) are stored in Solid State Drives (SSDs), also using the database, to the largest tier, the existing Hadoop block storage layer for large files. Our approach is based on extending HopsFS with an inode stuffing technique, where we embed the contents of small files with the metadata and use database transactions and database replication guarantees to ensure the availability, integrity, and consistency of small files.

#### Updates from V3.0:

- No updates were introduced.

### 3.5.2.3. APIs and exposed outcomes

#### Updates from V3.0:

- No updates were introduced.

The small files are handled transparently by the client and the file system without involving the users. It is recommended to interact with the data in HopsFS from the AEGIS user interface.

However, HopsFS can be accessed using the command line, Java client APIs, and RESTful APIs.

<b>Technical Interface</b>	
<b>Reference Code</b>	EDS#01
<b>Function</b>	HopsFS FileSystem
<b>Subsystems</b>	HopsFS
<b>Type, State</b>	
RPC, Synchronous	
<b>API Documentation</b>	
<a href="https://hadoop.apache.org/docs/stable/api/org/apache/hadoop/fs/FileSystem.html">https://hadoop.apache.org/docs/stable/api/org/apache/hadoop/fs/FileSystem.html</a>  Unsupported Calls: <ul style="list-style-type: none"> <li>• (get set list remove)XAttr : At the moment adding extended metadata is done from Hopsworks</li> <li>• (create rename delete) Snapshot</li> </ul>	
<b>Input Data</b>	
Multiple formats (depending on the chosen interface)	
<b>Output Data</b>	
Multiple formats (depending on the chosen interface)	

**Table 3-27: AEGIS Data Store technical interface 1**

<b>Technical Interface</b>	
<b>Reference Code</b>	EDS#02
<b>Function</b>	HopsFS WebHDFS
<b>Subsystems</b>	HopsFS
<b>Type, State</b>	
RESTful-API, Synchronous	
<b>API Documentation</b>	

<a href="https://hadoop.apache.org/docs/stable/hadoop-project-dist/hadoop-hdfs/WebHDFS.html">https://hadoop.apache.org/docs/stable/hadoop-project-dist/hadoop-hdfs/WebHDFS.html</a>
Unsupported Call: <ul style="list-style-type: none"> <li>• (get set remove) XAttr</li> <li>• (create rename delete) Snapshot</li> </ul>
<b>Input Data</b>
Multiple formats (depending on the chosen interface)
<b>Output Data</b>
Multiple formats (depending on the chosen interface)

**Table 3-28: AEGIS Data Store technical interface 2**

### 3.5.3. AEGIS Metadata Service

#### 3.5.3.1. Overview

The AEGIS Metadata Service is responsible for storing the metadata associated with a particular dataset within the AEGIS platform. These metadata pose the foundation of the processing of the data within the AEGIS platform, since they offer detailed information about the semantic and syntax of the data. This allows choosing appropriate analysis and visualisation methods. The metadata will be stored as Linked Data, using the AEGIS ontology and vocabulary<sup>11</sup>, which is based upon the DCAT-AP specifications. It will be developed as a service and integrated into the AEGIS Data Store.

#### Updates from V3.0:

- There are no updates in the general objective.

#### 3.5.3.2. Concept

The metadata management of the AEGIS platform is based on Linked Data specifications and technologies. It heavily reuses existing Linked Data vocabularies and relies on the Resource Description Framework (RDF). The native store for RDF is a triplestore. In general, triplestores only offer basic management functionality, mostly restricted to SPARQL features. In order to integrate the metadata into the AEGIS platform additional functionalities are required. In the following paragraphs these features are described in detail:

<sup>11</sup> <https://github.com/aegisbigdata/aegis-ontology>

## URI Management

Triplestores do not apply constraints to the choice and creation of URIs. SPARQL allows to freely choose URIs for graphs, objects, predicates or subjects and stores them accordingly. The AEGIS Metadata Service will automatically provide consistent and comprehensive URIs for all entities upon creation time. Those will follow best-practices, especially promoted by the Semantic Interoperability Community (SEMIC) of the European Commission [2]. For instance, the URI of a dataset would be: <http://www.aegis-bigdata.eu/set/data/<id>>

## Access Control

Datasets may be not public and therefore the metadata needs to be protected too. RDF does not natively support access control, but multiple approaches are discussed in literature [3]. The cleanest and most common solution is query rewriting, where the metadata is “annotated” with additional properties, which indicate the ownership. This will be done by adding an appropriate property to the AEGIS ontology. A proxy service is then sited in front of the SPARQL endpoint to enrich the queries with the current user by applying the FILTER directive. The corresponding result will also be filtered to remove any sensible ownership information. Hence, the native SPARQL endpoint will not be accessible anymore, but the proxy service. This service is part of the CoreMetadataService.

## Multilingual Capacities

RDF natively supports the provision of multilingual literals. This feature will be harnessed appropriately by adding support for it to the RESTful API of the service.

## Data Store Synchronisation

By definition the metadata refers to actual data in the AEGIS Data Store. Since data and metadata are managed by distinct service, a reliable synchronisation and connection is required. The logical link between the services is managed by suitable properties in the metadata, linking the entities in one service to the other. E.g. the unique ID of a dataset in the AEGIS Data Store is stored in the corresponding metadata (datasetId). The physical link is established by adding hooks into the CRUD functionality of the AEGIS Data Store, which calls the Metadata Service respectively. The synchronisation will not be bi-directional, but led by the AEGIS Data Store. Hence, a deletion of metadata will not delete the linked dataset, but deleting a dataset will delete the metadata.

## Recommendation Service

The core strength of Linked Data is the ability to discover (previously unknown) relations between entities. This is harnessed in a recommendation service. Based on suitable and advanced SPARQL queries similar and/or relevant datasets can be recommended, e.g. data from the same geographical area or with a similar semantics in the table.

## Thesauri and Vocabulary Management

The metadata should follow highest quality standards, avoid redundancies and reuse existing vocabularies. Therefore, an additional service for managing ontologies and vocabularies will

be applied. It acts as a helper service for simplifying the metadata provision process through auto-completion and recommendation features in order to:

- Provide an autocomplete functionality for metadata properties, classes and named individuals, defined by the AEGIS core vocabulary and DCAT-AP.
- Assist in the provision of metadata properties, which may be independent of the AEGIS core vocabulary. For example, propose appropriate properties or classes for a given keyword.

In order to fulfil the additional functionalities, the AEGIS Metadata Service consists out of three basic components:

- A standard triplestore as basic store.
- A core metadata service for managing URIs, access control, multiple languages and synchronisation.
- The LinDA Vocabulary and Metadata Repository for managing existing vocabulary and thesauri.

### 3.5.3.3. List of microservices

The AEGIS Metadata Service is a combination of two microservices.

Component Name	Microservice Name	Functionalities
AEGIS Metadata Service	CoreMetadataService	<ul style="list-style-type: none"> <li>• Creating and modifying the Linked Data metadata</li> <li>• Transform simple JSON to Linked Data</li> <li>• Recommendation engine for getting similar or suitable additional data</li> <li>• Proxy service for access control</li> </ul>
	OntologyManagementService	<ul style="list-style-type: none"> <li>• Management of the AEGIS Linked Data vocabularies and ontologies</li> <li>• Exposes reusable namespaces for generating the metadata</li> <li>• Based on the LinDA Vocabulary and Metadata Repository</li> </ul>

**Table 3-29: AEGIS Metadata service list of microservices**

**Updates from V3.0:**

- The exact features of each microservice are described in detail. The list of microservices did not change.

#### 3.5.3.4. Technologies to be used

The AEGIS Metadata service is implemented with the following technologies:

- Triplestore (Apache Fuseki and Virtuoso are supported)
- Vert.x framework for the CoreMetadataService
- LinDA Vocabulary and Metadata Repository

#### Updates from V3.0:

- The technology stack was changed and migrated from the previous version.

#### 3.5.3.5. APIs and exposed outcomes

#### Updates from V3.0:

- LinDA was added as an official artefact.

The AEGIS Metadata service will be exposed as two artefacts and services.

Technical Interface	
Reference Code	AL#01
Function	Managing the AEGIS metadata
Subsystems	Triplestore
Type, State	
RESTful-API, SPARQL endpoint	
Endpoint URI	
<a href="http://aegis-metadata.fokus.fraunhofer.de/">http://aegis-metadata.fokus.fraunhofer.de/</a>	
Input Data	
Metadata as JSON or RDF	

<b>Output Data</b>
Metadata as JSON or RDF

**Table 3-30: AEGIS Metadata Service technical interface**

<b>Technical Interface</b>	
<b>Reference Code</b>	AL#02
<b>Function</b>	Thesauri and Vocabulary Management
<b>Subsystems</b>	Elasticsearch
<b>Type, State</b>	
RESTful-API	
<b>Endpoint URI</b>	
<a href="http://linda.epu.ntua.gr/coreapi/recommend/">http://linda.epu.ntua.gr/coreapi/recommend/</a>	
<b>Input Data</b>	
GET Parameter	
<b>Output Data</b>	
JSON	

**Table 3-31: AEGIS Metadata Service technical interface 2**

### 3.6. AEGIS Integrated Services

#### 3.6.1. Overview

The AEGIS platform provides a multi-tenant data management and processing services for Big Data. The multi-tenancy behaviour allows different users and services to securely and privately access and process their data. The AEGIS platform enables users to share their data with other users on the platform and allow access for specific services. In addition, users can use different data processing services that are supported by the platform to process and visualise their data.

Under the hood, the data are mainly stored in the AEGIS Data Store; however, the AEGIS Data Store APIs are kept hidden from users. Instead, the AEGIS platform provides a Project/Dataset service to allow users to upload/download, explore, and do analysis on their data in a secure way without interacting with the AEGIS Data Store directly. To ensure secure and private access to the data, each user has an x509 certificate per project as well as a specific project user



for the Data Store per project. The certificate has a CN field which contains the project specific username and that gives the platform the possibility to provide application level authorisation at the RPC server. For instance, any application executed within a YARN container will access the Data Store (HopsFS) as the user running this application. YARN acts as a proxy user for the user and accesses HopsFS (HDFS) through user impersonation. Thus, all accesses are seen as being done by the running user and storage access is limited to the files that can be accessed by this user. Moreover, each user has a specific project user for each of the projects that he/she can access. This means that any YARN application can only access files that are accessible for the running project and cannot normally access files cross projects, even if the project belongs to the same user as the one running the application. All applications running on top of YARN such as Spark, will be governed by the same storage access as described for a YARN container.

Under the hood, AEGIS builds upon Hopsworks to provide integrated support for different services such as interactive notebooks with Jupyter that are used mainly by the Query Builder, Algorithm Execution Container and the Visualiser components. Other services such as Kafka, and ELK stack are also supported.

In addition to this, Tensorflow is also supported. TensorFlow<sup>12</sup> is an open source software library released in 2015 by Google to make it easier to design, build, and train deep learning models. Tensorflow is providing an ecosystem suitable for high performance numerical computation following a flexible architecture that enables easy deployment of computation across a variety of platforms (CPUs, GPUs and TPUs). Although TensorFlow is only one of several options available, we choose to use it here because of its good design, ease of use and large community of adopters.

#### Updates from V3.0:

- No updates were introduced in terms of design and functionalities.
- Removed Zeppelin support due to the issues documented in section 3.6.3.
- Several bug fixes and performance improvements were introduced.

#### 3.6.2. List of microservices

Hopsworks provides different integrated services that interact with each other and with users using the Hopsworks REST API.

Component Name	Microservice Name	Functionalities
Users	Auth	<ul style="list-style-type: none"> <li>• Provides authentication functionality for users to login, logout, register, and recover password</li> </ul>
	User	<ul style="list-style-type: none"> <li>• Provides information about the current user</li> </ul>

<sup>12</sup> <https://www.tensorflow.org/>

	Messages	<ul style="list-style-type: none"> <li>Provides an inbox functionality for users where they receive/send share requests for Datasets</li> </ul>
Projects	Projects	<ul style="list-style-type: none"> <li>Provides information about the projects for a user, as well as details on each of the projects such as list of datasets, description, and team members</li> </ul>
	Datasets	<ul style="list-style-type: none"> <li>Provides information about the datasets for a user.</li> <li>It provides upload, download, and explore functionalities on the data</li> </ul>

**Table 3-32: AEGIS Integrated services list of microservices****Updates from V3.0:**

- No updates were introduced.

*3.6.3. Technologies to be used*

AEGIS builds upon Hopsworks to provide multi-tenant data management and processing services for BigData. Hopsworks is a project-based multi-tenant platform for secure collaborative data science that runs on top of HopsFS. It provides an integrated support for different data parallel processing services such as Spark, Flink, and MapReduce, as well as a scalable messaging bus with Kafka, and interactive notebooks with Jupyter. Hopsworks introduces new abstractions called Projects and Datasets that provide the basis for which users can securely upload and privately process data and securely collaborate with other users on the platform. A Dataset is a directory subtree in HopsFS that can be shared between projects. A Project is a collection of datasets, users, and notebooks (Jupyter). In the AEGIS platform, Jupyter is mainly used by the Query Builder, the Algorithm Execution Container and the Visualiser components of the platform. In the course of the project, several updates were introduced, that were mainly bug fixes for usability and performance of the platform and the certificates handling.

TensorFlow is an ecosystem for developing deep learning models, containing all the tools from building to deployment. TensorFlow has 3 main components:

1. TensorFlow (API) - contains the API's to define the models and train the models with the data.
2. TensorBoard - helps to analyse, visualize, and debug TensorFlow graphs.
3. TensorFlow Serving - helps to deploy the pre-trained models.

Initially Apache Zeppelin was also considered for running computations in the form of notebooks, however multiple problems have been observed with the Zeppelin tool since the beginning of the project:

- A slow down on development, bug fixes and keeping up to date with third party dependencies.
- Jupyter has been found as a better fit with the Hopsworks security model rather than Zeppelin.
- The Jupyter community is more vibrant and active than the Zeppelin community.

One critical issue that was encountered in the development process of Hopsworks, was that Zeppelin uses a number of third-party dependencies, such as the Elastic Stack, but Zeppelin upgrade process to newer versions is performed with a very slow time plan. Since Zeppelin is not the only tool depending on Elastic Search, the slow upgrade time plan to newer versions would also force Hopsworks to slow down its updates on other services until Zeppelin would upgrade their own dependencies. For all these reasons, the consortium decided to remove the support of Zeppelin and proceed only with Jupyter.

**Updates from V3.0:**

- Removed the support for Zeppelin.

#### 3.6.4. APIs and exposed outcomes

**Updates from V3.0:**

- The AEGIS platform APIs have not changed since most of the updates are bug fixes for platform usability and performance that will not require change of user facing APIs.

##### 3.6.4.1. Users API

Hopsworks provides a RESTful API to create users and to login to the platform as documented in <https://app.swaggerhub.com/apis/maimail/hopsworks-user-api/1.0.0>

##### 3.6.4.2. Projects and Datasets

Once logged in, users can create Project/Dataset, add member to a Project, share their Dataset, or upload/download/analyse their data. The RESTful API is documented in <https://app.swaggerhub.com/apis/maimail/hopsworks-core-api/1.0.0>

##### 3.6.4.3. Interactive Notebooks

Users can create an interactive notebook in their Project using Jupyter. Jupyter is a web-based notebook service that allows users to interactively analyse and visualise their data using

different frameworks such as Spark. It only provides some basic charts; however, different JavaScript libraries could be loaded to support a more complex visualisation or the AEGIS Visualiser component could be utilised.

### 3.7. Query Builder

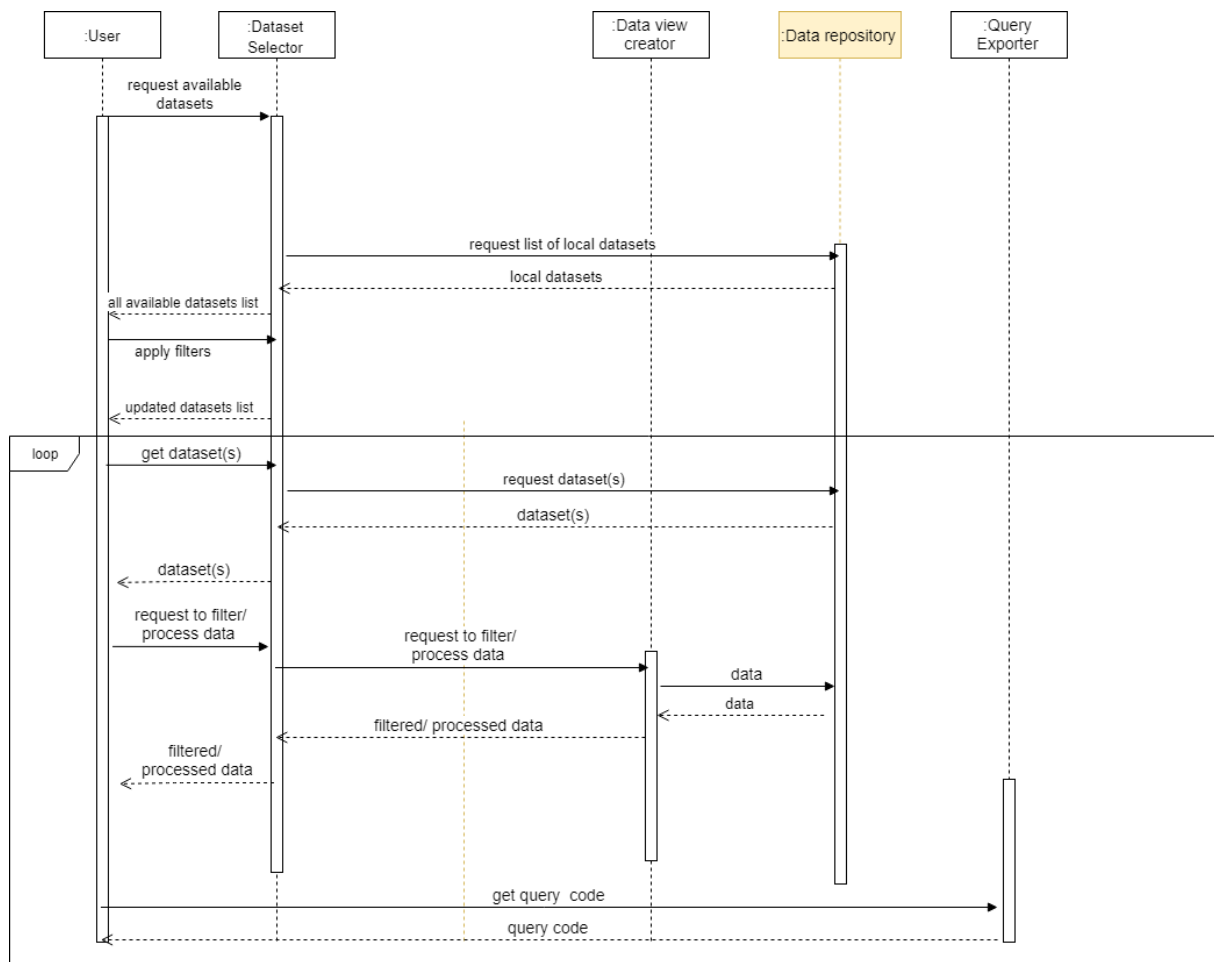
#### 3.7.1. Overview

Query Builder is the component that provides the capability to interactively define and execute queries on data available in the AEGIS system. Query Builder is primarily addressed to the AEGIS users with limited technical background, but is potentially useful for all, as it will simplify and accelerate the process of retrieving data and creating views on them, which could be then saved as new datasets or used as input for more high-level AEGIS tools, like the Visualiser and the Algorithm Execution Container.

The tool is developed as interactive notebook inside Jupyter, offering intuitive data browsing, selection and manipulation. As explained in section 3.2, the functionalities of Query Builder are not limited to the retrieval and combination of various datasets, but also support certain necessary processing tasks that cannot be known a priori, in terms of data filtering and cleansing. Thus, the tool incorporates functionalities that may conceptually be more relevant to the data cleansing tool. However, by integrating them in the current tool, there is a two-fold advantage: (a) the user is offered a more intuitive workflow, since data cleansing requirements may be not known prior to and independently of the query creation process and (b) the computational power of the AEGIS system is fully leveraged, as cleansing may be a very heavy process when dealing with big data.

The high-level functionalities offered by the Query Builder user interface are as follows:

- Dataset browsing
- Dataset selection and data preview
- Dataset merging and appending
- Data filtering, both row-wise and column-wise
- Various data manipulation and cleansing tasks, e.g. value replacement, fill-in of null values, changing column names, combining columns, removing duplicate entries etc.
- Save created view on data as new dataset
- Export the Python code that can be used to achieve the same data manipulation results that were created through the user's interaction with the UI
- Provide the Spark/Pandas DataFrame that corresponds to the created data view as input to Visualiser and/or Algorithm Execution Container
- Provide the Spark/Pandas DataFrame that corresponds to the created data view to the tech-savvy user that wants to directly use it in his/her code



**Figure 3-7: Query building and execution workflow**

**Updates from V3.0:**

- Metadata repository support was removed
- Support for Zeppelin notebook was removed as it is no longer part of the AEGIS integrated services, as described in section 3.6.3.

### 3.7.2. List of microservices

Query Builder is one of the components developed inside a notebook, as explained in Section 2.3. As such, the microservices of which it is composed correspond to specific functionalities also implemented inside the same note of the notebook. The microservices interact directly through Python code, as, in the normal workflow, they are all executed as parts of the same underlying process/job. Hence, the distinction of the five underlying microservices mostly corresponds to the conceptually separate tasks that are performed by each of them and the fact that in another context, i.e. externally to the notebooks, they would constitute different services.

The first and last microservices (namely the Dataset Selector Service and the Query Exporter Service) correspond directly to sub-components of the Query Builder, which are shown in the

Sequence diagram in Figure 3-7. The other three microservices (Cleanser Service, Merger Service and Dataset Creator Service) are integrated under the Data view creator part of the Query Builder (also shown in the corresponding sequence diagram).

Component Name	Microservice Name	Functionalities
Query Builder	Dataset Selector Service	<ul style="list-style-type: none"> <li>Acquire the list of available and accessible datasets from HopsFS of AEGIS Data Store</li> <li>Adjust the list of datasets shown to the user based on the performed choices</li> <li>Retrieve the selected dataset from the filesystem and load it into a DataFrame</li> </ul>
	Dataset Cleanser Service	<ul style="list-style-type: none"> <li>Remove/Replace missing values</li> <li>Perform data interpolation</li> <li>Provide a preview of the applied actions</li> <li>Provide aggregations and other statistics that help examine data integrity</li> <li>Apply rule-based data transformations</li> </ul>
	Dataset Merger Service	<ul style="list-style-type: none"> <li>Merge/Join datasets</li> <li>Perform approximate joins</li> </ul>
	Dataset Creator Service	<ul style="list-style-type: none"> <li>Apply aggregations on datasets</li> <li>Select/Drop columns</li> <li>Apply value replacing</li> <li>Rename columns</li> <li>Perform data interpolation</li> <li>Apply selectors/filters to dataset to refine the retrieved data</li> <li>Save a dataset as a file in the filesystem</li> <li>Provide a preview of the applied actions</li> <li>Load created dataset in a DataFrame</li> </ul>
	Query Exporter Service	<ul style="list-style-type: none"> <li>Translate data processing/filtering/merging actions performed so far into Python code that can be used externally to the tool to produce the same results</li> <li>Export dataset to new file</li> </ul>

**Table 3-33: Query Builder list of microservices**

**Updates from V3.0:**

- Metadata repository support was removed.

*3.7.3. Technologies to be used*

Query Builder is developed as a preconfigured notebook in Jupyter. The Jupyter version of the tool is implemented in Python and PySpark for the data processing and JavaScript for the user interface.

In order to provide effective big data querying and processing functionalities, Query Builder leverages the power of Apache Spark, which is available inside AEGIS Integrated Services (presented in section 3.6).

**Updates from V3.0:**

- Support for Zeppelin notebook was removed as it is no longer part of the AEGIS integrated services, as described in section 3.6.3.

*3.7.4. APIs and exposed outcomes***Updates from V3.0:**

- Support for Zeppelin notebook was removed as it is no longer part of the AEGIS integrated services, as described in section 3.6.3.

Query Builder has two types of exposed outcomes:

1. The Python code that corresponds to the actions performed by the user through the tool's user interface. The generated code can be used then by the user independently in order to achieve the same results without having to repeat the performed steps and can also be directly edited by more tech-savvy users.
2. The DataFrame that contains the data view that was created through all the data manipulation tasks performed by the user. The DataFrame can be passed to the Visualiser and the Algorithm Execution Container or be directly manipulated inside Jupyter through the user's custom code.

Although these are the main outcomes of the tool, there are also two more possible outcomes that may be produced through the user's usage of the tool:

1. New files may be created and stored in the local filesystem
2. The Jupyter notebook that is created may itself serve as an outcome if the user chooses to save and keep it for future reference.

## 3.8. Visualiser

### 3.8.1. Overview

The Visualiser is the component enabling the visualisation capabilities of the AEGIS platform for the output of the querying and filtering results coming from the Query Builder as well as the output of the analysis results as produced from the Algorithm Execution Container. More specifically, the Visualiser is undertaking the necessary actions to address the advanced visualisation requirements of the AEGIS platform by offering a variety of visualisation formats, which span from simple static charts to interactive charts offering several layers of information and customisation.

The Visualiser is implemented as predefined Jupyter notebook, which is part of the AEGIS Integrated Services enabling the interactive web-based notebook functionality in the AEGIS platform. The Visualiser component consists of a set of functionalities which support the execution of the visualisation process. This set of functionalities includes the dataset selection, the dataset preview generation, the visualisation type selection, the visualisation configuration, the visualisation generation and the interactive dashboard. In the following list, the functionalities of the Visualiser are elaborated:

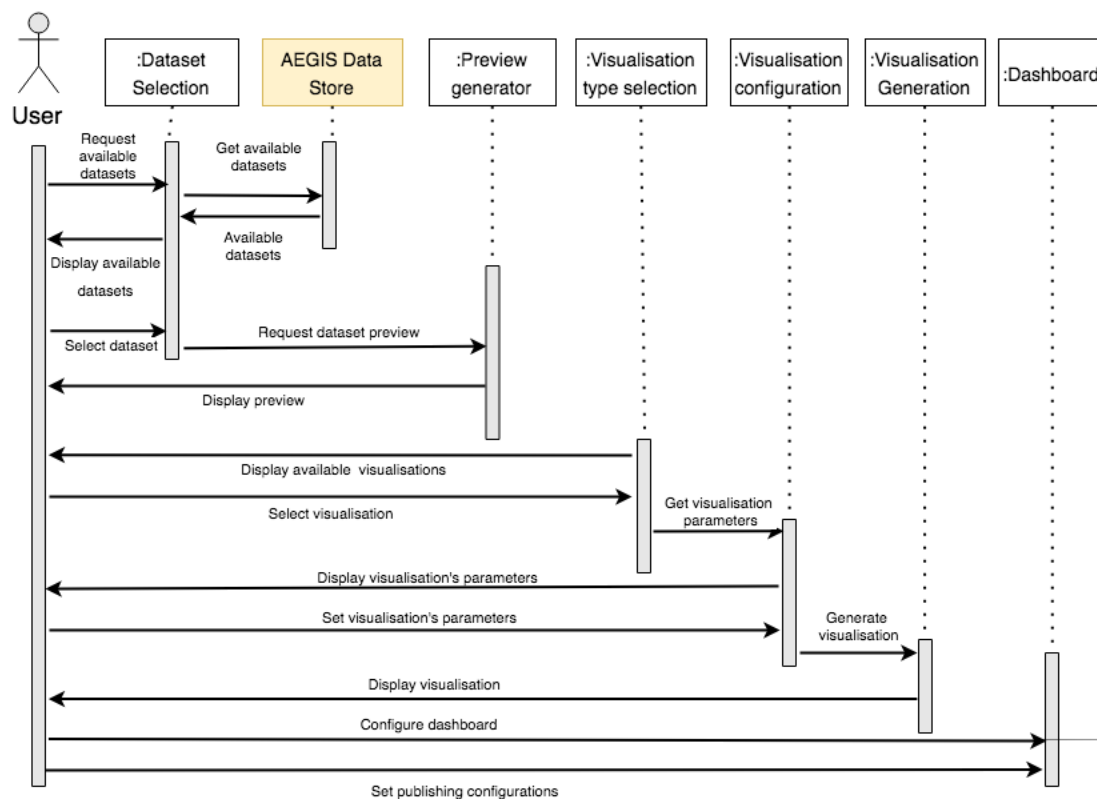
- Dataset selection: The list of available datasets within the project are presented to the user for selection<sup>13</sup>.
- Dataset preview generation: Upon selecting the dataset, a preview of the dataset is displayed<sup>13</sup>.
- Visualisation type selection: The list of available visualisations for the selected dataset is presented to the user for selection.
- Visualisation configuration: Based on the visualisation type selected for the desired dataset a set of parameters are displayed to the user to trim the visualisation.
- Visualisation generation: Once the visualisation type along with the parameters are set for the desired dataset, the visualisation generation is triggered. The results can be used in the current session for creating an interactive dashboard.
- Dashboard: The result of the visualisation generation is presented to the user into an interactive dashboard. This dashboard can contain also multiple generated visualisation results.

Figure 3-8 depicts the execution of the visualisation process.

---

<sup>13</sup> This functionality is hidden in the holistic notebook where the AEGIS notebooks are integrated, as described in Section 2.3. Within this holistic notebook, the Visualiser will receive the input for visualisation directly from the Query Builder or the Algorithm Execution Container.





**Figure 3-8: Sequence diagram of the visualiser component**

Besides the basic functionalities described above, the Visualiser incorporates two additional functionalities. The first one is the exception handling that is supported at each step of the process described above. If any failure occurs during the visualisation process, the user is informed with the proper message. The second one is the support of exporting the generated visualisation in the form of an HTML file that can be saved as a new asset in the user's project in the AEGIS Data Store.

The Visualiser supports a variety of visualisation types. Especially in the case of Maps visualisations, the Visualiser is offering advanced visualisations with the support for Heatmaps on top of Maps, as well as the enhanced map visualisation with support for markers with custom labels and colours and the support for FastMarkerCluster on Maps. Additionally, the Visualiser is providing an enhanced visualisation for Data Tables. The Visualiser supports the following visualisation types:

- Scatter plot
- Pie chart
- Bar chart
- Line chart
- Box plot
- Histogram
- Time series
- Heatmap
- Bubble chart

- Map (with support for HeatMaps on Maps, markers with custom labels and colours and FastMarkerCluster)
- Data Tables

#### Updates from V3.0:

- Support for exception handling during the execution of the visualisation process.
- Support for exporting the generated visualisations as an HTML file that can added as an asset in the user's project.
- Advanced Map visualisations with the support for Heatmaps on top of Maps, as well as markers with custom labels and colours and FastMarkerCluster on Maps.

### 3.8.2. List of microservices

The Visualiser component, as explained in Section 2.3, is developed as a predefined Jupyter notebook and it is following the microservices architecture. The designed microservices are enabling the advanced visualisation capabilities of the AEGIS platform and are orchestrated towards the execution of the visualisation process, as described in the previous section. Each microservice is assigned with a specific functionality within the visualisation process and is implemented as a note in the same notebook, interacting through Python code.

In particular, the Dataset Selection and the Preview Generator, as shown in Figure 3-8, are undertaken by DatasetSelectionService microservice. Additionally, the microservice VisualisationSelectionService is responsible for the Visualisation Type Selection and the Visualisation Configuration. The Visualisation Generation is handled by the ChartBuildingService and ChartCreationService microservices. Finally, the microservice VisualisationService is handling the Dashboard functionality.

In total five microservices are developed and are described in the following table:

Component Name	Microservice Name	Functionalities
Visualiser	DatasetSelectionService <sup>14</sup>	<ul style="list-style-type: none"> <li>• Acquire the list of available and accessible datasets from HopsFS of AEGIS Data Store</li> <li>• Provide the list of available datasets for selection</li> <li>• Provide a preview of the selected dataset</li> </ul>

<sup>14</sup> This microservice is not used in the holistic notebook where the AEGIS notebooks are integrated, as described in Section 2.3. Within this holistic notebook, the Visualiser will receive the input for visualisation directly from the Query Builder or the Algorithm Execution Container.

	VisualisationSelectionService	<ul style="list-style-type: none"> <li>• Provide the list of available visualisation types</li> <li>• Provide and manage the parameters (such as axis variables and titles) for each visualisation type</li> </ul>
	ChartBuildingService	<ul style="list-style-type: none"> <li>• Prepare the data in the appropriate format based on the selection of the visualisation type and set parameters</li> </ul>
	ChartCreationService	<ul style="list-style-type: none"> <li>• Generate the appropriate visualisation based on the data, visualisation type and parameters</li> </ul>
	VisualisationService	<ul style="list-style-type: none"> <li>• Display the generated visualisation as a UI component</li> </ul>

**Table 3-34: Visualiser list of microservices****Updates from V3.0:**

- No updates were introduced.

*3.8.3. Technologies to be used*

As already described the Visualiser component is implemented as a predefined Jupyter notebook, a multipurpose interactive web-based notebook service for running Spark code on Hops YARN, which is part of the AEGIS Integrated Services. Jupyter offers functionalities for data visualisation out-of-the box in addition to data ingestion, data discovery and data analytics functionalities. In addition to Jupyter, the user interface is implemented using Python, JavaScript and HTML with the support of two Python libraries, namely the Folium<sup>15</sup> and the highcharts<sup>16</sup> libraries. These specific libraries were selected as they provide state-of-the-art visualisations specialised in charts and data visualisation. Additionally, the Qgrid Jupyter widget<sup>17</sup> is utilised in order to enable support for Data Tables with intuitive scrolling, sorting and filtering controls.

<sup>15</sup> <http://folium.readthedocs.io/en/latest/>

<sup>16</sup> <https://www.highcharts.com/>

<sup>17</sup> <https://github.com/quantopian/qgrid>

**Updates from V3.0:**

- The Jupyter widget Qgrid is utilised for enhanced Data Tables visualisation support.

*3.8.4. APIs and exposed outcomes***Updates from V3.0:**

- The produced visualisation can be saved also as an HTML file.

The Visualiser is providing the generated visualisation as an exposed outcome. The Visualiser is generating the visualisation tailored by the user, taking as input either the results of the query processing as facilitated by the Query Builder or the analysis results that are provided as the outcome of Algorithm Execution Container. The produced visualisation can be saved as an image, as an HTML file or introduced in an interactive dashboard.

**3.9. Algorithm Execution Container***3.9.1. Overview*

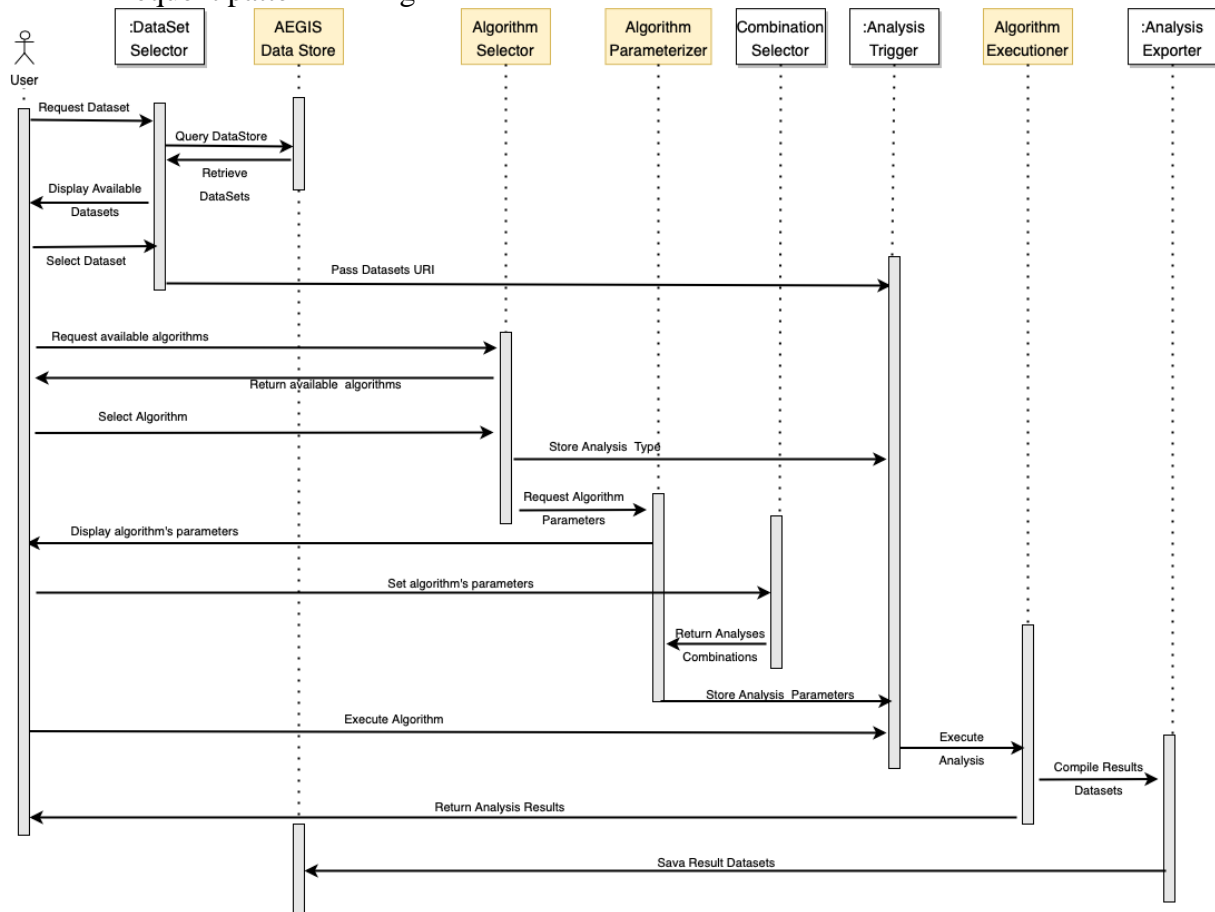
Analytics in AEGIS are to be constructed with the use of the Algorithm Execution Container, which is a module that runs on top of a web-notebook and can be used either on its own, or as a follow-up to the Query Builder notebook.

With the aim to provide extra functionalities to both novel and non-expert users, this component features a UI that consists of an algorithm selection template, offering to users some basic information regarding each algorithm available in the big data analysis platform of AEGIS.

Initially, the user has to select a Dataset which will be the basis for the analysis, and the “Dataset Selector” microservice is triggered to retrieve the dataset from the AEGIS Data Store. Following this, the user proceeds with the selection of an algorithm (out of an algorithm family), specific parameters of each algorithm are presented, to provide to users the option to fill in all variables of the algorithm and perform an analysis over the platform. In case a user selects so, he can utilise a parameter grid where range of values are selected instead of exact values. These actions are collected by the “Analysis Trigger” microservice, which is composed of a series of nested or interlinked notebook paragraphs and that is passing over the analysis request to the underlying Spark engine. In the case of the grid parameter choice, multiple combinations of the parameters are chosen for execution. The output of an analysis is then generated by the Algorithm Execution Container and the performance of each algorithm is being previewed in the same notebook and is saved back into the AEGIS Data Store using the “Analysis Exporter” microservice. At the same time, error messages may appear in case errors occurred.

The current design of the Algorithm Execution Container supports the following categories of analyses, and at presents counts 20 algorithms:

- Dimensionality Reduction/Feature Extraction/Selection
- NLP Functions
- Recommenders
- Clustering
- Classification/Regression
- Frequent pattern mining



**Figure 3-9: Algorithm Execution Container sequence diagram**

#### Updates from V3.0:

- New algorithms included and redesigned process for Jupyter notebook
- Files with executed analysis are richer in terms of information stored
- Error Messages included in the output text areas
- Parameter grid included for allowing algorithms to run on a range of values for improved model training
- Incorporation of feature to apply models on unseen values

### 3.9.2. List of microservices

The Algorithm Execution Container is a component that is developed inside the Jupyter Notebook and is part of the overall analysis functionality of the platform. The microservices interact with the backbone through Python code.

Component Name	Microservice Name	Functionalities
Algorithm Execution Container	Dataset Selector	<ul style="list-style-type: none"> <li>Interacts with the AEGIS storage and directly selects a dataset</li> </ul>
	Parameter Combination Constructor	<ul style="list-style-type: none"> <li>Gets input from the parameters grid and constructs combinations of parameters to run multiple versions of the same analysis</li> </ul>
	Analysis Trigger	<ul style="list-style-type: none"> <li>Selects the analysis to be performed by the user</li> <li>Passes the analysis parameters to the analytics function</li> <li>Triggers the analysis to be performed in the AEGIS Spark</li> </ul>
	Analysis Exporter	<ul style="list-style-type: none"> <li>Stores the created analysis model in AEGIS Data Store</li> <li>Stores the dataframes that contain the results of the model application on the input data in AEGIS Data Store</li> </ul>

**Table 3-35: Algorithm Execution Container list of microservices**

#### Updates from V3.0:

- DataSet Selector now works by directly displaying available datasets instead of URI entry
- New microservices to perform combinations of parameters in case of parameter grid selection by the user (called Parameter Combination Constructor)

### 3.9.3. Technologies to be used

The Front-End of the Algorithm Execution Container is based on AngularJS framework running on top of a Jupyter notebook as already pre-defined paragraphs that present the UI to the user. MLlib is the core algorithm library that is supported. The interoperation between AngularJS and MLlib under the environment of Jupyter is facilitated by Python, that is being used for configuring and executing the algorithms that are selected by the user.

#### Updates from V3.0:

- Completely based on Jupyter notebook
- Support for Zeppelin notebook was removed as it is no longer part of the AEGIS integrated services, as described in section 3.6.3.
- Python used as the underlying algorithm execution language, instead of Scala which was used for the Zeppelin notebook.

#### 3.9.4. APIs and exposed outcomes

##### **Updates from V3.0:**

- No updates were introduced.

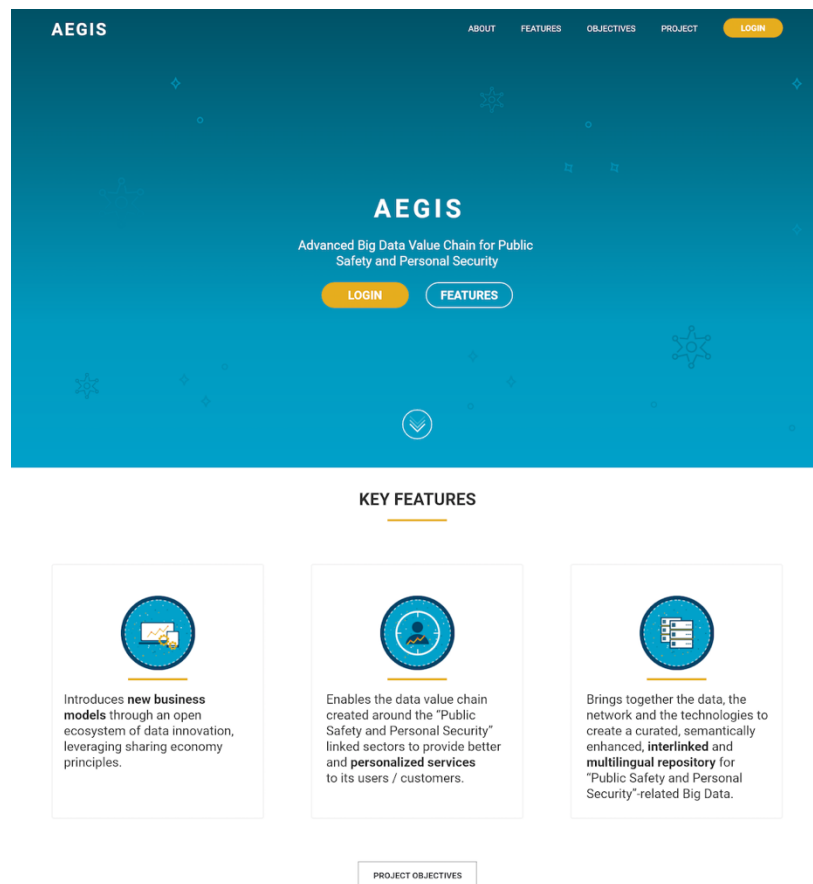
No APIs are being provided by this component, as it directly interacts with the Jupyter Notebook. The pre-generated Python code that corresponds to the actions performed by the user through the tool's user interface can be used then by the user independently in order to achieve the same results without having to repeat the performed steps and can also be directly edited by more tech-savvy users. The outcomes of the component are passed back to the AEGIS storage facility, if the user chooses to save and keep them for future reference.

### 3.10. AEGIS Front-End

#### 3.10.1. Overview

The AEGIS Front-End is the upper layer of the whole AEGIS architecture, receiving and sending the outputs/inputs from/to the AEGIS API layer. The AEGIS Front-End received a restyling in order to address the feedback collected and simplify the user interface towards the aim of making it more effective and complete. For this reason, a series of mock-ups of the AEGIS Front-End were designed that are driving the implementation of the redesigned AEGIS Front-End. The Landing Page has been introduced in order to introduce the platform to the users, providing a short overview of the platform and its key features, as well as a small description of the AEGIS project (see Figure 3-10).

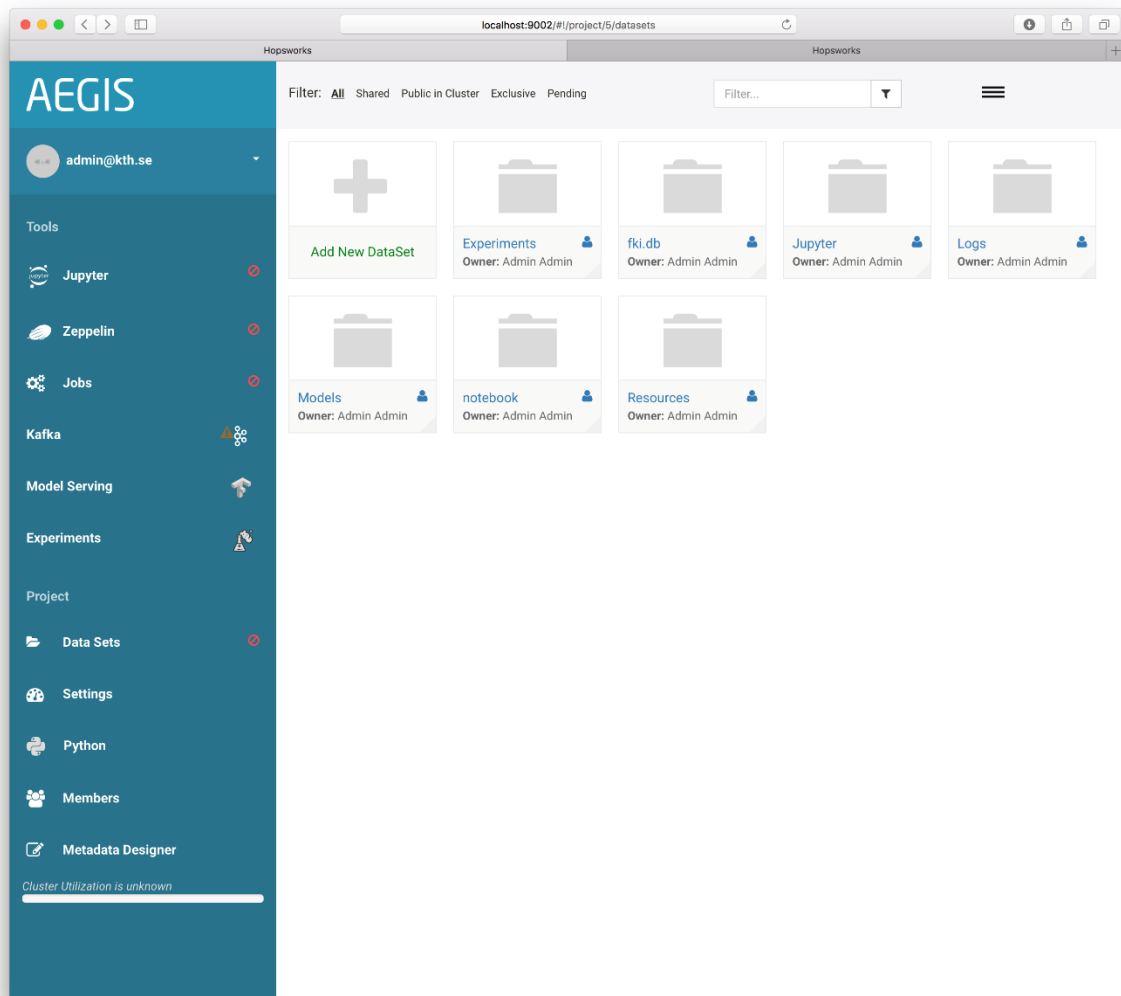
The first step to access the platform is the creation/authentication of an account. Users can be assigned with the roles of data owner/data scientist, with different permissions for the management of the projects/datasets. Through the main page, that is now principally focused on the projects, the users are able to browse, search or create new projects with the option to specify the related members. Moreover, in the main page the user will be able see additional information such as a Getting Started tutorial and to explore the most recent news about the AEGIS project. In addition to this, the main page offers a full search functionality.



**Figure 3-10: Mock-up of the AEGIS platform landing page**

The Front-End facilitates all the AEGIS components which have an interaction with the user (e.g. Query Builder, Visualiser, Algorithm Execution Container). A common graphical environment has been improved, still according to the look and feel of the AEGIS institutional web site, including direct links to the single web components corresponding to the AEGIS main functionalities, here integrated in the form of notebooks. In particular, the main menu of the AEGIS platform, as depicted in Figure 3-11, for each project presents the following items: Tools (Jupyter, Jobs), Kafka, Model Serving, Experiments, Project (Data Sets, Settings, Python, Members, Metadata Designer). The applications related to Queries, Visualisations and Analytics have been implemented within Jupyter in the form of notebooks.





**Figure 3-11: Mock-up of the main menu of the AEGIS platform**

As described above, “Getting Started” feature is provided within the platform, introducing the users to the different tools offered by AEGIS in order to address the diverse needs of its users. At the same time, a “Getting Started” workflow is also provided, describing the “first visit” step by step and introducing the advanced options offered by the platform. Moreover, a detailed user guide has been attached to each of the following tools: Query Builder, Visualiser, Algorithm Execution Container.

Finally, as explained in detail in section 3.11, the platform will feature support to multilingualism on multiple levels. This will impact on the AngularJS Frontend, which will integrate translations of the static content of the platform as well as an online machine translation service, allowing an on-the-fly translation of data.

#### **Updates from V3.0:**

- Full restyling of the graphical user interface has been designed with the help of mock-ups and it is currently under development.

### 3.10.2. List of microservices

A list of microservices have been developed in order to handle the selection/creation of a project within the repository. In total two microservices have been developed and are described in the following table:

Component Name	Microservice Name	Functionalities
Front-End	GetProjects	<ul style="list-style-type: none"> <li>Return the sorted list of all the available projects in the platform</li> </ul>
	CreateProject	<ul style="list-style-type: none"> <li>Call the service for the creation of a new project within the platform</li> </ul>

**Table 3-36: AEGIS Front-End list of microservices**

#### Updates from V3.0:

- Two new microservices were introduced, namely the GetProjects and the CreateProjects, to support the new functionalities of the AEGIS Front-End.

### 3.10.3. Technologies to be used

The AEGIS Front-End is built on top of the Hopsworks platform. Hopsworks is a self-service User Interface for Hops Hadoop, which introduces new concepts needed for project-based multi-tenancy: projects, users, and datasets. All jobs and interactive analyses are run from the HopsWorks UI and Jupyter Notebooks, an iPython notebook style web application. While developing the AEGIS platform, a central role has been taken by the notebook technology, providing the technology required to implement in particular the following components: Query Builder, Visualiser and Algorithm Execution Container. Jupyter is the selected notebook framework (more details in section 3.6).

Another important framework which has been used for the development of the graphical user interface is AngularJS<sup>18</sup>. AngularJS is a very powerful JavaScript based development framework to create Rich Internet Application (RIA<sup>19</sup>). It is used mostly in Single Page Application (SPA<sup>20</sup>) projects. It extends the HTML DOM with additional attributes and makes it more responsive to user actions. AngularJS is open source, completely free and used by thousands of developers around the world. It is licensed under the Apache License version 2.0. Applications written in AngularJS are cross-browser compliant. AngularJS automatically

<sup>18</sup> <https://angularjs.org/>

<sup>19</sup> [https://en.wikipedia.org/wiki/Rich\\_Internet\\_application](https://en.wikipedia.org/wiki/Rich_Internet_application)

<sup>20</sup> [https://en.wikipedia.org/wiki/Single-page\\_application](https://en.wikipedia.org/wiki/Single-page_application)

handles JavaScript code suitable for each browser and allows to implement the Model-View-Controller (MVC<sup>21</sup>) pattern on the client side using JavaScript.

**Updates from V3.0:**

- No updates were introduced.

*3.10.4. APIs and exposed outcomes*

Not applicable.

**3.11. Multilingualism Support***3.11.1. Overview*

The AEGIS platform is built to handle cross-domain and cross-country datasets. Therefore, it is required to support multilingualism on multiple levels in order to facilitate an easier data discovery and data integration process. It should be noted at this point that the multilingualism support is referring to the set of technologies and tools that are utilised within the components of the AEGIS platform in order to enable multilingualism, rather than a stand-alone component of the platform.

*3.11.2. Approach*

Multilingualism support can be basically differentiated between three different entities, where it should be applied: static content, metadata and data. Those different entities need to be dealt with separately in order to accommodate their respective characteristics:

**Static Content**

Static content refers to all fixed literals or editorial content of the AEGIS platform, especially the frontend. This includes menu items, buttons, labels etc. This content needs to be translated and integrated into the frontend accordingly. For the scope of this project the translations will be limited to the most important menu items and labels. The AEGIS platform will provide the means for offering the static content in the languages: English German, Greek, Italian and Swedish.

**Metadata**

The multilingualism of the metadata poses the most important aspect, since it supports users in finding suitable datasets. The metadata is highly dynamic and with potentially thousands of datasets a static translation not feasible. Hence, an automatic process will be applied utilising available machine translation services. On creation or update time of the metadata, literals will

---

<sup>21</sup> <https://en.wikipedia.org/wiki/Model%E2%80%93view%E2%80%93controller>

be sent to suitable service and its results will be stored. The underlying technology of the AEGIS Metadata Service natively allows the storing of multiple languages for literals. It is planned to support all 24 official languages of the European Union. The translation will be limited to literals of the DCAT-AP properties. If linked entities offer the selected language, an automatic resolution will be attempted. The multilingual metadata will also be indexed by the AEGIS search component and enabling a cross-language search.

## Data

A static or automatic translation of the data itself does not make sense due to its volume. In addition, (automatic) translations may corrupt the underlying and hidden insights of the data. Still, it is desired that users be able to comprehend the content, or at least get a broad idea of the content. Therefore, an online machine translation service will be integrated in the frontend, allowing an on-the-fly translation of a small excerpt of a selected file.

### 3.11.3. Technologies to be used

## Static Content

The translations of the static content will be part of the AngularJS frontend and be applied by using an established module for enabling multilingualism support in AngularJS. The module is named `angular-translate`<sup>22</sup> and can be integrated into the AEGIS frontend without issues.

## Metadata

The general support for multilingual metadata will be implemented by the AEGIS Metadata Service. It will also handle the required calls to a machine translation service and offer a respective callback if needed. For the scope of the project the eTranslation service of the European Commission (EC) will be used<sup>23</sup>. The service offers machine translation for all 24 European languages and can be used by EU institutions and public administrations. The AEGIS project is granted access to the service. An additional middleware is applied to handle the communication with the eTranslation service, which was developed in the scope of the European Data Portal<sup>24</sup>. This middleware acts as an aggregator and persistence layer for translations jobs. However, the integration will be generic in order to support a straight-forward switch of providers if necessary.

## Data

For on-the-fly translation and frontend integration the Microsoft Bing Translator-Widget<sup>25</sup> will be applied and integrated into the AngularJS frontend.

---

<sup>22</sup> <https://angular-translate.github.io/>

<sup>23</sup> [https://ec.europa.eu/info/resources-partners/machine-translation-public-administrations-ettranslation\\_en](https://ec.europa.eu/info/resources-partners/machine-translation-public-administrations-ettranslation_en)

<sup>24</sup> <https://gitlab.com/european-data-portal/translation-service>

<sup>25</sup> <https://www.bing.com/widget/translator>

**Updates from V3.0:**

- Concrete services for machine translation were determined: EC eTranslation and Bing Translator-Widget.

*3.11.4. List of microservices*

The Translation Middleware is composed by one microservices and the following functionalities:

Component Name	Microservice Name	Functionalities
Translation Middleware	translation-service	<ul style="list-style-type: none"> <li>• Handles communication to EC eTranslation service</li> <li>• Creating batches of translations</li> <li>• Handle return callback from to EC eTranslation service</li> </ul>

**Table 3-37 - Translation Middleware microservice**

*3.11.5. APIs and exposed outcomes*

The following table documents the API of the Translation Middleware.

Technical Interface	
Reference Code	MLM#01
Function	Translation Middleware for EC eTranslation
Subsystems	Translation Middleware
Type, State	
RESTful API	
Endpoint URI	
-	
Input Data	
Text data	
Output Data	

Translated text data

**Table 3-38: Translation Middleware technical interface**

### 3.12. Holistic Security Approach

#### 3.12.1. Overview

In the high-level architecture of the AEGIS platform the consortium identified the need for a holistic security approach that should be incorporated throughout the AEGIS platform and that will be applied in the whole lifecycle of the data exploitation. Towards this aim, within the holistic security approach the security aspects of data in storage and data in transit are safeguarded, the platform operations are protected with the use of certificates and the technical interfaces are secured with a token-based authentication. It should be noted at this point that the holistic security approach is not a standalone component, but rather a set of technologies and tools that are utilised within the components of the AEGIS platform in order to enable cross-platform security.

In the AEGIS platform architecture, the main decision taken was the adoption of the Hopsworks<sup>26</sup> platform as the Big Data Processing Cluster of the AEGIS platform. Hopsworks is providing out-the-box the HopsFS, a new implementation of the Hadoop Filesystem (HDFS), covering the storage solution of the AEGIS platform. With respect to the security aspect for data in storage HopsFS is offering advanced security with a plethora of authentication mechanisms as well as data access control, data integrity and data consistency mechanisms. HopsFS is making use of checksum to ensure security and integrity control of the data in storage covering the envisioned by the consortium security aspect for data in storage. Checksum is used to verify the integrity of data by verifying that the data have not been altered or corrupted. The value of checksum is usually calculated using cyclic redundancy check (CRC) algorithms and cryptographic hash functions. The utilisation of checksum is the ideal solution for big data infrastructures and for platforms performing data analysis as it is not introducing any efficiency problems and delays in data handling, while also addressing the security, privacy and integrity of the data stored in the AEGIS Data Store, unlike any available cryptographic solution.

Concerning the security of data in transit or data in motion, which includes data transfer between the Hopsworks services and clients either within the internal network or through the internet, Hopsworks is providing data encryption via Secure Sockets Layer (SSL) and Transport Layer Security (TLS) at the RPC layer offering the required security level as envisioned by the AEGIS consortium.

Before being able to use the platform, a user must register itself and an administrator must approve it. Users authenticate themselves with a username and password. In order to start uploading files or launch jobs to the cluster the user must create a Project or be invited to an already existing one. At this point a *project specific* user is created, which is in the form of *ProjectName\_\_Username* and this will be the effective user in the context of a Project. With

---

<sup>26</sup> <http://hops.io>

this format we can distinguish between the same user belonging to a different Project. All operations in the platform are performed as the project specific user.

There are two types of certificates in Hopsworks:

1. Every node in the cluster is provisioned with *host certificates*. The CommonName (CN) field of the X.509 subject contains the fully qualified domain name (FQDN) of the host. This type of certificates is used by Hops daemons such as NameNode, DataNode, ResourceManager, NodeManager and by operations performed by system's superusers.
2. Every time a user creates a new Project or joins an already existing one, Hopsworks creates *project specific user certificates*. When a user interacts with Hops, either the filesystem or the scheduler, communication is protected using these certificates. The CN in that case will be the project specific username of that user.

Hopsworks comes with its own Certificate Authority (CA) which issues the aforementioned certificates. Finally, Hopsworks is a web application and the application server can be configured for HTTPS connections. During deployment of Hopsworks, we install the *Root CA*, also *HopsCA*. HopsCA does not issue directly the certificates, instead there is an *intermediate CA*, which is signed by HopsCA, and signs all the certificates for internal consumption. There can be multiple intermediate CAs.

The holistic security approach also covers the security aspects for the technical interfaces (e.g. REST) provided by the platform. This includes the interfaces provided by the components of the platform in regards to the authorisation, authentication and access approval mechanisms. For the security of the technical interfaces the consortium decided to introduce a token-based authentication with JSON Web Token (JWT)<sup>27</sup>. JWT is an open standard (RFC 7519) that defines a compact and self-contained way for securely transmitting information between parties as a JSON object. For the introduction of the JWT as the authentication and secure information exchange mechanism, a series of actions were taken. At first, a new service has been introduced and was integrated in the backend of the AEGIS platform, undertaking the token generation upon successful login, as well as the token verification. The methods included within this service implement the access control mechanism that is integrated in the AEGIS platform. In addition to the service, a new filter method has been introduced. This filter is utilised in every technical interface and performs the token verification for each incoming request.

The following table presents the holistic security approach of AEGIS platform for the data lifecycle security as described above.

Security Aspect	AEGIS Holistic Security Approach	Remarks
Data in Storage	HopsFS mechanisms for authentication, authorisation and access control of stored data.	Checksum is utilised as the best solution for the security, privacy and integrity of the data in

<sup>27</sup> JSON Web Tokens, <https://jwt.io/>

	Usage of checksums for data integrity and data consistency.	storage with efficiency and high performance.
Data in Transit	<p>Hopsworks provides Secure Sockets Layer (SSL) and Transport Layer Security (TLS) data encryption and authentication at the RPC layer.</p> <p>The AEGIS platform utilises SSL certificates in order to secure the cross-platform communication and operations.</p>	SSL and TLS encryption are the de-facto standard in the security of data in transit.
Technical Interfaces	A token-based authentication and authorisation mechanism with JSON Web Token (JWT) is utilised in all technical interfaces.	JWT is considered the dominant solution for the authentication and secure exchange information in technical interfaces.

**Table 3-39: Holistic Security Approach summary****Updates from V3.0:**

- No updates were introduced.

*3.12.2. Technologies to be used*

The holistic security approach is based on the following three technologies, one for each aspect of the approach. Concerning the data in storage aspect the usage of checksum was selected. Checksum is a small-sized datum derived as the outcome of the cryptographic hash function or checksum algorithm on a block of data or file. This outcome is utilised to identify data corruption errors or modifications and overall data integrity since even small changes will produce a different outcome.

With regards to the data in transit or data in motion security aspect the SSL and TLS cryptographic protocols are the de-facto standard for secure communication over the network. It ensures the secure connection by eliminating the unauthorised read and modification of the data in transit. TLS is an updated more secure version of SSL, introducing the symmetric cryptography with unique keys based on a shared secret for each connection. Each communicating parties is using a public-key to authenticate and the data integrity evaluation is performed with the use of message authentication code.

For the security of the technical interfaces of the platform JSON Web token (JWT) is utilised for authentication and secure information exchange purposes. JSON Web Token (JWT) is an open standard (RFC 7519) utilising digitally signed using JSON Web Signature (JWS) and/or



encrypted using JSON Web Encryption (JWE) JSON objects as a safe way to represent a set of information between two parties. As a consequence, this token is composed by a header, a payload and a signature. JWT is used for authentication purposes, as the token produced during the login authentication is defining the access level for the routes, services and resources of the platform. JWT can be also used for secure information transfer between communication parties.

**Updates from V3.0:**

- No updates were introduced.

*3.12.3. API*

Not applicable.

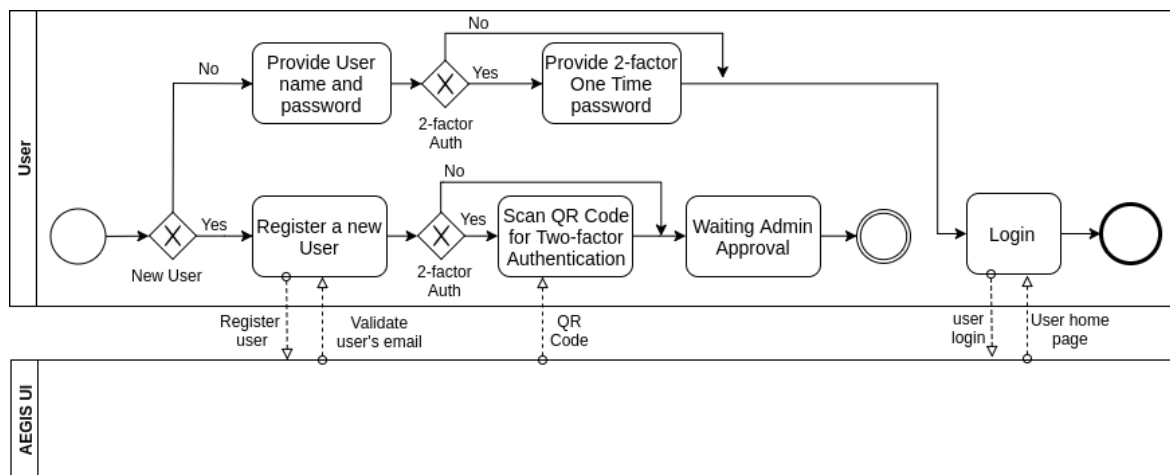
## 4. USER INTERACTION WORKFLOWS

In the current section, the main workflows that facilitate the data-driven innovation in the PSPS domains are presented, as documented in section 4 of deliverable D3.2, with the necessary adaptations based on the updates of the components of the AEGIS platform.

All workflows are focusing on the user perspective and the purpose of this section is to hide the technical details on how the AEGIS components are interacting and on the internal processes of each component but rather illustrate the provided functionalities of AEGIS platform. All workflows are modelled in BPMN diagrams and on each workflow, a specific functionality is presented involving one or more components described in section 3. By chaining and combining these workflows, all AEGIS scenarios and identified user requirements are covered.

### 4.1. Sign-up and Login

The figure below shows the interactions between a user and the AEGIS user interface. A new user can create a new account by providing his/her name and password, and then wait for admin approval before being able to use the platform. In addition, a 2-factor authentication password could be used if enabled.

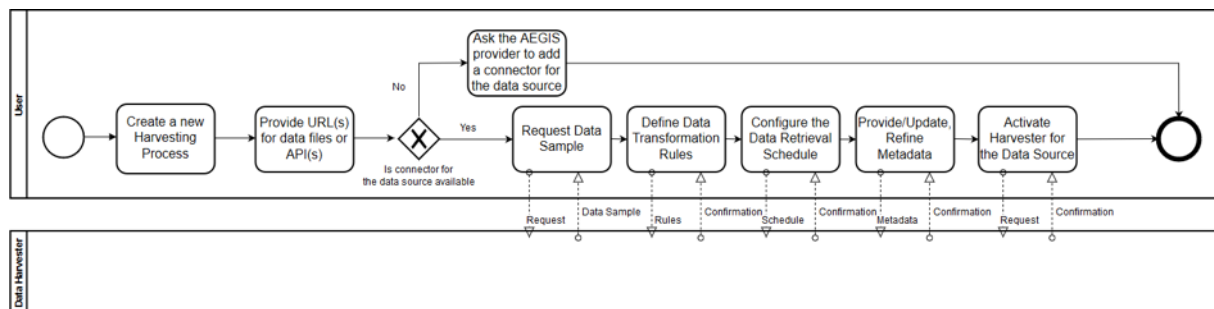


**Figure 4-1: Sign-up and Login workflow**

### 4.2. Data import

#### 4.2.1. Importing data for a new dataset

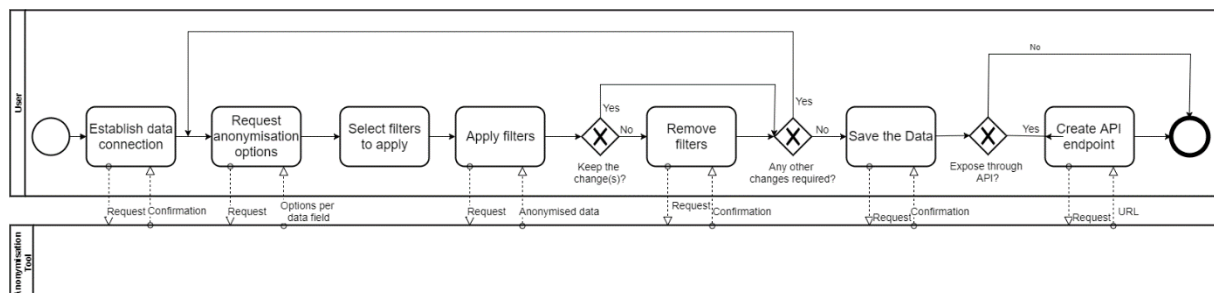
The figure below presents a workflow of user interaction with the AEGIS harvester for registering a new dataset in AEGIS. The workflow shows the required user actions for configuring the harvester for a dataset metadata registration/import as well as its data import and transformation to the target format.



**Figure 4-2: Importing data and metadata and registering them as a part of a new dataset**

#### 4.2.2. Anonymisation workflow

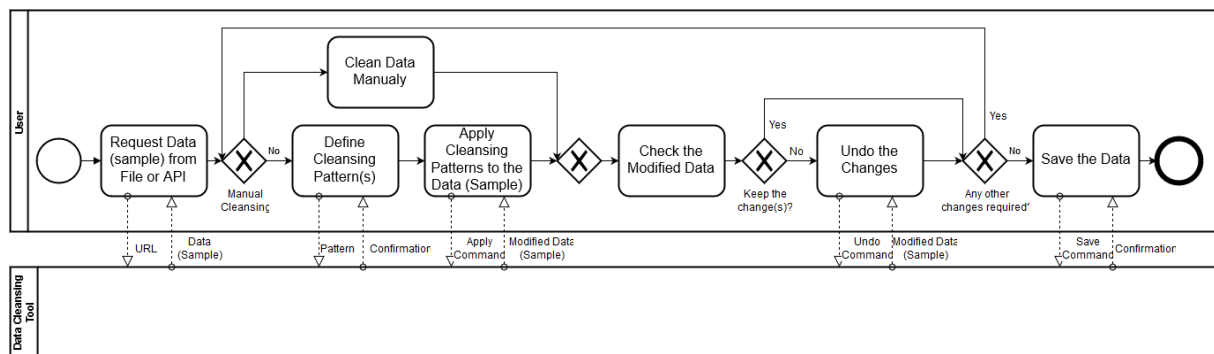
As explained in Section 3.3, anonymisation is performed offline, prior to uploading any potentially sensitive data to the core AEGIS platform. Anonymisation may be an iterative process, as several actions may be required until all personal information has been stripped off the original dataset. The figure below shows the actions undertaken by the user in order to anonymise their data through the provided anonymisation tool, prior to importing them to the AEGIS web platform.



**Figure 4-3: Data anonymisation workflow**

#### 4.2.3. Data cleansing workflow

As explained in section 3.2, AEGIS offers both online and offline cleansing functionalities, which may span from simple value replacements to more complex and computationally intense data manipulations. The offline data cleansing is performed through a dedicated AEGIS offline tool that offers an intuitive data cleansing workflow which is presented in the following figure.



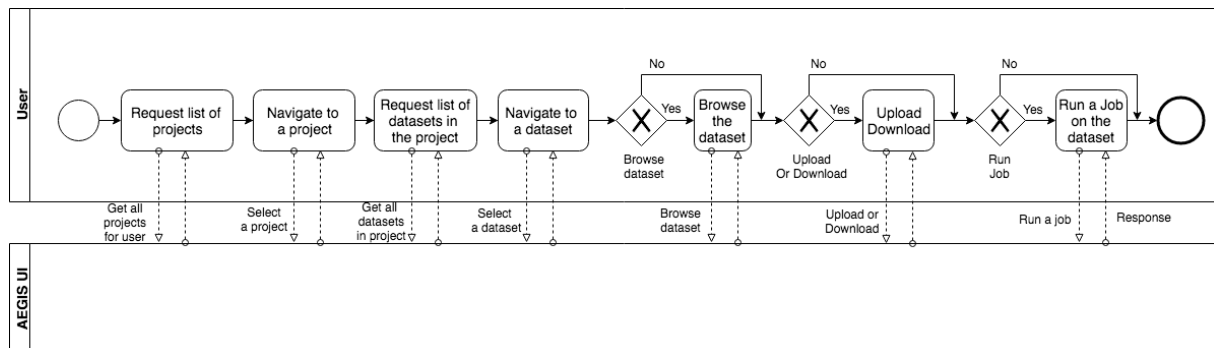
**Figure 4-4: Data cleansing workflow**

Regarding the online data cleansing, due to the flexibility offered by the Notebooks, there is no unique workflow to follow. However, the workflows that include the usage of notebook-based components (e.g. the ones in sections 4.3.2 and 4.6 ) provide some insights on the expected user interaction.

### 4.3. Data and service exploration (search)

#### 4.3.1. From the main AEGIS platform

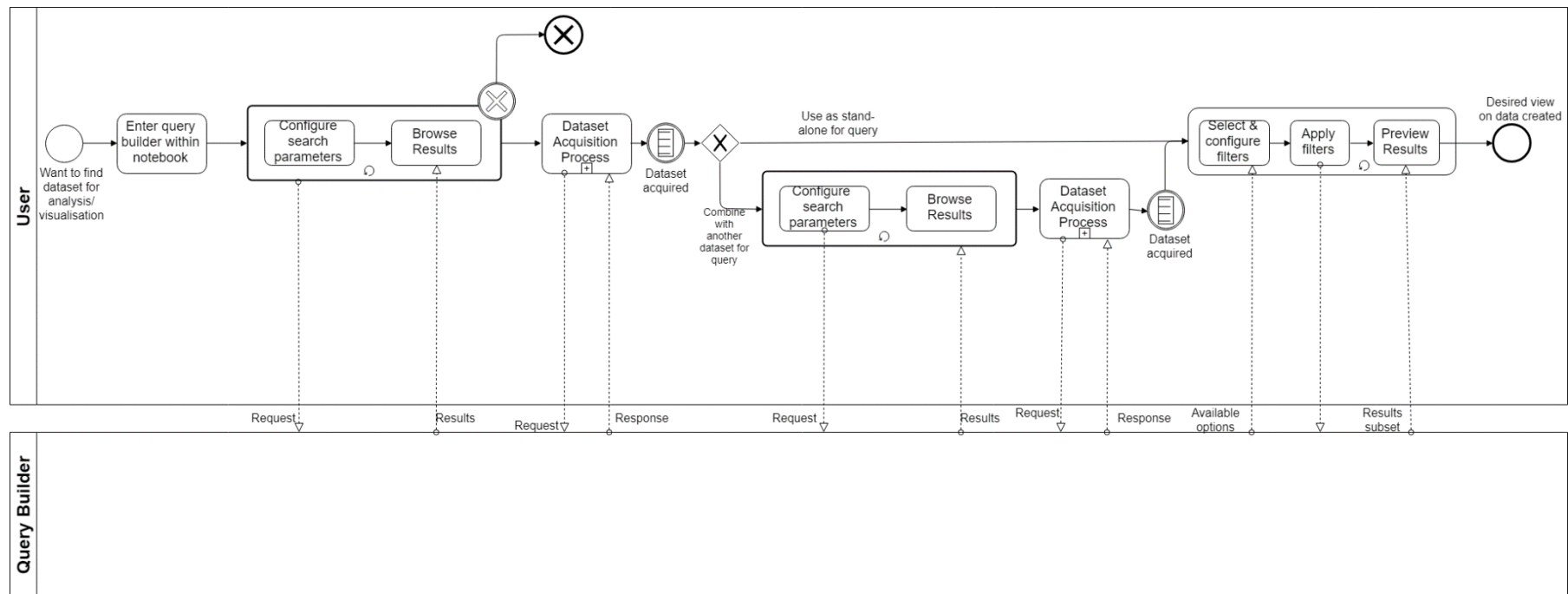
The figure below presents the main actions the users can take to explore the data on the AEGIS platform. The user can request all his/her projects and datasets, and navigate to any project or dataset. Once in a dataset, the user can browse, upload, or download files.



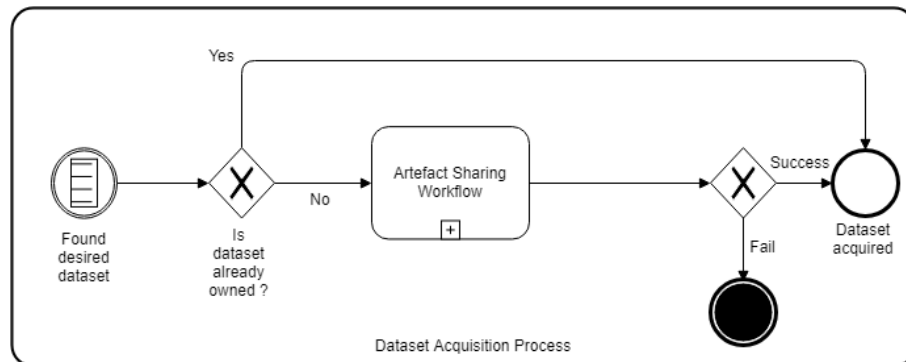
**Figure 4-5: Data and service exploration workflow**

#### 4.3.2. Using query builder

The following two figures show the user's perspective when using query builder to find data and create an appropriate dataset (more accurately create a view on selected data) to be fed into analysis and/or visualisation or to be saved as a new dataset. The process of "creating an appropriate dataset" includes also the cleansing functionalities that have been integrated in the Query Builder (through the selection, configuration and application of the available filters). Although the user primarily interacts with the Query Builder component, other components are utilised in the background. The Brokerage Engine is involved in the dataset acquisition sub-process shown in the diagram, which is an instance of the artefact sharing process described in Section 4.5 and is external to the Query Builder utilisation process.



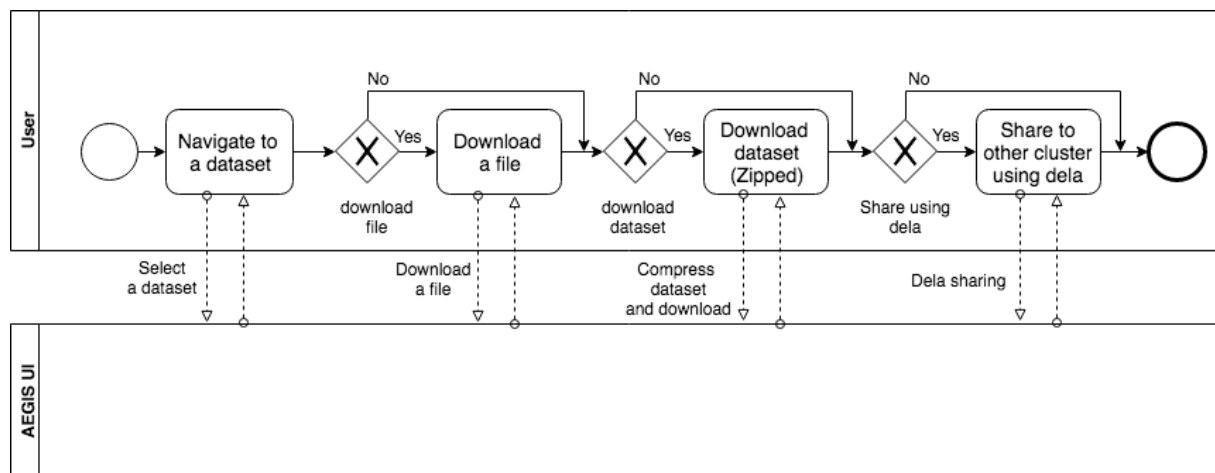
**Figure 4-6: Dataset exploration through query builder workflow**



**Figure 4-7: Data acquisition sub-process workflow**

#### 4.4. Data export from AEGIS

The figure below presents the different ways for the user to export their data from the AEGIS platform. The user can download his/her files or share the whole dataset with other users within AEGIS.



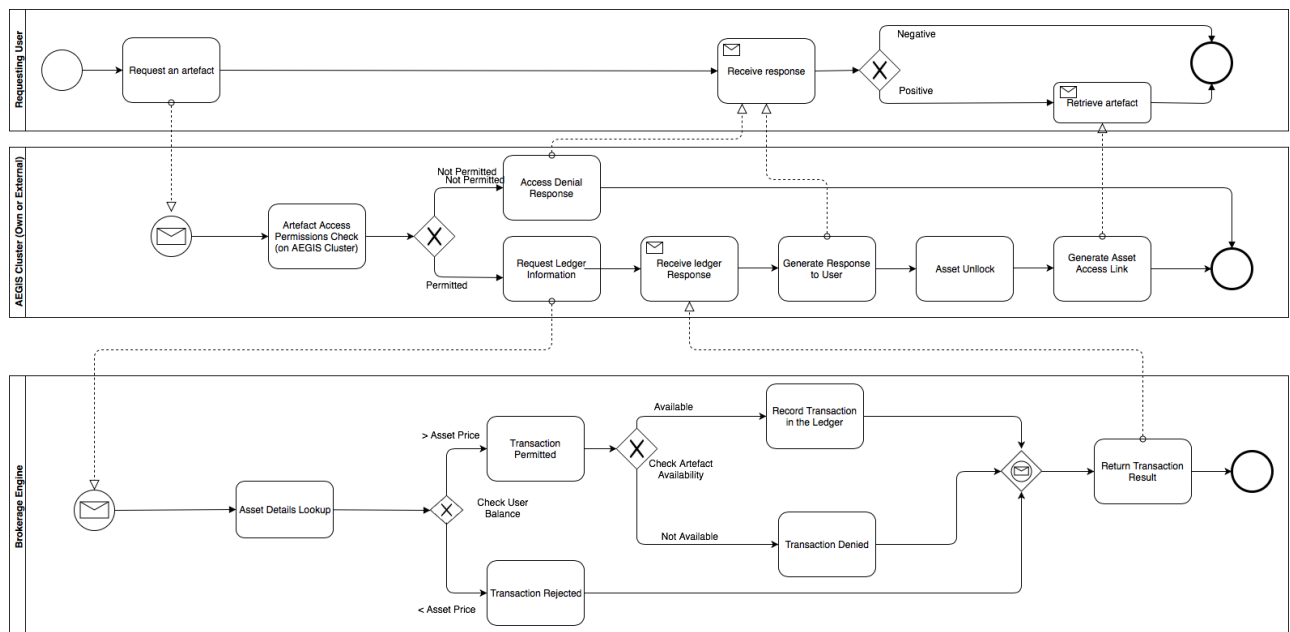
**Figure 4-8: Data export workflow**

#### 4.5. Artefact sharing/reuse

The following figure shows the workflow for data sharing over the AEGIS platform. The operation to be performed, refers only to assets that are not public/free on the platform. It involves both the core AEGIS platform as well as the Brokerage Engine, which will check if artefact sharing/reuse can be performed. At first level, the AEGIS platform checks whether the operation at high level is permitted (e.g. if the data asset exists, if the user has the right credentials to view the data artefact, if the user is logged in, etc.). If access is possible and is permitted, then the Brokerage Engine is invoked. The Brokerage Engine checks the ledger to resolve the following situations:

- Identify whether the user requesting the data is capable of receiving it (e.g. if he/she has enough “coins” in case the data asset is not free), and
- Verify the availability of the data asset, comparing previous records in the ledger with the DPF elements that are attached to the data asset. This is essential only in case a data asset is provided with exclusivity rights (either permanently or within a specific timeframe), so that there is a check that no exclusivity rights have been transferred at the moment.

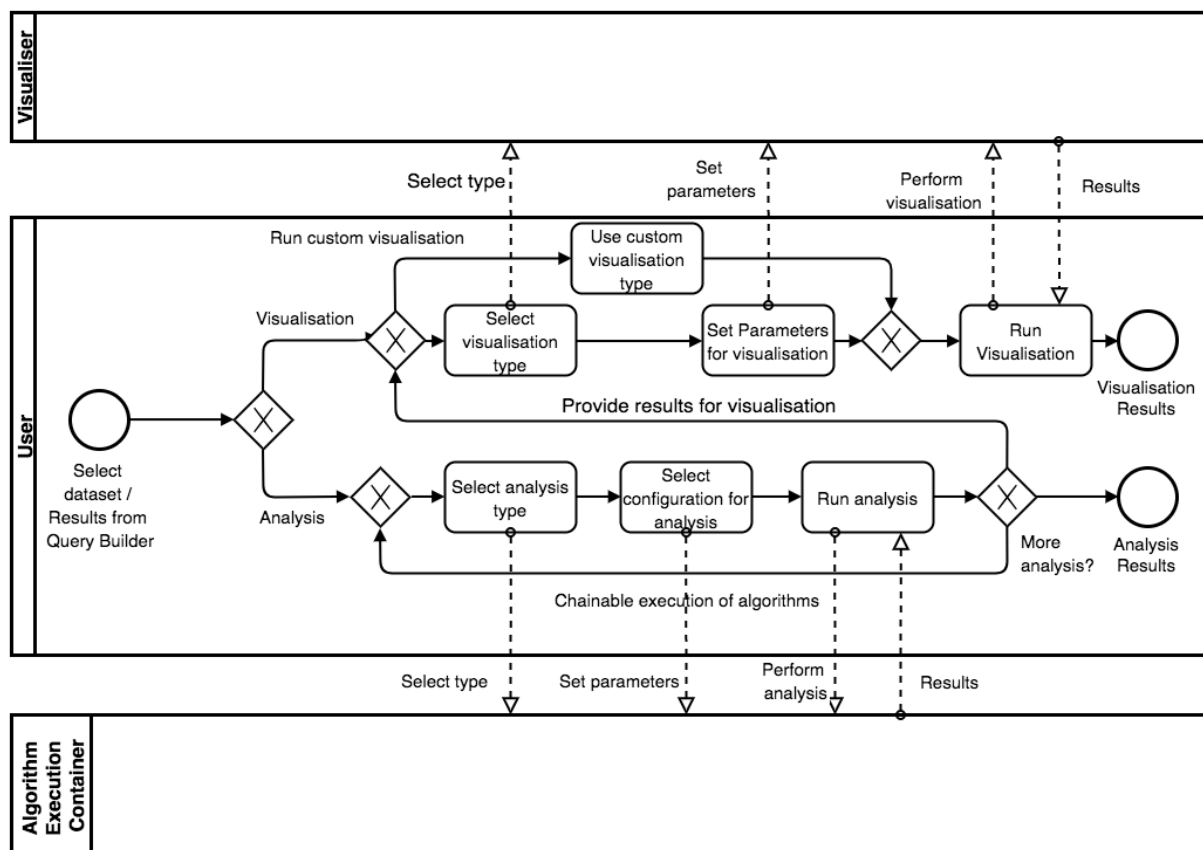
In case the above-mentioned check resolve that the data asset can be shared, then the transaction is inserted to the ledger, and the AEGIS platform is notified to release the data asset to the user.



**Figure 4-9: Artefact Sharing Workflow**

#### 4.6. Service creation

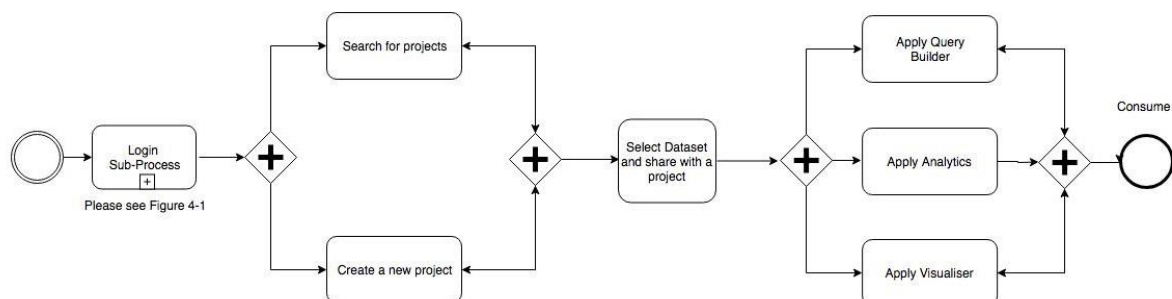
The following figure presents the workflow of the user’s perspective when he/she will use the AEGIS services in order to either perform an analysis or generate a visualisation on an available dataset. In particular, the user interacts with the Algorithm Execution Container and the Visualiser components. The user is offered the option to request visualisation from the list of available visualisations in the Visualiser or a custom visualisation on an available dataset or on a dataset created as a result of the Query Builder execution. Additionally, the user is offered the option to perform a new analysis, multi-level analysis by chainable execution of algorithms, and request visualisation on the analysis results.



**Figure 4-10: Service creation workflow**

#### 4.7. Service consumption

The following figure shows a workflow of the user's perspective when he/she will use the AEGIS platform to perform a general service consumption, which in this specific case includes the account creation/authentication, the search functionalities related to projects, latest assets and offers, the selection of a Dataset (together with its association with a Project) and the application to AEGIS main functionalities (Query Builder, Analytics, Visualiser).



**Figure 4-11: AEGIS Service consumption workflow**



## 5. CONCLUSION

The current deliverable documents the efforts undertaken within the context of the tasks 3.1, 3.2, 3.3, 3.4 and 3.5 of WP3. The scope of this deliverable was to provide final detailed documentation with regard to the high-level and technical architecture of the platform, the components of the platform, as well as the workflows of the platform. Hence, the deliverable built directly upon the information provided from the previous deliverable, namely the deliverable D3.4 in order to describe all the updates and enhancements that were introduced in the course of the development of the platform and provide the up-to-date complete documentation.

More specifically, in the current deliverable documented the final version of AEGIS platform high-level architecture which includes all the improvements and refinements, highlighting the role of each component within the platform along with the responsibilities and functionalities of each component. Additionally, the deliverable provided insights on the positioning of each component within the architecture. Besides the high-level architecture, the current deliverable presented the final version of the technical architecture of the platform with the focus being on the functional decomposition of components, the relationship between them and the designed data flow.

Following the platform's architecture, the current document provided the final detailed documentation of the design and functionalities of each component, as well as the details for the enhancements and refinements that were introduced in order to enable or enhance the offerings of the platform. Additionally, for each component the list of the designed microservices that support these functionalities is documented. The detailed documentation includes also the list of technologies that were utilised for the implementation of the component, as well as the implemented technical interfaces or exposed outcomes of each component that were used in the smooth integration of the various components in order to realise the designed workflows of the platform.

Within the context of this deliverable, the final AEGIS platform workflows in the form of BPMN diagrams were presented. The focus on the workflows description was on hiding the technical details and on highlighting the functionalities of the platform from the user's perspective.

The outcomes of this deliverable will drive the implementation activities of the project towards the implementation of the final release of the AEGIS platform, namely the AEGIS Platform Release 4.00, that will be delivered in M30 and will be used as guidance towards the strengthening of the AEGIS platform offerings.

**APPENDIX A: LITERATURE**

- [1] S. Niazi, S. Haridi and J. Dowling, “Size Matters: Improving the Performance of Small Files in HDFS,” in *EuroSys*, 2017.
- [2] P. Archer, S. Goedertier and N. Loutas, Study on persistent URIs, with identification of best practises and recommendations on the topic for the MSs and the EC, ISA Programme, 2012.
- [3] S. Kirrane, A. Mileo and S. Decker, “Access control and the resource description framework: A survey.,” in *Semantic Web*, 2017, pp. 311-352.