



Policy Cloud
Cloud for Data-Driven Policy Management

CLOUD FOR DATA-DRIVEN POLICY MANAGEMENT

Project Number: 870675

Start Date of Project: 01/01/2020

Duration: 36 months

D2.6 CONCEPTUAL MODEL & REFERENCE ARCHITECTURE

Dissemination Level	PU
Due Date of Deliverable	30/6/2021, M18
Actual Submission Date	30/06/2021
Work Package	WP2 Requirements, Architecture & Innovation
Task	T2.2 Definition of Target Conceptual Model & Reference Architecture
Type	Report
Approval Status	
Version	V2.0
Number of Pages	p.1 – p.70

Abstract: This document provides the first update to the Conceptual Model and Reference Architecture of PolicyCLOUD (the initial document has been submitted as Deliverable D2.2). This second version (Deliverable D2.6) provides the definition of the overall architecture and its components, focuses on the integration aspects of the architecture and presents example scenarios and associated data sources from the project's Use Cases including end-to-end data path analysis. The final update of the deliverable will be published in M30.

The information in this document reflects only the author's views and the European Community is not liable for any use that may be made of the information contained therein. The information in this document is provided "as is" without guarantee or warranty of any kind, express or implied, including but not limited to the fitness of the information for a particular purpose. The user thereof uses the information at his/ her sole risk and liability. This deliverable is licensed under a Creative Commons Attribution 4.0 International License.



PolicyCloud has received funding from the European Union's Horizon 2020 research and innovation programme under grant agreement No 870675.

Versioning and Contribution History

Version	Date	Reason	Author
1.1	4/09/2020	Deliverable D2.2 Conceptual Model & Reference Architecture (submitted)	Please refer to the submitted document
1.2	13/05/2021	Integration enabling cloud infrastructure (section 7.3.2)	Giuseppe La Rocca (EGI)
1.3	06/06/2021	<ul style="list-style-type: none"> Data Acquisition and Analytics Integration (section 7.6.11) Data Path Analysis (From Cloud Gateways to LXS Database) based on the implemented Use Case scenarios (section 8.3) Integration with the Kubernetes cluster (section 7.6.11.2) Integration of the Social Dynamics component (section 7.6.11.3) 	Ofer Biran (IBM), Oshrit Feder (IBM), Yosef Moatti (IBM), Nikitas Sgouros (UPRC)
1.3	06/06/2021	<ul style="list-style-type: none"> Ethical and Legal Compliance Framework (updated section 7.5.1) Ethical and Legal Compliance Framework Integration with Use Cases and technology (section 7.5.2) 	Alberto Bettiol (ICTLC), Martim Taborda Barata (ICTLC)
1.4	07/06/2021	<ul style="list-style-type: none"> Data Governance model, protection and privacy enforcement mechanisms Integration - Architecture Integration, Integration with the Kubernetes cluster (section 7.10.2) 	Konstantinos Oikonomou (UBI)
1.4	07/06/2021	<ul style="list-style-type: none"> Cloud Gateways (updated section 7.4.1) Data Marketplace (updated section 7.9.1) 	Thanos Kiourtis (UPRC), Argyro Mavrogiorgou (UPRC), George Manias (UPRC), Nikitas Marinos Sgouros (UPRC), Dimosthenis Kyriazis (UPRC)
1.5	03/06/2021	<ul style="list-style-type: none"> Policy Development Toolkit and Data Visualization (updated section 7.8.1) PDT Architecture (updated section 7.8.2) PDT Architecture Integration (section 7.8.3) 	Chris Maragkos (OKS), Kostas Moutselos (ICCS), Armend Duzha (MAG)

1.6	17/06/2021	<ul style="list-style-type: none"> Data Path Analysis (from the LXS database backend to visualization of result) (section 8.4) 	Jose Maria Zaragoza (LXS), Jacob Roldan (LXS), Patricio Martinez (LXS)
1.7	22/06/2021	<ul style="list-style-type: none"> Incentives Management (updated section 7.4.2) Incentives Management Architecture Integration (section 7.4.2.1) 	Maria Angeles Sanguino Gonzalez (ATOS), Ricard Munné (ATOS)
1.8	29/06/2021	<ul style="list-style-type: none"> Peer review 	Yosef Moatti (IBM), Nikos Achilleopoulos (MAG)
1.9	29/06/2021	<ul style="list-style-type: none"> Quality Check 	Argyro Mavrogiorgou (UPRC)
2.0	30/06/2021	<ul style="list-style-type: none"> Technical Coordination, Architecture Integration, Editing of document 	Panayiotis Tsanakas (ICCS), Panayiotis Michael (ICCS), Vrettos Moulos (ICCS)

Author List

Organisation	Name
ATOS	Maria Angeles Sanguino Gonzalez
ATOS	Jorge Montero Gomez
ATOS	Tomas Pariente Lobo
ATOS	Ricard Munné
EGI	Giuseppe La Rocca
IBM	Ofer Biran
IBM	Oshrit Feder
IBM	Yosef Moatti
OKS	Chris Maragkos
ICCS	Kostas Moutselos
ICCS	Vrettos Moulos
ICCS	Panayiotis Tsanakas
ICCS	Panayiotis Michael
ICTLC	A. Bettiol
ICTLC	M. Taborda Barata
ITA	Rafael del Hoyo
LON	Ebenezeer Williams
LON	Sarah Frost
LON	Adil Mohammed Ali
LXS	Jose Maria Zaragoza
LXS	Jacob Roldan
LXS	Patricio Martinez
LXS	Javier López Moratalla
LXS	Sadra Ebro
MAG	Armend Duzha
MAG	Nikos Achilleopoulos
OKS	Petya Bozhkova
OKS	Konstantinos Nasias
SARGA	Javier Sancho
SOF	Iskra Yovkova
UBI	Konstantinos Oikonomou
UBI	Giannis Ledakis
UPRC	Thanos Kiourtis
UPRC	Ilias Maglogiannis

UPRC	Argyro Mavrogiorgou
UPRC	George Manias
UPRC	Nikitas Marinos Sgouros
UPRC	Dimosthenis Kyriazis

Abbreviations and Acronyms

Abbreviation/Acronym	Definition
ABAC	Attribute-based access control
API	Application Programming Interface
CMF	Cloud Management Framework
CFREU	Charter of Fundamental Rights of the European Union
DB	Database
DPIA	Data Protection Impact Assessment
DSS	Decision Support System
EC	European Commission
ECHR	European Convention on Human Rights
EOSC	European Open Science Cloud
EU	European Union
GDPR	General Data Protection Regulation
GTD	Global Terrorism Database
IaaS	Infrastructure as a Service
JDBC	Java Database Connectivity
JSON	JavaScript Object Notation
KPI	Key Performance Indicators
ML	Machine Learning
NLP	Natural Language Processing
NoSQL	Non-Structured Query Language
OLAP	Online analytical processing
OLTP	Online transaction processing
PaaS	Platform as a Service
PDT	Policy Development Toolkit
PMO	Policy Model
PME	Policy Model Editor
PM	Policy Maker
PP	Public Policy
PR	Pattern Recognition

REST	Representational state transfer
SaaS	Software as a Service
SKA	Situational Knowledge Acquisition
SKM	Situational Knowledge Model
SOA	Service Oriented Architecture
SPA	Single Page Application
SQL	Structured Query Language
TRL	Technology Readiness Level
UI	User Interface
VM	Virtual Machine

Contents

Versioning and Contribution History.....	2
Author List.....	4
Abbreviations and Acronyms	6
1 Executive Summary.....	11
2 Introduction	12
2.1 Summary of Changes	14
3 Terminology.....	15
4 PolicyCLOUD offerings.....	16
5 PolicyCLOUD capabilities	17
6 PolicyCLOUD Conceptual Model	19
6.1 Conceptual Model	19
7 PolicyCLOUD Architecture	22
7.1 Architecture Building Blocks	22
7.2 Architecture Overview	25
7.3 Layer 1a - Cloud Based Environment.....	27
7.3.1 The EGI Federated Cloud.....	27
7.3.2 Integration enabling cloud infrastructure	28
7.4 Layer 1b - Data Management and Data Stores.....	28
7.4.1 Cloud Gateways.....	28
7.4.2 Incentives Management.....	30
7.4.3 Data Management and Data Stores	32
7.5 Ethical and Legal Compliance Framework	35
7.5.1 Ethical and Legal Compliance Framework.....	35
7.5.2 Ethical and Legal Compliance Framework Integration with Use Cases and technology.....	37
7.6 Layer 2 - Data Acquisition and Analytics	40
7.6.1 Data Acquisition and Analytics – Positioning & Goals.....	40
7.6.2 Extensibility and Reusability of Analytic Functions	41
7.6.3 Data Cleaning.....	42
7.6.4 Data Interoperability	42
7.6.5 Data Fusion with Processing and Initial Analytics.....	43

7.6.6	Seamless Analytics on Hybrid Data at Rest.....	44
7.6.7	Situational Knowledge Analysis	45
7.6.8	Opinion Mining.....	46
7.6.9	Sentiment Analysis	46
7.6.10	Social Dynamics	47
7.6.11	Data Acquisition and Analytics Integration	47
7.7	Layer 3 – Policy Management Framework.....	50
7.7.1	Policy Modelling & KPIs Identification	50
7.7.2	Middleware for Policies	50
7.8	Layer 4 - Policy Development Toolkit.....	51
7.8.1	Policy Development Toolkit and Data Visualization	51
7.8.2	PDT Architecture	51
7.8.3	PDT Architecture Integration.....	53
7.9	Layer 5	54
7.9.1	Data Marketplace.....	54
7.10	Data Governance Model, Protection and Privacy Enforcement.....	56
7.10.1	Data Governance Model, Protection and Privacy Enforcement	56
7.10.2	Data Governance model, protection and privacy enforcement mechanisms Integration	57
8	Use Case examples for end-to-end data path analysis.....	59
8.1	Use Case 1: Participatory Policies Against Radicalization.....	60
8.1.1	Scenario A: Radicalization incidents	60
8.1.2	Main Objective.....	60
8.1.3	Key Performance Indicators	61
8.1.4	Data Sources.....	61
8.2	Use Case 2: Intelligent policies for the development of the agrifood industry	62
8.2.1	Scenario B: Visualization of negative and positive opinions on social networks for different products.....	62
8.2.2	Main Objective.....	62
8.2.3	Key Performance Indicators	62
8.2.4	Data Sources.....	64
8.3	Data Path Analysis (From Cloud Gateways to LXS Database) based on the implemented Use Case scenarios.....	65

8.4	Data Path Analysis (from the LXS database backend to visualization of result).....	66
9	Conclusion.....	69
	References.....	70

List of Tables

Table 1 – UC1 business KPI9.....	61
Table 2 – Data Sources list for Scenario A of the Participatory Policies Against Radicalization Use Case	61
Table 3 – UC2 technical KPI4	62
Table 4 – UC2 technical KPI6	63
Table 5 – UC2 business KPI8.....	63
Table 6 – UC2 business KPI9.....	63
Table 7 – UC2 business KPI10	64
Table 8 – UC2 business KPI11	64
Table 9 – Links to Aragon use case data stores	64

List of Figures

Figure 1 – The PolicyCLOUD Conceptual Model.....	19
Figure 2 – PolicyCLOUD Architecture Building blocks.....	22
Figure 3 – PolicyCLOUD Architecture Implementation over the European Cloud Initiative infrastructure offered by EGI	24
Figure 4 – PolicyCLOUD overall Architecture	25
Figure 5 – Incentives Management architecture integration	31
Figure 6 – Part of the PolicyCLOUD Overall Architecture Diagram relevant to WP4.....	40
Figure 7 – WP4 interface with WP3 and WP5.....	41
Figure 8 – The streaming data path.....	43
Figure 9 – Seamless analytics on ingested data	45
Figure 10 - WP4 Constituent Architecture	48
Figure 11 – Policy Development Toolkit Communication Components.....	52
Figure 12 – Policy Model development Integration Schema	53
Figure 13 – Data Marketplace architecture	55
Figure 14 – Data Governance model, protection and privacy enforcement mechanisms – Extracted views (a), (b) and (c) from the diagram of PolicyCLOUD Overall Architecture.	57
Figure 15 - Data Governance & Privacy Enforcement Mechanism integration flow	58
Figure 16 - Visualization on PolicyCLOUD of the result of Scenario A: Radicalization incidents of Use Case 1	60
Figure 17 - Data Path Analysis.....	65
Figure 18 - Sequence diagram for PDT-DAA interaction	67

1 Executive Summary

A second version of the PolicyCLOUD Conceptual Model & Reference Architecture (originally submitted as Deliverable D2.2) is presented in this document. The final update of the deliverable will be submitted in M30.

The PolicyCLOUD Conceptual Model presents the overall project concept along 2 main axes. Along the first data axis PolicyCLOUD delivers Cloud Gateways and APIs to access data sources and adapt to their interfaces so as to simplify interaction and data collection from any source. Along the second main axis, the Policies Management Framework of PolicyCLOUD allows the definition of forward-looking policies as well as their dynamic adaptation and refocusing to the population they are applied on.

Based on the project's offerings along the main two axes of the Concept, five main building blocks (in a layered manner) define its Architecture: (1) The Cloud Based Environment and Data Acquisition, (2) Data Analytics, (3) the Policies Management Framework, (4) the Policy Development Toolkit and (5) The Marketplace.

The architecture also includes a Data Governance Model, Protection and Privacy Enforcement and the Ethical Framework as depicted in figure 2.

The architecture allows for integrated data acquisition and analytics. It also allows data fusion with processing and initial analytics (see 7.6.5) as well as seamless analytics (see 7.6.6) on hybrid data at rest.

The main focus of the updates of the document is on the Integration in PolicyCLOUD which follows three directions: (i) architecture integration, (ii) integration with the cloud infrastructure and (iii) integration with Use Case scenarios through the implementation of end-to-end scenarios. Additional integration activities take place along the two frameworks of PolicyCLOUD, (a) the Data Governance model, protection and privacy enforcement mechanism and (b) the Ethical and Legal Compliance framework.

For end-to-end data path analysis we have used two Use Case scenarios: (i) the scenario of Use Case 1: "Radicalization incidents" and the scenario of Use Case 2: "Visualization of negative and positive opinions on social networks for different products".

2 Introduction

The definition of the Conceptual Model and Reference Architecture is a continuous, dynamically changing task, following the development of the project from M1 to M30. This document is the first update of the Conceptual Model and Reference Architecture of PolicyCLOUD. The initial document has been submitted as Deliverable D2.2 and this second version is Deliverable D2.6. The final update will be submitted in M30.

The document is structured as follows: The PolicyCLOUD Conceptual Model explaining the overall project concept through 2 main axes is presented in Section 6, while the PolicyCLOUD Architecture consisting of five main building blocks (five Layers) that realize the project's offerings along the main two axes of the Concept, is presented in Section 7.

More specifically an overview of the overall architecture as presented and discussed (i) during the Kick-Off meeting, (ii) during the development of the preliminary specification as an internal report made available to partners and (iii) during specialized workshops integrating constituent architectures, is presented in section 7.2. In sections 7.3-7.9 the five layers of the architecture are presented as follows:

- **Layer 1a-Cloud Based Environment** is presented in Section 7.3.
- **Layer 1b-Data Management – Data Stores** is presented in Section 7.4.
- **Layer 2-Data Acquisition and Analytics** is presented in section 7.6.
- **Layer 3-Policy Management Framework** is presented in section 7.7.
- **Layer 4-Policy Development Toolkit and Visualization** is presented in section 7.8.
- **Layer 5-Data Marketplace** is presented in Section 7.9.

The **Ethical and Legal Compliance Framework** presented in Section 7.5 is included in the architecture from the very beginning of the project in order to provide extensive and in-depth analysis of relevant legal, regulatory, societal and ethical aspects.

The **Data Governance Model, Protection and Privacy Enforcement** used to protect data and ensure decisions across the complete path that follow specific guidelines and legislations, is presented in Section 7.10.

The architecture allows for integrated data acquisition and analytics. It also allows data fusion with processing and initial analytics (see 7.6.5) as well as seamless analytics (see 7.6.6) on hybrid data at rest.

For end-to-end data path analysis and in order to demonstrate the characteristics of the integrated architecture we have used two Use Case scenarios (section 8): The scenario of Use Case 1: "Radicalization incidents" (8.1.1) and the scenario of Use Case 2: "Visualization of negative and positive opinions on social networks for different products" (8.2.1).

The main focus of the updates of the document is on the Integration in PolicyCLOUD which follows three directions: (i) architecture integration, (ii) integration with the cloud infrastructure and (iii) integration with Use Case scenarios through the implementation of end-to-end scenarios.

Additional integration activities take place along the two frameworks of PolicyCLOUD, (a) the Data Governance model, protection and privacy enforcement mechanism and (b) the Ethical and Legal Compliance framework. More specifically, in this document we report integration activities that take place for the integration of the mechanisms provided by the Data Governance model, protection and privacy enforcement and (b) for the implementation of the recommendations received from the Ethical and Legal Compliance framework in an integrated manner to the PolicyCLOUD technology and to the end-to-end Use Case scenarios.

The new sections of the updated document (also summarized in section 2.1) address the following integration subjects:

1. Integration enabling cloud infrastructure (section 7.3.2) addresses the integration of the PolicyCLOUD components on the Kubernetes cluster provisioned by RECAS-BARI, through EGI.
2. Data Acquisition and Analytics Integration (section 7.6.11) and Integration with the Kubernetes Cluster (section 7.6.11.2) presents the Work Package 4 related constituent architecture implemented on a Kubernetes cluster demonstrating the integration of OpenWhisk, Kafka, Leanxcale datastore and a Cloud Gateway.
3. Various alternatives on how to integrate the Social Dynamics component (section 7.6.11.3) are examined.
4. The incentives management tool uses as integration point the PDT interface where the component frontend is integrated as an additional entry point for the policy makers (section 7.4.2.1).
5. A section referring to the Integration of the Ethical and Legal Compliance Framework with the Use Case scenarios and technology has been authored (section 7.5.2).
6. The integration of Data Governance model, protection and privacy enforcement mechanisms is presented (section 7.10.2) as also the integration of these mechanisms with the Kubernetes cluster (section 7.10.2.2).
7. A section presenting the PDT architecture integration has been authored (section 7.8.3).
8. Data Path Analysis based on the implemented use case scenarios (section 8) explains the data path followed for visualizing a heatmap showing the frequency of occurrence of radicalization incidents for the scenario A of Use Case 1. The data path analysis consists of two subpaths: (i) the subpath from the Cloud Gateways to LXS database (section 8.3) and (ii) the subpath from the LXS database backend to the visualization of result (section 8.4).

2.1 Summary of Changes

A summary of changes is provided in the following list:

1. Integration enabling cloud infrastructure (section 7.3.2).
2. Data Acquisition and Analytics Integration (section 7.6.11) including Integration with the Kubernetes Cluster (section 7.6.11.2).
3. The process of integrating the Social Dynamics component (section 7.6.11.3).
4. The architecture integration of the incentives' management tools (section 7.4.2.1).
5. A special analysis on the Integration of the Ethical and Legal Compliance Framework with the Use Case scenarios and technology (section 7.5.2).
6. The integration of Data Governance model, protection and privacy enforcement mechanisms (section 7.10.2) including the integration of these mechanisms with the Kubernetes cluster (section 7.10.2.2).
7. PDT architecture integration (section 7.8.3).
8. Data Path Analysis based on the implemented use case scenarios (section 8) which includes: (i) the subpath from the Cloud Gateways to LXS database (section 8.3) and (ii) the subpath from the LXS database backend to the visualization of result (section 8.4).

3 Terminology

Policies KPIs are the key performance indicators (i.e. metrics/parameters) included in the structural representation of policies. These indicators are used to model the policies as well as to monitor and evaluate them.

Platform as a Service Orchestrator allows to coordinate the provisioning of virtualized compute and storage resources on Cloud Management Frameworks, both private and public (like OpenStack, OpenNebula, AWS, etc.) and the deployment of dockerized long-running services and batch jobs on Apache Mesos clusters [1].

PDT (Policy Development Toolkit) is a framework which incorporates the visualization workbench and provides a unique point of interaction with the policy makers. Through the toolkit the policy makers are able to state their questions and perform policy modelling and policy making.

Object Storage is designed to support exponential data growth and cloud-native workloads. It provides cross-region offerings, and integrated services. Depending on the access frequency of the data, storage can be provided in three “**smart tiers**”: **Hot, Cool and Cold** [2].

4 PolicyCLOUD offerings

PolicyCLOUD offerings are materialized through five main building blocks, supported by the **Ethical and Legal Compliance Framework** and the **Data Governance Model, Protection and Privacy Enforcement Framework**.

In summary these offerings are the following:

1. **The Cloud Capabilities & Data Collection Engine** that incorporates technologies for interfacing and acquiring data from various sources.
2. **The Reusable Models & Analytical Tools Engine** that incorporates all data services / technologies provided by PolicyCLOUD for the data path/lifecycle.
3. **The Policies Management Framework**.
4. **The Policy Development Toolkit** providing an interactive environment and the Front-End of the system.
5. **The Data Marketplace** which enables data and knowledge to be exploited as assets while keeping conformance with legal and ethical requirements and privacy protection.

Finally, the **Ethical and Legal Compliance Framework** assures that all PolicyCLOUD offerings conform to the required ethical, legal and security aspects while the **Data Governance Model, Protection and Privacy Enforcement Framework** protects data and ensures decisions across the complete path following specific guidelines and legislations.

The details of the PolicyCLOUD offerings listed above are provided in section 7.1 with title Architecture Building Blocks.

5 PolicyCLOUD capabilities

PolicyCLOUD provides an innovative suite of state-of-the-art technology capabilities and management frameworks over a Cloud environment as presented in the following list:

- **Cloud Based environment** to support the development of PolicyCLOUD using Platform as a Service (PaaS) solutions.
- **Unified Cloud Gateway** moving streaming and batch data from data owners into PolicyCLOUD layers while performing data source reliability.
- **Incentives identification and management** offering a set of tools to identify and manage incentives able to engage different participants on the policy making process.
- **Access to heterogeneous data stores.**
- **Scalable Database** with the ability to scale out over hundreds of nodes.
- **Polyglot capabilities** enabling the Querying of Heterogeneous Data Sources in a Unified manner.
- **Ability to combine analytics on streaming data and on data at rest.**
- **Transparent to the user movement of colder data** to the Object Store tier.
- **Data Cleaning** for the detection and correction of corrupted or inaccurate records received from Cloud Gateways.
- **Data Interoperability** based on data-driven design, coupled with linked data technologies, in order to improve both semantic and syntactic data and dataset interoperability.
- **A business process for clearing private data, as well as "open data"** evaluating if and to which extent personal data in terms of the GDPR is allowed to be processed by PolicyCLOUD.
- **Data Fusion** task permitting the merging of data coming from disparate sources into a single data set, integrated with initial analytics and data processing tasks.
- **Seamless Analytics** permitting undifferentiated access and query capabilities both to hot (in the DB) and cold (in the object storage) data.
- **Situational knowledge** from data from sensors, social media and datasets offering feature extraction, clustering and categorization.
- **Opinion Mining** providing social attitude regarding specific topics, identifying specific entities and generating a "contributor graph" based on discussions of various policies from citizens.
- **Sentiment analysis** based on the input received from the pilots about their policies.
- **Social Dynamics** providing a concurrent, web-based environment for social simulation. The environment allows users to create graph-based population models online.

- **Framework for Cloud usage by Public Authorities** examining (a) the different mechanisms, methods and technologies used for policy lifecycle and (b) a proposition of a set of adaptable techniques towards the utilization of cloud environments for policies creation.
- **Middleware for modelling and designing of Policies** providing a mechanism for policies to be modelled and designed based on specific structural representations, allowing users to create a policy by selecting a schema of data, applying well known Key Performance Indicators.
- **Policy Development Toolkit (PDT)** constituting the Front-End of PolicyCLOUD environment. It integrates several sub-components to enable policy makers to create, update and validate policy models.
- **Integrated cloud-based framework** designed for the Cloud, structured over five layers including an Ethical and Legal Compliance Framework and a Data Governance Model providing all the above capabilities.

6 PolicyCLOUD Conceptual Model

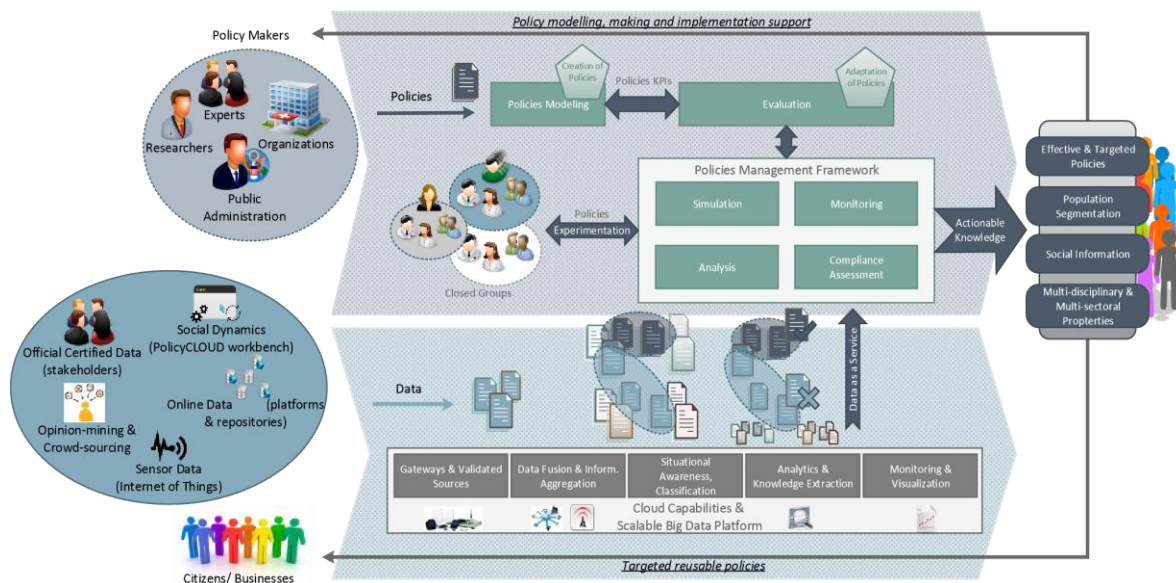


FIGURE 1 – THE POLICYCLOUD CONCEPTUAL MODEL

6.1 Conceptual Model

PolicyCLOUD architecture delivers a set of innovative technologies with an overall goal to enable data-driven management of policies lifecycle, from their modelling and implementation, to optimization, compliance monitoring, adaptation and enforcement.

As depicted Figure 1, PolicyCLOUD architecture enables the compilation of multi-disciplinary, and multi-sectoral optimized policies. Multi-disciplinary policies aim at addressing different spatiotemporal levels. In terms of time scales, different policies are proposed to be applied in long-term, while these policies could address a specific area (e.g. city), a region, or even a country. The combination of these properties of policies are optimized through PolicyCLOUD according to the modelling and evaluation of different policies and their corresponding KPIs.

Additionally, data emerging from policies “collections” / clusters (e.g. all policies in a city, environmental policies in different cities, health policies for specific age groups, etc) provide additional information for the optimization of policies in the aforementioned scale. Furthermore, PolicyCLOUD architecture enable multi-sectoral optimization of policies.

As shown in Figure 1, policies effectiveness is assessed and optimized based on their KPIs (vertical optimization) while KPIs of policies from other sectors are also taken into consideration (horizontal, cross-sector optimization). To realize the overall multi-sectoral effectiveness of policies, PolicyCLOUD architecture includes technologies for correlation of the policies and the data used to compile these policies through reusable and scalable models and analytic tools.

The architecture serves the overall project concept of PolicyCLOUD and it is realized through 2 main axes: the data axis and the policies axis (Figure 1).

Along the first data axis PolicyCLOUD delivers Cloud Gateways and APIs to model the data sources and adapt to their interfaces so as to simplify interaction and data collection from any source.

Some of these sources may not provide reliable information and thus before taking it into consideration, gateways are enhanced with the functionality of validating the data in order to develop trust and reliability profiles and patterns of the sources and exploit only the reliable ones.

In terms of data sources, PolicyCLOUD obtains open data from the ecosystem stakeholders (e.g. public authorities), sensor data from Internet of Things infrastructures (e.g. environmental sensors), data from online platforms, opinion-mining and crowd-sourcing data (both from online platforms and from the proposed PolicyCLOUD living lab approach), as well as data related to social dynamics and behaviour through the corresponding analytical tools.

The ethical and legal compliance framework included in the architecture enables a process we name “data clearance” which examines available open-data for privacy issues (even if some data are characterized as “open” they may include private data). Data clearance processing combines legal expertise with technology (e.g. access control at critical points) in order to safeguard that data are efficiently used in a legal and ethical manner.

Based on the above, data fusion and information aggregation enable the compilation of information into new data and metadata structures which are interlinked and analyzed. This information along with existing policies provide a network of knowledge which is dynamically exploited for improving the effectiveness of existing policies and facilitating the creation and adoption of new policies.

PolicyCLOUD architecture delivers mechanisms for clustering, classification and situational awareness on big datasets and the corresponding policies. Core element in this process is the delivery of a powerful Reusable Models & Analytical Tools offering for cleaning datasets, modelling and representing them, as well as harnessing information and enabling knowledge extraction. This is performed by taking into consideration data and existing policies that correspond to target groups / public authorities with specific goals and population characteristics.

Given the wealth of information and the different administrative and legal domains under which data will be governed and managed, PolicyCLOUD includes a data governance model (based on RACI) that governs the complete data lifecycle (e.g. who has access, to which data, etc).

Along the second main axis the Policies Management Framework of PolicyCLOUD is exploited for the definition of forward-looking policies which are dynamically adapted and methodically focused on the population that are applied on.

Initially the policies are modelled in order to extract quantitative and qualitative information from them, such as KPIs, operational and functional dependencies for analysis and evaluation.

The architectural framework employs the knowledge incorporated into the clusters of data and policies for a) assessing and stratifying the risks of policies, b) monitoring and assessing their compliance and c) forecasting the effectiveness of policies, including variations and combinations of policies.

The process is supported both by simulation methodologies and techniques, as well as by analysing the results of applying the policies to closed groups – i.e. evidence-based. Evaluation is not based on policy-level but on KPI-level per policy and across sectors (addressing different verticals including environment, migration, employment, etc.). In addition, through the mechanisms described in the architecture, the policies strengths and weaknesses are identified and analyzed while when it comes to policies adoption, their effectiveness on different conditions, populations, methodologies etc. is effectively assessed.

Therefore, the policies not only are evaluated, but they are also differentiated with different parameter sets, applicable to certain groups, locations and conditions, with in advance knowledge of the risk and performance trade-offs. Identification of the exact elements of policies that can affect their outcomes, across all policies, will also enable the creation of policies taking advantage of the excellence of the particular elements on better and more targeted results, minimizing in parallel the uncertainty when integrating them in the public policy strategy.

The outcomes - as actionable knowledge - are delivered to policy makers as evidence-based targeted strategies for policy making (including the most relevant population segmentation and evidences to maximize the policies efficiency).

7 PolicyCLOUD Architecture

7.1 Architecture Building Blocks

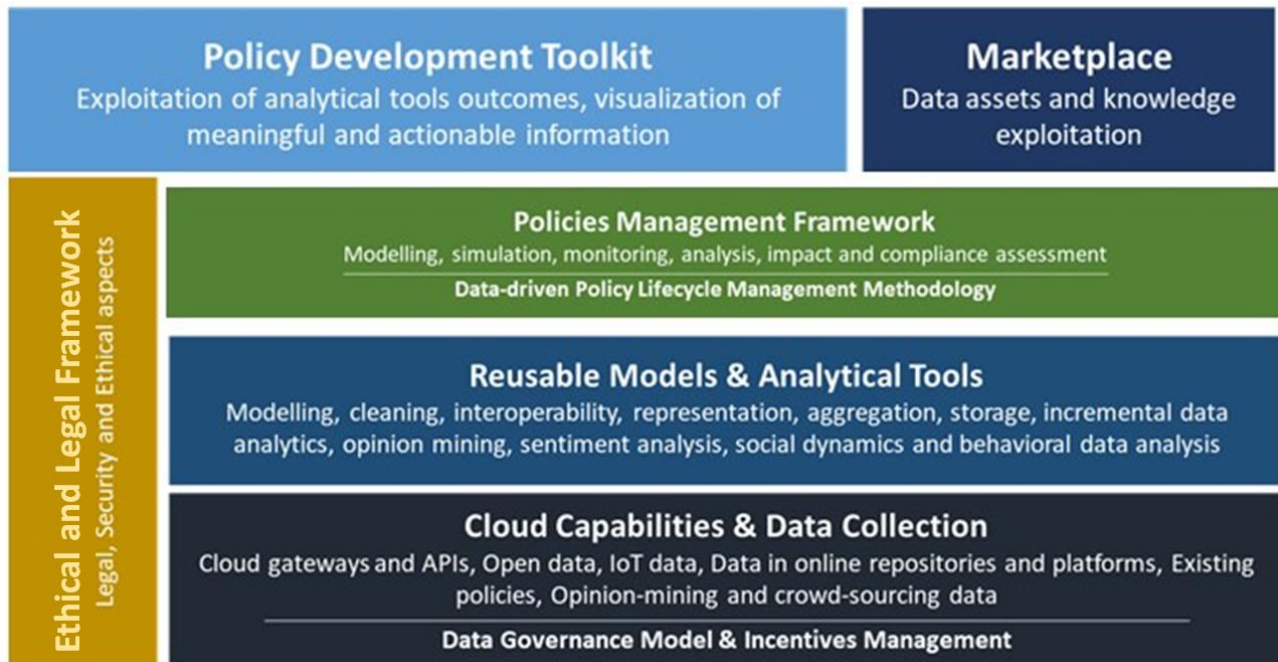


FIGURE 2 – POLICYCLOUD ARCHITECTURE BUILDING BLOCKS

The architecture of PolicyCLOUD includes five main building blocks that realize the project's offerings (Figure 2) along the main two axes of the Concept described in the previous section. These building blocks are presented in the following paragraphs in a bottom-up manner:

1. The **first building block** of the PolicyCLOUD architecture is the **Cloud Capabilities & Data Collection Engine** block that incorporates technologies for interfacing and acquiring data from different sources (through unified cloud gateways and APIs), assessing their reliability and attaching the corresponding metadata to the sources and ensuring privacy enforcement for the collected data, using the developed cloud infrastructure management. This block also includes mechanisms for identifying attributes of data and stakeholders in order to ensure that all data decisions are according to the data governance rules specified by the data owners, while it integrates techniques for managing the incentives in order to ensure citizens participation.

2. The **second building block** of the architecture is the **Reusable Models & Analytical Tools Engine** that incorporates all data services / technologies provided by PolicyCLOUD for the data path / lifecycle: modelling, cleaning, interoperability, linking / aggregation, storage and incremental analytics, for constructing the required reusable models. Moreover, this engine will also offer techniques for sentiment analysis from different online platforms, and tools for opinion-mining allowing stakeholders to “develop” through the provided toolkit, in an automated way, different means (such as aspect ranking) in order to acquire and analyse the corresponding information from citizens.
3. The **third building block** refers to the **Policies Management Framework** that incorporates services for the identification of the required KPIs in order to model the policies and identify potential interdependencies with other policies within and across sectors at different levels (e.g. local, national, etc). The framework also includes tools for collecting evidence monitoring information both from the engaged citizens and from the population targeted by the policies, while also assessing the compliance to these policies and thus assessing the policies impact (based on the identified KPIs).
4. The **fourth building block** (the interactive environment) provides the **Policy Development Toolkit** allowing policy makers to interact with the models and analytical tools as well as to specify their requirements and constraints with respect to different policies (e.g. specification of the need for policies that can have a real-time impact due to emergencies). In addition, the toolkit facilitates visualization of policies monitoring in an adaptive and incremental way.
5. The **fifth building block** of the architecture is the **Data Marketplace** which enables data and knowledge to be exploited as assets. Data Marketplace has two goals: (a) the usage of data in different contexts (scenarios for policy making) and (b) the identification of market opportunities.

The **Ethical and Legal Compliance Framework** assures that all the PolicyCLOUD offerings conform to the required ethical, legal and security aspects, thus ensuring the sustainability of the modelled policies. The framework is **vertically** depicted in the figure given that it obtains information from the **Cloud Capabilities & Data Collection Engine** (such as social networks data), while it communicates in a **bi-directional way** with the **Interacting Environment** by obtaining data from the **Policy Development Toolkit** and the **data marketplace**, and by specifying analytics tasks through this toolkit.

The architecture building blocks will be implemented over the European Cloud Initiative infrastructure offered by EGI (Figure 3).

The PolicyCLOUD Marketplace is part of the infrastructure and offers the solutions in terms of models and analytical tools that can be exploited by the end-users (i.e. policy makers and public authorities) through the PolicyCLOUD Policy Development Toolkit.

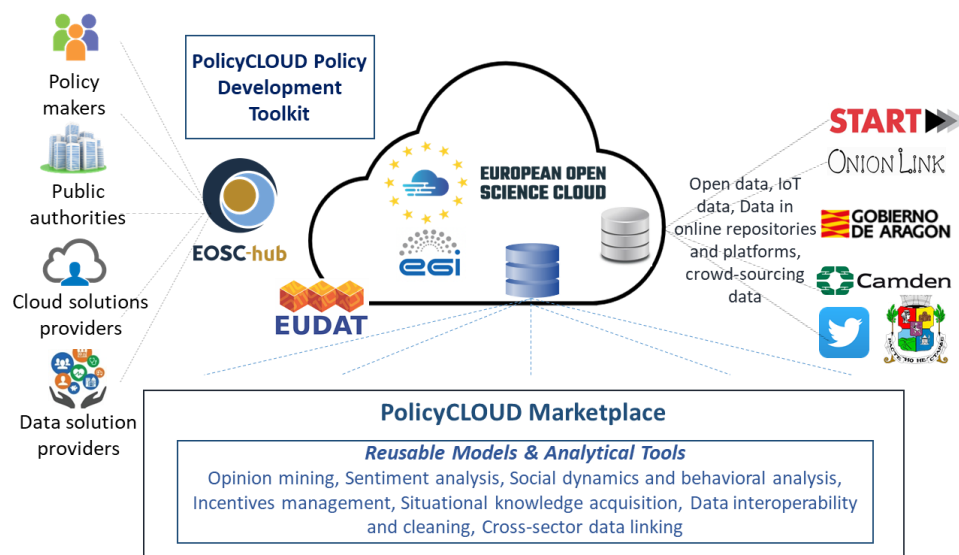


FIGURE 3 – POLICYCLOUD ARCHITECTURE IMPLEMENTATION OVER THE EUROPEAN CLOUD INITIATIVE INFRASTRUCTURE OFFERED BY EGI

7.2 Architecture Overview

The Overall Architecture (Figure 4) has been discussed and further developed (i) during the Kick-Off meeting, (ii) during the development of the preliminary specification as an internal report made available to partners and (iii) during specialized workshops integrating constituent architectures. The final version of the architecture will be published in M30. The architecture's layers and frameworks will be analyzed in the sections that follow.

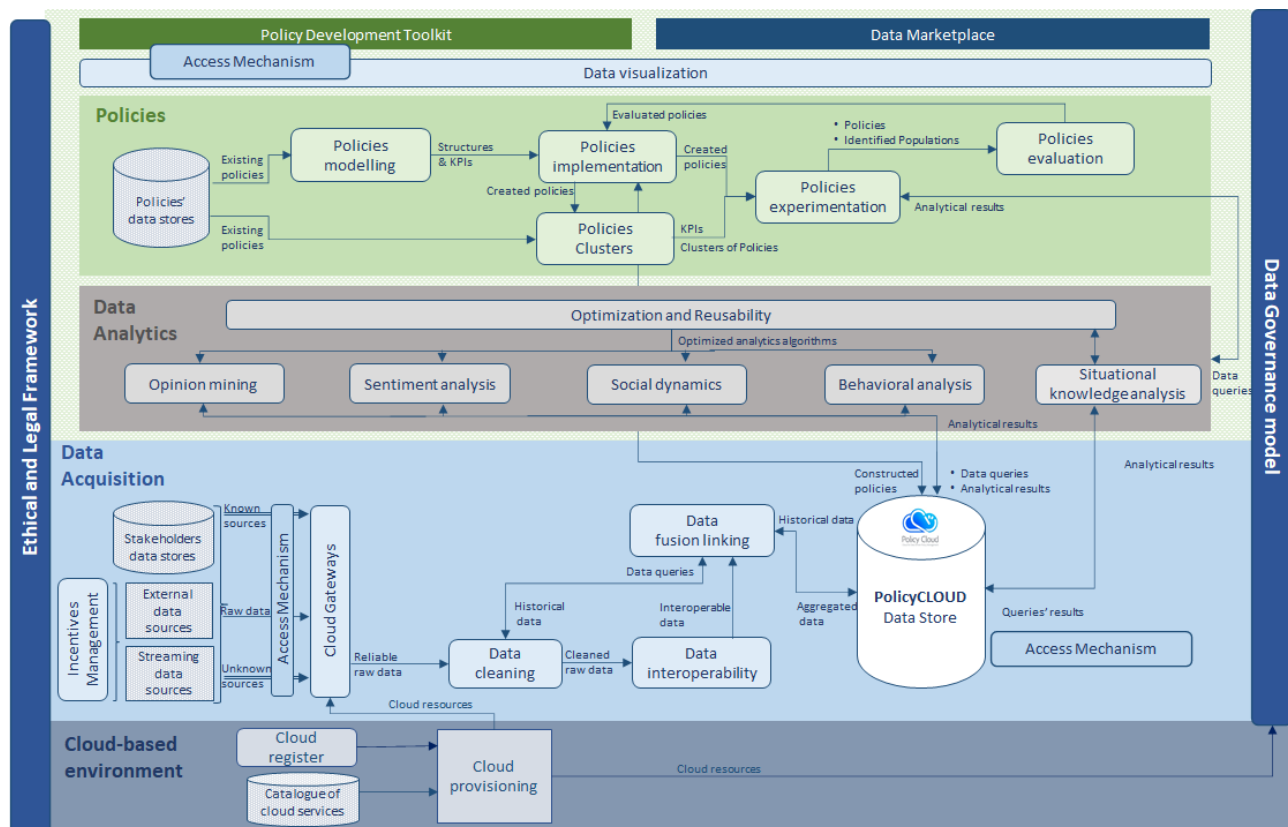


FIGURE 4 – POLICYCLOUD OVERALL ARCHITECTURE

As a complete environment, the proposed architectural approach is presented in Figure 4. The overall flow is initiated from various data sources, as depicted in the figure through the respective *Data Acquisition* block. Data sources can be data stores from public authorities or external data sources (e.g. mobile devices, IoT sensors, etc.) that contribute data following the provision of incentives, facilitated through the *incentives management* mechanism.

A set of APIs incorporated in a gateway component, enable data collection by applying techniques to identify the reliable sources and for these sources obtain the data and perform the required *data quality assessment and cleaning*. *Semantic and syntactic interoperability* techniques are utilized over the cleaned data providing the respective interoperable datasets to the PolicyCLOUD datastore following the required *data linking and aggregation* processes.

The datastore is accessible from a set of machine learning models represented through the *Data Analytics* building block. Machine learning models may incorporate opinion mining, sentiment and social dynamic analysis, behavioural analysis and situational / context knowledge acquisition. The data store and the analytics models are hosted and executed in a *cloud-based environment*. For this purpose, a catalogue in which an extensible set of services are registered has been implemented. Furthermore, all the analytics models are realized as services, thus enabling their invocation through a proposed policy development toolkit – realized in the scope of the *Policies* building block of the proposed architecture as a single point of entry into the PolicyCLOUD platform.

The toolkit allows the compilation of *policies as data models*, i.e. structural representations that include key performance indicators (KPIs) as a means to set specific parameters (and their target values) and monitor the implementation and effectiveness of policies against these KPIs along with the list of analytical tools to be used for their computation. According to these analytics outcomes, the values of the KPIs are specified resulting to *policies implementation / creation*. It should be noted that PolicyCLOUD also introduces the concept of *policies clusters* in order to interlink different policies, and identify the KPIs and parameters that can be optimized in such policy collections.

Across the complete environment, an implemented *data governance and compliance model* is enforced, ranging from the provision of cloud resources regarding the storage and analysis of data to the management of policies across their lifecycle.

7.3 Layer 1a - Cloud Based Environment

7.3.1 The EGI Federated Cloud

The EGI Federated Cloud is an IaaS-type cloud, made of academic private clouds and virtualized resources and built around open standards. Its development is driven by requirements of the scientific community. The Federation pools services from a heterogeneous set of cloud providers using a single authentication and authorisation framework that allows the portability of workloads across multiple providers and enables bringing computing to data. The current implementation is focused on IaaS services but can be easily applied to PaaS and SaaS layers. The EGI Federated Cloud architecture is based on the concept of an abstract Cloud Management Framework (CMF) that supports a set of cloud interfaces to communities.

Each resource centre of the infrastructure operates an instance of this CMF according to its own technology preferences and integrates it with the federation by interacting with EGI core components:

- Service registry for configuration management of federated cloud services.
- EGI AAI for authentication and authorisation across the whole cloud federation.
- Accounting for collecting, and displaying usage information.
- Information discovery about capabilities and services available in the federation.
- Virtual Machine image catalogue and distribution, replicating VM images as needed by the user communities in a secure way.
- Monitoring, performing service availability monitoring and reporting of the distributed cloud service end-points.

Users of the EGI Federated Cloud infrastructure can interact with cloud providers in several ways:

- Directly using the IaaS APIs of the resource centres to manage individual resources.
- Leveraging federated IaaS provisioning tools that allow managing and combining resources from different providers enabling the portability of application deployments between them. The EGI Federated Cloud task force is currently in the process of evaluating and selecting the best tools for this task.
- Using PaaS solutions such as the Infrastructure Manager (IM)¹, a Federated IaaS Provisioning tool, or the PaaS orchestrator developed within INDIGO-DataCloud².

In the context of the PolicyCLOUD project, EGI contributes to the provisioning of the needed computing resources to set-up the PolicyCLOUD infrastructure. This cloud infrastructure will help policy makers, public authorities and different stakeholders, to analyse a plethora of datasets from different data sources, and facilitate policy making. EGI offering for the project includes a federated IaaS cloud to run compute - or data -intensive tasks and host online services in virtual machines or Docker containers on

¹ See: <https://www.grycap.upv.es/im/index.php>

² See: <https://www.indigo-datacloud.eu/>

IT resources accessible via a uniform interface. More details about the federated EGI Cloud infrastructure and the solutions offered to address the needs of the project will be highlighted in D3.1 - Cloud Infrastructure Incentives Management and Data Governance: Design and Open Specification 1.

7.3.2 Integration enabling cloud infrastructure

Considering that most of the components of the PolicyCLOUD infrastructure are dockerized and distributed as Docker containers, to facilitate the provisioning of compute and storage resources, and the orchestration of distributed Kubernetes clusters, the access to the cloud infrastructure is also enabled via the PaaS orchestrator developed in the context of the INDIGO-DataCloud project³. For more details about the main architecture of the INDIGO-Data Cloud PaaS Orchestrator⁴, please refer to deliverable D3.2 “Cloud Infrastructure Incentives Management and Data Governance Software Prototype 1”.

A trial phase was planned in August 2020, during which the 10-20% of the full capacity of the PolicyCLOUD infrastructure was configured to allow technical partners to run tests and assess its performance. More specifically, during this phase, LXS, IBM and ICCS managed to test the deployment of a GitLab instance for the project and the deployment of OpenWhisk on Kubernetes. In October 2020 the project signed a pay-for-use agreement with the cloud resource provider (RECAS-BARI) and the full needed capacity was allocated to the project.

7.4 Layer 1b - Data Management and Data Stores

Components: Cloud Gateways (T3.3), Incentives Management (T3.4), Data Store (Figure 4).

7.4.1 Cloud Gateways

In the context of PolicyCLOUD, the Cloud Gateway and API component developed by UPRC seeks to enhance the abilities and services offered by a unified Gateway to move streaming and batch data from data owners into PolicyCLOUD data stores layers, which support both SQL and NoSQL data stores and public and private data. Based on the specifications provided in D3.1 Cloud Infrastructure Incentives Management and Data Governance Design and Open Specification 1 [3] of the PolicyCLOUD project, the effort related to Cloud Gateways & APIs component focuses on providing a complete and “smart” entryway into PolicyCLOUD project, allowing multiple APIs or microservices to act cohesively and thus provide a uniform, gratifying experience to each stakeholder. The provided Gateway API allows building scalable and robust APIs, while simplifying the interaction and data collection from various sources and providers. The main goal of this component is to handle a request by invoking multiple microservices

³ <https://cordis.europa.eu/project/id/653549>

⁴ <https://indigo-paas.cloud.ba.infn.it/home/login>

and aggregating the results. Hence, it enhances the design of resources and structure, add dynamic routing parameters and develop custom authorizations logic. PolicyCLOUD's Cloud Gateway and API component supports scalability, high availability and shared state without compromising performance. Moreover, it supports client side load balancing, so that the overall system can apply complex balancing strategies and do caching, batching, fault tolerance, service discovery and handle multiple protocols. To this end, MoleculerJS⁵, a framework that bases its functionality on microservices architecture methodology, is being utilized as the core element of Cloud Gateway component. MoleculerJS framework has built-in microservices that can support the above characteristics, such as load balancing [4] or fault tolerance [5]. The latter is also being addressed though the integration with the Kafka [6] event streaming platform, one of the main tools utilized in the PolicyCLOUD project and which is used as a buffering mechanism and a message bus for providing and moving data across the whole data pipeline and across all different analytical components.

Through this ability the gateway is also able to directly ingest incoming data into the appropriate data store based on their privacy level. Therefore, it makes easy to differentiate the queries/requests having to be redirected to the overall data management, analysis and storage system of the project. On top of all these, this component examines and capture the reliability levels of both all the available data sources and their incoming data, thus "feeding" into the PolicyCLOUD platform only the reliable data that comes from only reliable data sources. In this context, the Gateway is able to map all the incoming data sources to specific levels of trustfulness, and thus capturing their reliability. As a result, all the data sources that do not meet the trustfulness criteria are excluded, ensuring the origination of the data sources' incoming data, the adaptive selection of all these available data sources in order to be kept connected into the PolicyCLOUD platform. The component ensures also that the collection of the data comes only from reliable data sources so as to be used for further analysis. Furthermore, in terms of integration with other internal components and mechanisms of the PolicyCLOUD platform, the Cloud Gateway & APIs component has been successfully integrated with the Access Mechanisms in order to ensure that all the required security standards are being met and that specific roles and privileges are being defined precisely. On top of this, the Cloud Gateway & APIs component has also been integrated with OpenWhisk tool which is utilized in the scope of PolicyCLOUD platform in order to provide a serverless, holistic, integrated and end-to-end pipeline of the Data Acquisition and Analytics layer. Finally, the gateway includes an API documentation page, by using Swagger UI, in terms of providing a graphical interface for interacting with the API. The latter facilitates the exploration of all available requests and responses that are listed including also the required parameters.

⁵ <https://moleculer.services/>
www.policycloud.eu

7.4.2 Incentives Management

The overall idea of Incentives Management is to offer a set of tools to identify, declare, track and manage incentives activities to engage the different participants on the policy making process, understanding their motivations in the light of the context. Therefore, this task will provide tools to manage the incentives activities performed with the policies stakeholders, either through closed groups, like some communities evaluating some proposed policies, or even engaged citizens.

The different incentives may be of different types, social, cultural, political or other types. For this purpose, the Incentives Management component may provide to the policy maker access to results from policies on PolicyCLOUD in order to create and manage incentives that relate to these results.

Similarly, the component will manage and keep track of the incentive actions proposed by the policy maker in order to involve the participants and evaluate the outcomes of these actions.

More specifically, and from a theoretical point of view, Incentives Management activities pursue to provide an individual incentives plan that will define a set of rewards corresponding to specific participant actions.

Following the four dimensions introduced by [7] (Malone, 2010): what, who, why, and how, the incentives plan will be pre-established as follows:

- **WHO (participants/requesters):** The Incentives Management task will be focused on engaging citizens, organizations who may be affected by the introduction of policies defined in PolicyCLOUD. In the case of PolicyCLOUD the exact group of citizens and/or organizations will be settled attending the existing use cases and drive by the policy maker.
- **WHAT (actions/tasks):** Is the information exchange, contributions and collaboration expected by the participants.
- **HOW (way or manner):** Define how the participants collaboration is expected. In the case of PolicyCLOUD, the way of collaboration will be established in the context of the existing use cases and drive by the policy maker.
- **WHY (rewards/incentives):** It is aimed to the establishment of different types of incentives (e.g. social, cultural, political, etc.) in return for the participant collaborations done through the execution of existing tasks (what) performed in a concreted way (how). In the case of PolicyCLOUD, the incentives will be established in the context of the existing use cases and drive by the policy maker.

Citing the description included in D3.1 deliverable, the Incentives Management activity will be focused on the following: *“Provide the maximum support to the policy maker... toward a twofold aim: support the policy maker in the incentives identification and help the policy maker in the incentives management”*.

[3]

7.4.2.1 INCENTIVES MANAGEMENT ARCHITECTURE INTEGRATION

The tool will provide policy makers the possibility to declare and track incentives actions. As pictured in Figure 5, the integration point of the Incentives component is the PDT interface where the component frontend will be integrated as an additional entry point for the policy makers, so as to have all the needed components accessible from the same access point. It may also be possible to show to the policy maker information from the policy models' KPIs they have already declared in order to better shape and adapt the incentives actions.

Crowdsourcing tools may be used by the policy maker, but those will be kept totally separate from the Incentives Management component or any other components or modules from PolicyCLOUD. It will be on policy maker discretionary use of the results and information gathered through these Crowdsourcing tools that they may shape and track specific incentive actions with their corresponding policy stakeholders.

As per the backend, the component will be managing the access to the different entities with a corresponding data storage tied to it.

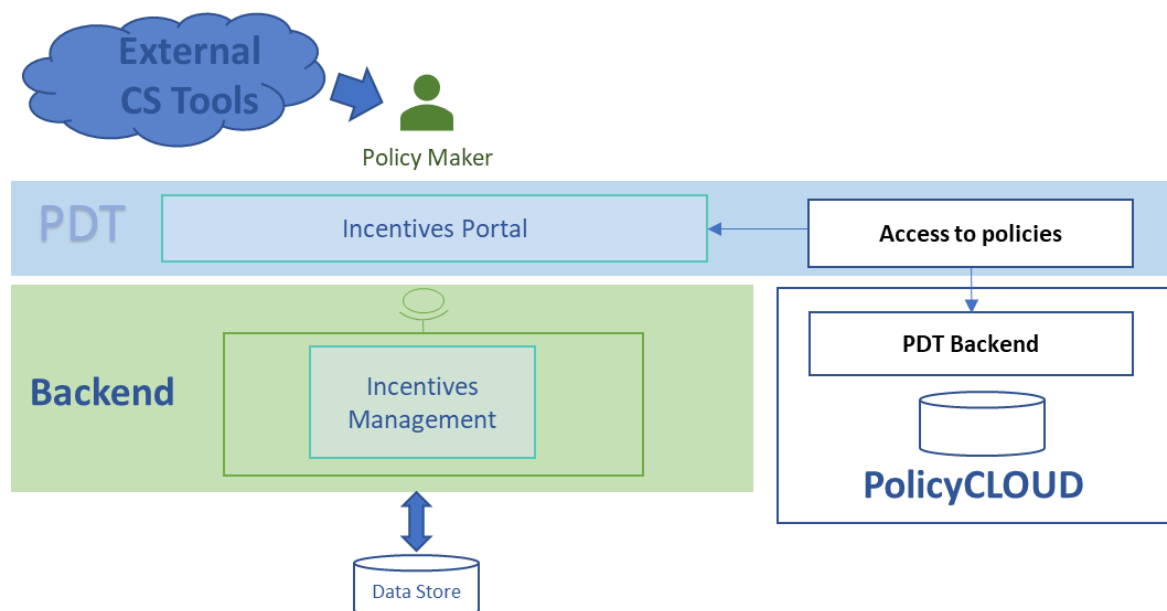


FIGURE 5 – INCENTIVES MANAGEMENT ARCHITECTURE INTEGRATION

For more details, please refer to next Deliverable D3.4 “Cloud Infrastructure Incentives Management and Data Governance: Design and Open Specification 2” [8].

7.4.3 Data Management and Data Stores

In the scope of the PolicyCLOUD project, different challenges are being raised regarding data management, an internal part of the data acquisition process, as data stored into the data repository of the platform are being accessed by different and heterogeneous manners. Firstly, as part of the project itself, two scenarios from two different use case providers have already been integrated to the common platform, while more scenarios are planned to be integrated to the common platform from all four use case scenarios. At the same time, the platform itself is envisioned to be exploited in the future by other cases. Each one of these independent organizations (from the four use cases) is currently using its own data management systems, relying on different types of data schemas, while there is need for a central data repository to fit the needs of all. Secondly, each organization typically has different silos, relying on heterogeneous data stores for data persistency, using completely different data models: from traditional relational database systems, NoSQL databases, Hadoop datalakes etc. Moreover, the PolicyCLOUD vision is to deal with different in nature data, that is, data at-rest which typically refers to data that is permanently stored and various queries are being executed in order to retrieve the results, and streaming data that refers to data that are being continuously inserted to the system without always the need for persistent storage, but with the ability to apply automatic analytics on top of them. Streaming data refers also to external data that is remotely accessed upon demand. Nevertheless, according to the requirements defined in Deliverable D2.1, there is the need for support of hybrid workloads, such as OLTP workloads for managing operational data and ensuring transactional semantics, and OLAP workloads in order to perform analytical queries over the operational data, while ensuring the data consistency. Finally, as operational data usually become obsolete after a certain point in time with rare modifications and in order to cope with analytics over big data, typically the data are being transferred to a data warehouse such as Object Storage, that is more suitable for performing this type of analytics. The requirement in this scenario is to move the corresponding data slices while maintaining data consistency, transparently to the analytical tools, enabling them to use a common interface for accessing data, no matter whether this data resides in the operational data store or in the object storage.

With respect to the design of the overall architecture, the Data Store of PolicyCLOUD is conceptually a central component where data is being ingested (either via a streaming mechanism or with a static data acquisition from external sources) and is being accessed via a common interface by all analytical tools that require data retrieval for their analysis. An additional requirement is to access data that resides in external data sources that are not eligible to be physically imported to the central persistent storage of the platform remaining on premise due to data regulator constraints or due to excessive ingestion/maintenance costs. The central data store component has to provide access to such external sources, via the common interface used by the analytical tools.

At this point, it is very important to distinguish between the major three different types of data sources that the PolicyCLOUD will support: i) ingest-now data, ii) streaming data and iii) external data. With the term stakeholder data we refer to data that belongs to the organization that can be ingested to the

platform via the data acquisition mechanism. With the term streaming data we refer to data that is not static (or *data at rest*) but rather might be generated by IoT devices or coming from a social media feed such as tweeter, and requires a processing in real-time and accumulation for further analytics. Finally, with the term external data we refer to both data that is either not owned by the organization and cannot be retrieved and ingested to the platform, or it cannot be ingested due to privacy considerations, and to data that might be owned by the organization, but cannot be imported due to technological constraints, and thus they are considered as external to the platform. In the following, we provide specific details on how the technology provided by the data stores and data management building block will deal with these three types of data.

- ***Stakeholder data***

In order to address the challenges for data management and overcome the barriers imposed by the data constraints coming from the use cases, the PolicyCLOUD Data Store component will rely on the LXS data repository which natively provides characteristics that are relevant to those challenges, and will be further extended in the scope of the project to cover all aforementioned requirements. More information regarding the characteristics of the datastore can found in the document of Deliverable D4.1.

- ***External data***

The challenge on the isolated silos across different kinds of data stores at each organization is addressed by leveraging the polyglot capabilities of LXS that enables to integrate its query engine with different data stores. Using the CloudMdSQL query language, which is an extension of the standard SQL, the data user can write queries in a unified manner that targets heterogeneous data stores and let the query engine of the PolicyCLOUD datastore to retrieve and merge the intermediate results. This will overcome – to the extend possible - the need for accessing data that are stored in different silos inside an organization or in external sources. The polyglot capabilities of the data store are also important for the datalake capability of enabling query processing of unstructured data, which is typically used in the majority of the datasets provided by the existing and future scenarios.

- ***Streaming data***

Often it becomes necessary to manage streaming data combined with data at rest, in order to correlate events with operational data and/or update a dataset based on an event. This is a bottleneck for traditional databases when streams arrive at large scale, as they are incapable of dealing with those operational workloads at that high rate. Due to the scalable transactional processing provided by the LXS datastore and its additional interface that allows directly accessing its storage layer, it can support data ingestion coming from streams.

Moreover, due to its extended capabilities for live aggregations, it can support the combination of streaming events with data at rest which requires data expensive operations (i.e. average value of a field) that can be supported by traditional data management systems, where usually the solution of caching the results is preferred over the consistency of the result with respect to transactional

semantics. This is very important when pre-processing needs to take place over a stream, which is a typical scenario that PolicyCLOUD targets.

Apart from dealing with those three different types of data, the data management mechanisms of the platform will benefit from the results for the EU H2020 project BigDataStack and its Seamless Analytical Framework, where similar scenarios with regards to the movement of “cold” data from an operational to an object store are being addressed. That will allow for data to be moved to the object storage at runtime, transparently to the user by ensuring data consistency and without the downgrade of the performance during the movement of the data. The data repository supports standard SQL statements via the common JDBC, and splits the data operations so that they can be executed in both underlying stores, and merges the intermediate data in order to return the same result as if the data was stored in a single database. By doing this, the data analyst will not have to alter its implementations in order to support scenarios where there is the need to combine data from both stores. In the scope of PolicyCLOUD, the prototype firstly developed in the EU BigDataStack project will be further developed to cope with the scenarios defined here, with the plan to increase its current technology readiness level (TRL).

7.5 Ethical and Legal Compliance Framework

7.5.1 Ethical and Legal Compliance Framework

To maximize societal acceptability and trust in PolicyCLOUD, and the policies developed with PolicyCLOUD's assistance, the PolicyCLOUD consortium is aware of the need to carrying out an extensive and in-depth analysis of relevant legal, regulatory, societal and ethical aspects, define appropriate requirements to address all relevant aspects identified, and see to an optimal embedding of those requirements into the design of the solution – including a thorough evaluation as to the extent to which this has been successful. Special attention must be paid to ethical and societal issues which may be triggered throughout the project.

Therefore, it is necessary to identify a set of controls – which may pertain, *inter alia*, to platform dimensions, features, and functionalities – and their links to the range of significant new practices enabled by the platform which are relevant from a legal, regulatory, ethical and/or social perspective. These controls must remain aligned with the various iterations of the development of the platform, and the specificities of the different use cases which serve to demonstrate the platform's capabilities. The **Ethical and Legal compliance framework** (task T3.5) thus aims at analysing and giving guidance on the legal, regulatory, ethical and societal requirements by which PolicyCLOUD should abide.

Within these controls, particular attention is currently being paid to the choice of data sources and data extracted from those sources, as well as to the admissibility of their use by the controllers/contributors bringing in data into the PolicyCLOUD infrastructure. This counts for “personal data”, as defined in the GDPR, but also for other types of data – such as “open data” – which may involve legal issues with regard to the protection of intellectual property, including the protection of databases and trade secrets. One such control being envisioned is a “clearing mechanism” of sorts, to allow specific categories of data to be ingested by the platform based on the requirements set by end-users and the extent to which legal permission for such ingestion exists – this mechanism could be implemented at the stage of “data acquisition”, though it is to be determined whether it is to operate before or after data is processed through the cloud gateway.

Furthermore, within the controls defined in task T3.5 are controls to address the (shared) responsibilities of PolicyCLOUD, the partners and stakeholders when processing “personal data” under the GDPR, as well as basic considerations for the admissibility of use of such data. This also includes the role of the cloud provider (i.e., the organization/entity which, ultimately, will act as the provider of the PolicyCLOUD platform) as a controller or as processor, depending on the specific processing activities at stake. Under the principle of data minimization, the possibility to rely on aggregated or anonymous data (as opposed to identified or identifiable “personal data”) is under exploration with the relevant technical and use case partners. Moreover, guidance is being given on the requirements of security and confidentiality, as well as data protection by design and default, taking also into account the work performed in the context of task T3.6 (Data Governance Model, Protection and Privacy Enforcement).

Task T3.5 also addresses other ethical, as well as societal requirements, in order to align the platform with these to the greatest extent feasible. This aims at maximizing societal acceptability and trust in PolicyCLOUD and, through it, in policies developed with the aid of PolicyCLOUD. As such, controls have been defined to address the reliability of data and prevention of false raw data – which is of primordial importance to mitigate the risk of incorrect conclusions being derived from the platform’s output, which may potentially lead to incorrect policy-making decisions. This includes considerations as to the potential for abuse of data in order to intentionally manipulate a decision-making process, as well as the potential for abuse of the platform as a whole. The need to ensure explainability of output provided by the platform, and that policy-makers retain critical judgment when interpreting platform output and making decisions – as they retain ultimate responsibility for such decisions – is also deemed an essential ethical and societal concern, reflected in appropriate controls. Further analysis will be necessary with regard to the data marketplace (in particular, on the liability for data in a later period of commercial exploitation) and the incentives management system.

In the context of deliverable D3.3 [9]⁶, an analysis was carried out of the relevant legal, regulatory, ethical and societal issues detected in relation to PolicyCLOUD, with a synthetic review of the existing debate and literature provided. These issues (the main of which were expressed above) were addressed in general terms, and also from the perspective of the specific platform components and use cases, based on the information available as of the date of completion of the deliverable. From this deliverable, an initial set of controls was defined, to be used to ensure the platform’s adherence to the principle of compliance by design (as better explained in the following section). This set of controls – or checklist – is to be adapted over time, considering the relevant developments which are underway for several components of the PolicyCLOUD platform (not least of which are the Data Marketplace and Incentives Management components), as well as for the four uses cases (with the use case partners defining further relevant data sources and use case scenarios).

In subsequent iterations of deliverable D3.3 (D3.6 and D3.9, due on M22 and M34 respectively), to the extent that the development of the platform suggests that this may be relevant and necessary, deliverable D3.3 will be expanded to include an analysis of additional or particular concerns relevant at a local level – considering, in particular, the jurisdictions in which each of the use case partners of the currently envisioned four use cases are established. Additionally, subsequent iterations of deliverable D3.3 will be duly updated to reflect all these relevant developments, with sections being expanded on or amended as needed. Finally, in an effort to increase the effectiveness of this framework, subsequent

⁶ With regards to legal and regulatory issues, the scope of the analysis, in the context of this deliverable, was generally limited to EU and international law, without exploring in detail the specific national and/or local requirements related to the countries and jurisdictions in which the use cases are implemented. Nevertheless, where specific analysis on local and/or national regulations shall result as appropriate and/or necessary, this has been highlighted as a field for which further research is needed and that will be consequently developed in the next versions of the deliverable, to be released at M22 and M34.

iterations will rely on feedback provided by the Consortium to improve on the practical requirements identified (including in terms of feasibility of their implementation), further tailoring them to the specificities of the platform and use cases. It is expected that the final iteration of this framework (D3.9) will reflect a set of requirements which are understood and accepted by the Consortium as vital to ensure the ethical and legal soundness of the project.

7.5.2 Ethical and Legal Compliance Framework Integration with Use Cases and technology

7.5.2.1 APPLICATION OF THE COMPLIANCE BY DESIGN PRINCIPLE

In this section, we will analyze how, during the development of the PolicyCLOUD project, compliance with the identified legal, regulatory, ethical and societal requirements is being assessed. With specific regards to data protection and privacy issues, we will also define a methodology for the implementation of a DPIA which will be conducted with regards to each of the four use cases.

7.5.2.2 COMPLIANCE BY DESIGN APPROACH

To address all the relevant ethical, legal, regulatory, and societal risks related to the project, a compliance by design approach is being adopted.

Compliance by design means applying a systematic approach to integrating relevant compliance requirements into tasks and processes. The effective implementation of this principle is based on the detailed and structured analysis of all the applicable requirements (as initially identified in deliverable D3.3), followed by translation of those requirements into workable compliance processes [10].

A three-stage approach is being applied:

1. The **first stage** is dedicated to the **identification and the assessment of relevant requirements**.

This was accomplished with deliverable D3.3, at the end of which an initial list of requirements/controls was defined, based on the assessment developed within the deliverable of the applicable legal, regulatory, ethical and societal issues potentially triggered by the platform's technical components and use cases.

2. The **second stage** includes the **analysis on how the rules apply to individual processes**.

This is under development – the initial list of requirements/controls has been broken down into smaller “checklists” – i.e., Legal/Ethical Checklists – in which individual controls are allocated to different Work Packages within the project. As of the date of this deliverable, Checklists have been developed for the main components under the responsibility of WP3, WP4 and WP5, and for the use cases under the responsibility of WP6. Further Checklists are to be created for specific components which are still at an early stage of development at present (e.g., the Data Marketplace and Incentives Management components).

The relevant Checklists have been shared with the WP Leaders, with a view to refining each Checklist in terms of the relevance and feasibility of the identified controls. The goal is also to identify specific owners for each refined control, so as to follow-up with those owners on the definition of specific, practical measures to implement each control, and on the actual implementation of those measures. Feedback has been exchanged with the WP Leaders, and the initial consolidation of the mentioned Checklists is nearing its completion.

3. The **third stage** focuses on the **design and implementation of a roadmap**.

Once the Checklists have been initially consolidated with the assistance of the WP Leaders, and specific owners have been identified for each refined control, those owners are to be engaged so as to (1) identify specific, practical measures to implement each control assigned to them, and (2) define a roadmap for the implementation of those measures.

The final step in this process is to maintain continuous engagement with the WP Leaders and control owners, so as to assess the implementation process in accordance with the defined roadmaps, addressing any concerns which may arise over time (e.g., adjusting controls as needed to fit new platform developments).

7.5.2.1.2 Data protection and privacy

To address the issues related to the project and the use cases concerning data protection and privacy (i.e., for use cases where “personal data”, under the GDPR, are to be processed), data protection impact assessments (“**DPIAs**”, under Art. 35 GDPR) will be implemented, and the results of those assessments will be presented in the context of deliverables D3.6 and D3.9 of the project.

Through these DPIAs, PolicyCLOUD will assess the processing operations to be performed, as well as the technologies, tools, and systems to be used, in relation to each specific use case scenario in which the processing of personal data is envisioned, to identify inherent risks in a structured manner. Furthermore, these DPIAs will be used to identify measures which can be implemented to bring those risks down to acceptable levels. The DPIA reports will contain a systematic description of the envisaged processing operations, the purposes for which personal data will be processed, an assessment of the legitimate interests pursued (where applicable), an assessment of the necessity and proportionality of the operations in relation to those purposes, an assessment of the risks to the rights and freedoms of data subjects, and a description of the measures envisaged to address those risks, as noted in Art. 35, par. 7 GDPR [11].

These DPIAs will be performed according to the methodology defined in the international standard ISO/IEC 29134.

The process for the performance of the DPIAs will include:

1. A **preparation phase**, during which the DPIA teams will be set and provided with direction, the DPIA plan will be prepared, the necessary resources will be determined, and the relevant stakeholders will be engaged.
2. A **performance phase**, during which the information flows of personal data will be identified, the implications of the relevant use case scenario (in the context of the project) will be analyzed, the data protection and privacy risks will be assessed, a risk treatment plan will be defined and relevant privacy safeguards will be determined.
3. A **follow-up phase**, during which a DPIA report will be prepared and published, and the risk treatment plan will be implemented. In this context, also a review and/or reaudit program of the DPIA will be defined, to monitor both the correct implementation of the risk treatment plan and of the potential changes to the previously assessed personal data processing activities.

7.6 Layer 2 - Data Acquisition and Analytics

Components: Data Cleaning (T4.2), Data Interoperability (T4.2), Data Fusion (T4.1), Situational Knowledge Analysis (T4.3), Opinion Mining (T4.4), Sentiment Analysis (T4.4), Social Dynamics (T4.4), Behavioral Analysis (T4.5), Optimization and Reusability (T4.6)

7.6.1 Data Acquisition and Analytics – Positioning & Goals

In this section we provide the high level architecture of the Data Acquisition and Data Analytics tasks, which is responsible for ingesting the data from various sources while applying filtering and initial analytics, and preparing it for deeper analytics on longer term storage (DB, object storage).

The relevant part from overall architecture is shown in Figure 6 for convenience. This part focuses on Data Acquisition and Data Analytics over which the integrated processing will be applied.

More specifically, data fusion tasks are integrated with the initial analytics and data processing tasks (e.g. filtering, validation and cleaning). Applying deeper analytic tasks are performed in collaboration with the continued data fusion (e.g. moving older data from DB to object storage).

From the aspect of work packages partitioning, this layer is under the responsibility of WP4 (Reusable Models & Analytical Tools) and its tasks, with a strong relation to Task 3.3 (Cloud Gateways) and Task 3.6 (Data Governance Model, Protection and Privacy Enforcement) of WP3. In Figure 7 we show the conceptual model from the work packages partitioning point of view and the WP4 interfaces to WP3 below and WP5 above, as provided in the Grant Agreement document.

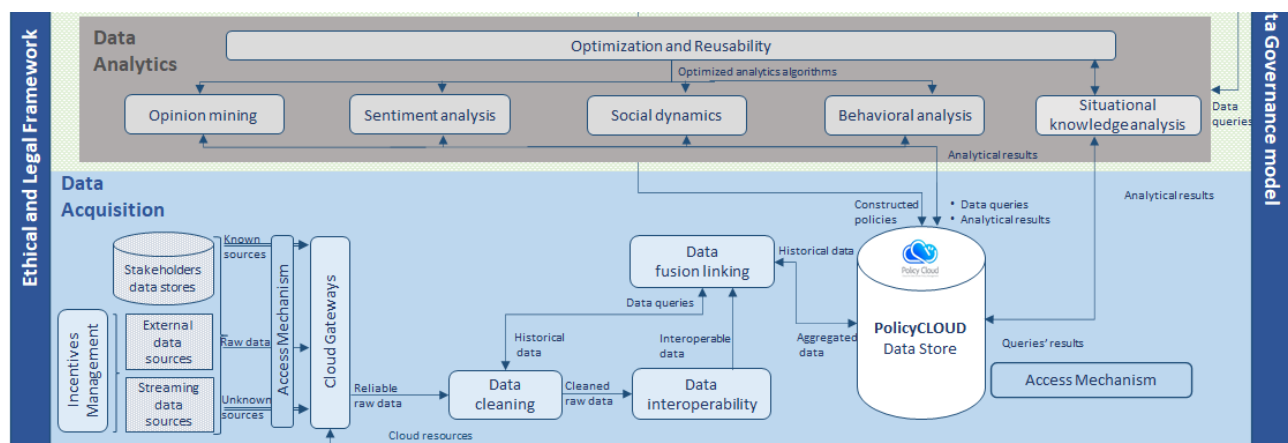


FIGURE 6 – PART OF THE POLICYCLOUD OVERALL ARCHITECTURE DIAGRAM RELEVANT TO WP4

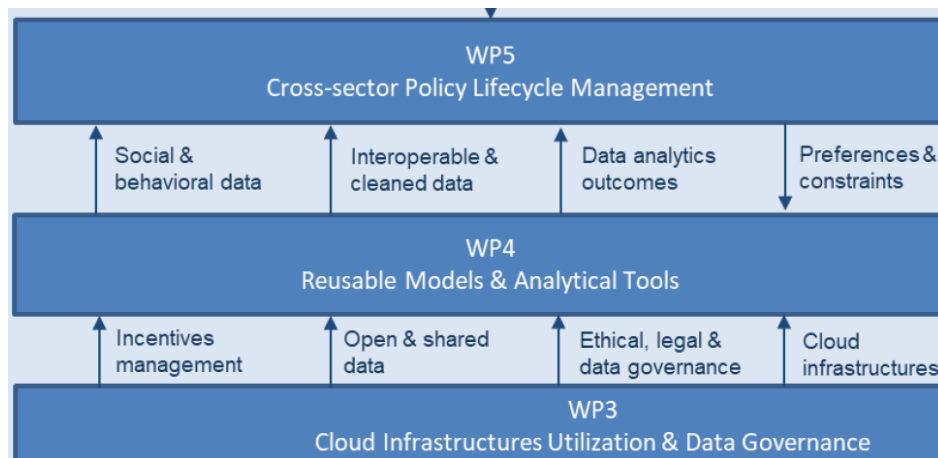


FIGURE 7 – WP4 INTERFACE WITH WP3 AND WP5

The major goals of Data Acquisition and Analytics layers are on par with WP4 defined goals:

- Data fusion and aggregation – for different data sources types.
- Data cleaning ensuring quality of information, sources reliability assessment, reliability-based selection of information sources.
- Sentiment analysis techniques for policy assessment.
- Analysis of the social and behavioural data and requirements provided by social science experts for data selection in a given case.
- Decoupling of the analytical models and tools from the underlying infrastructure and datastores, assuring their reusability.

7.6.2 Extensibility and Reusability of Analytic Functions

The architecture of the Data Acquisition and Analytics layers will provide extensibility and reusability of analytic functions. New analytics functions (services) can be registered into PolicyCLOUD and reused for applying analytics on new and existing registered data sources. The decided alternative at this point is a registration as serverless functions that are activated on demand, either by a direct PolicyCLOUD user request or by event/rule. There are two types of functions:

1. Ingest analytics / transformation function, which will be used to apply initial analytic and/or transformation on the data fusion path of data sources.
2. Rest data analytic function which will be activated upon PolicyCLOUD user action on specified data source (which was already ingested) to provide analytic results for policy decisions.

The design details of analytic functions registration and activation are provided in deliverable D4.1.

7.6.3 Data Cleaning

The Data Cleaning component will offer all the appropriate algorithms and techniques for detecting and correcting (or removing) corrupt or inaccurate records from all the collected data that will be retrieved as an input from the Cloud gateways component. More specifically, this component will be responsible for identifying all the incomplete, incorrect, inaccurate or irrelevant parts of this data, and then replacing, modifying, or deleting the dirty or coarse data. Thus, possible missing, irregular, unnecessary, or inconsistent data will be found and totally cleaned. Especially dealing with missing data is one of the most tricky but common parts of the data cleaning process since most of the models do not accept missing data. To this context, the Data Cleaning component will detect and totally clean all the missing data by combining techniques such as the Missing Data Heatmap, the Missing Data Percentage List, as well as the Missing Data Histogram, thus extracting quite accurate and reliable results. With regards to irregular data, cleaning is made possible by using techniques such as the Histogram and the Descriptive Statistics for the numeric values, and by exploiting the Bar Chart for categorical values.

Regarding the unnecessary data, since it refers to data that will not add any value to the PolicyCLOUD overall platform, by constructing the corresponding rules and constraints, all the uninformative/repetitive, irrelevant values, as well as the duplicates will be automatically detected and erased. Finally, since any possible inconsistent data will be automatically corrected it is also crucial that all the collected datasets will follow specific standards to fit the corresponding PolicyCLOUD data models. As soon as all the data is fully cleaned it will be sent into the Data Interoperability component for further utilization.

7.6.4 Data Interoperability

The Data Interoperability component aims to enhance the interoperability of analytics processing in the PolicyCLOUD project based on data-driven design, coupled with linked data technologies, such as JSON-LD [12], and standards-based ontologies and vocabularies to improve both semantic and syntactic data and dataset interoperability. The provided Interoperability Component seeks to extract semantic knowledge and good quality information from the cleaned data that will be the input to its system, as shown in the initial architecture of the overall project. This knowledge, shaped in a machine-readable way, will be used in next tasks for Big Data analytics, Opinion Mining, Sentiment Analysis etc.

One of the preliminary steps of this component is to identify relevant, publicly available, and widely used classifications and vocabularies, such as the Core Person Vocabulary provided by DCAT Application Profile for Data Portals in Europe (DCAT-AP), that can be re-used to codify and populate the content of dimensions, attributes, and measures in the given datasets. Hence, this component aims to adopt standard vocabularies and classifications early on, starting at the design phase of any new data collection, processing or analytical components. Using for example NLP techniques and tools like Text Classification, NER, POS tagging and even Machine Translation [13] [14] we can identify and classify same entities, their metadata and relationships from different datasets and sources and finally create cross-domain vocabularies in order to identify every new incoming entity. Likewise, in order to create and enhance semantic interoperability between classifications and vocabularies this component seeks to engage in

structural and semantic harmonization efforts, mapping cross-domain terminology used to designate measures and dimensions to commonly used, standard vocabularies and taxonomies. Thus, by implementing a “JSON-LD context” to add semantic annotations to interoperability component’s output, the system will be able to automatically integrate data from different sources by replacing the context-dependent keys in the JSON output with URIs pointing to semantic vocabularies, that will be used to represent and link the data. This mechanism enhances information by connecting data piece by piece and link by link, allowing for any resource (authors, books, publishers, places, people, hotels, goods, articles, search queries) to be identified, disambiguated and meaningfully interlinked.

7.6.5 Data Fusion with Processing and Initial Analytics

In this section we present the architecture for integration of all the tasks relevant to data fusion. We demonstrate this integration by an end-to-end example data fusion scenario, from a Twitter social network data source. The data is fused, cleaned, validated and initially analysed for extracting the relevant knowledge insights which are then persistently stored for future deeper analytics and possibly generating immediate alerts.

The participating tasks in this scenario are:

- T3.3 Cloud Gateways.
- T4.1 Cross-sector Data Fusion Linking.
- T4.2 Enhanced Interoperability & Data Cleaning.
- Potential initial analytics by T4.3 Situational Knowledge Acquisition & Analysis T4.4 Opinion Mining & Sentiment Analysis and T4.5 Social Dynamics & Behavioral Data Analytics.

The framework for data fusion and analytics will either be based on Apache Spark Streaming open source (<https://spark.apache.org/streaming>), KSQL (<https://github.com/confluentinc/ksql>) or Serverless engine based on Apache OpenWhisk (<https://openwhisk.apache.org>). In Figure 8 we depict the end-to-end data path for this scenario.

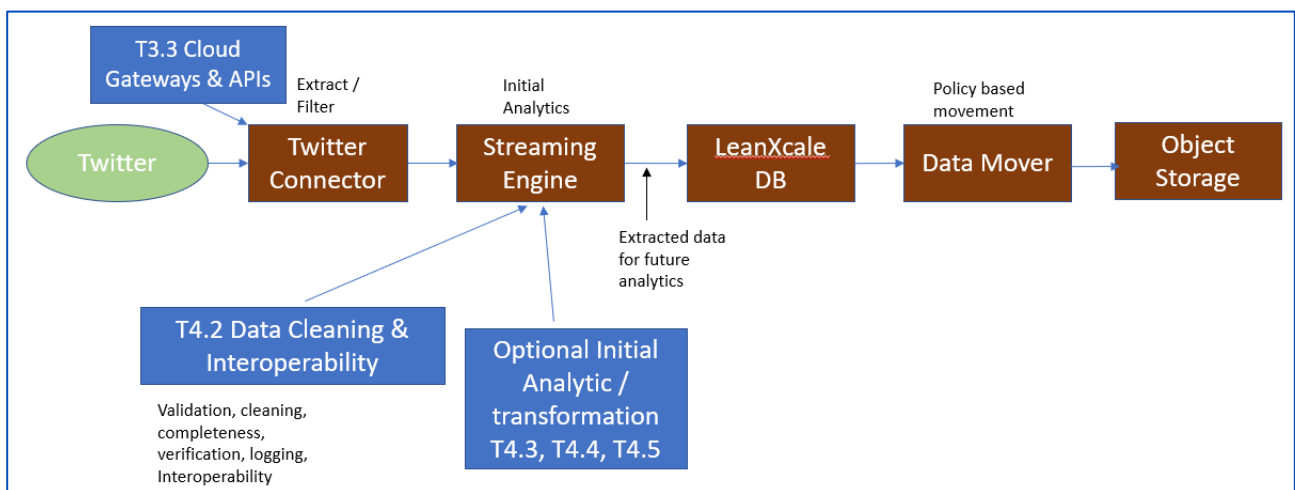


FIGURE 8 – THE STREAMING DATA PATH

Task 4.1 (Cross-sector Data Fusion Linking) is responsible for the overall data path and streaming framework in this scenario. The Twitter connector will be implemented by task T3.3 (Cloud Gateways) and will create the stream of relevant data into the Streaming engine. It is expected to apply basic filtering by policy rules that are active in the PolicyCLOUD framework (resulting from actual policies that are subject for validation). The data cleaning and reliability validation will be performed by Task T4.2 which will be running within the streaming engine. Optional initial analytics on the streamed data may be performed by tasks T4.3 (Situational Knowledge Acquisition & Analysis), T4.4 (Opinion Mining & Sentiment Analysis) and T4.5 (Social Dynamics & Behavioural Data Analytics).

At the end of the data path, the Data Mover is responsible for moving older data from hotter storage (DB) to a colder (object storage) periodically, according to certain policy rules (discussed more in details in the next section).

7.6.6 Seamless Analytics on Hybrid Data at Rest

In this section we provide the architecture for applying the analytics functions on the data at rest, which is combined of knowledge insights extracted within the data fusion, as well as more 'raw' data (however still after cleaning and validation processes). The "right" side of the data path in Figure 9 present a periodical movement of older data from hotter storage (DB) to a colder (object storage) according to policy rules, which addresses the scalability and cost aspects of dealing with big data. Object storage is the perfect platform for storing big data for analytic purposes when no future modification of the data is expected, while the DB platform is superior performance-wise for analytics on the hotter data. The requirement is to apply seamlessly analytics on both hot (in the DB) and cold (in the object storage) data. The basic technology of data movement and seamless analytics was developed by IBM and LeanXcale partners in the BigDataStack H2020 project (<https://bigdatastack.eu>) and will be exploited and adapted for PolicyCLOUD.

The participating tasks for the provided functionality are:

- T4.1 Cross-sector Data Fusion Linking
- T4.3 Situational Knowledge Acquisition & Analysis
- T4.4 Opinion Mining & Sentiment Analysis
- T4.5 Social Dynamics & Behavioral Data Analytics
- T4.6 Optimization & Reusability of Analytical Tools

As depicted in Figure 9 the framework for data movement and seamless analytics will be provided by overall task T4.1 (Cross-sector Data Fusion Linking). Task T4.6 (Optimization & Reusability of Analytical Tools) Optimization aspects (to be developed in the later phases of the project) will additionally provide the interface for seamlessly applying the analytic tasks as T4.4 (Opinion Mining & Sentiment Analysis) and T4.5 (Social Dynamics & Behavioral Data Analytics) on the data at rest.

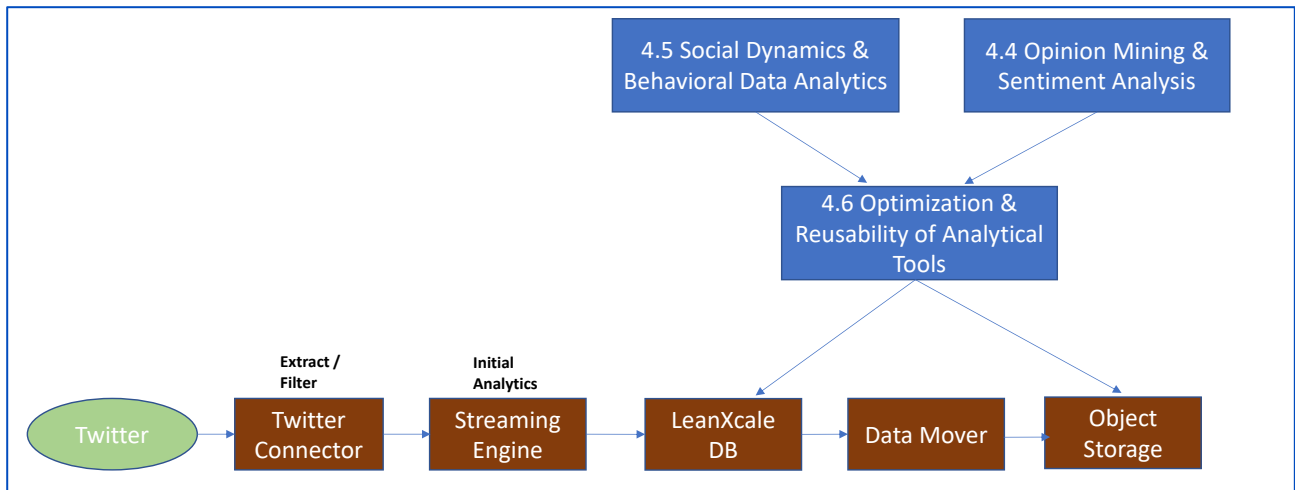


FIGURE 9 – SEAMLESS ANALYTICS ON INGESTED DATA

7.6.7 Situational Knowledge Analysis

In the context of PolicyCLOUD the Situational Knowledge Acquisition (SKA) component brings the capability of acquiring knowledge from the Data & Policy aspects of the platform. The extracted knowledge will be used to influence the decisions taking place based on the PolicyCLOUD system.

The following capabilities will be provided through the SKA component:

- It will deal with real-time facts (such as data from sensors) from which it will derive situational knowledge.
- A situational knowledge model (SKM) will be provided for structuring the knowledge acquired. This data model will contain a high-level description of real-world situations (context) which are the interest of the PolicyCLOUD system. The model will be defined by the use cases based on the types of situations/context to be acquired.

Some of the characteristics of the component will be:

- Feature Extraction. The extraction knowledge stage will be done through Feature Extraction (ML) techniques able to create/derive new situational features from existing ones. This extraction step will be enhanced by the situational knowledge model which will guide the derivation of new features or the abstraction of existing ones.
- Dataset clustering and categorization. Data categorization must be possible in a very flexible way according to the structure envisaged for formal descriptions of business fitted entities [15] (Olszewski, Robert, 2001).

7.6.8 Opinion Mining

The following tasks have been identified as being the basic activities to be performed in the context of opinion mining and sentiment analysis. The identification of these tasks is the result of internal conversations with use case owners, in order to extract information and needs for data analytics based on the various scenarios.

- *Opinion Mining.* Observe events and social attitude in respect to specific topics.
- *Named-entities recognition.* Identification of specific entities (users, locations, groups, ...) cited on text.
- *Graph Analysis.* This task will develop an additional component that will perform further analytics by generating a “contributor graph” based on the contributors that are talking about the policies. This graph can be built on top of any platform with enough information about the contributors (e.g. Twitter), in order to determine the main influencers and create groups of similar contributors. This requirement will be refined based on the data that will be provided by each pilot. Other mechanisms such as page-rank, will be developed to generate the common analysis on graphs.

A specific focus will be devoted to particularities of social networks, such as:

- *Hashtags Detection,* identification of Twitter style hashtags from text.
- *Twitter Hashtags and Mentions Tacking,* find and monitor mentions on Twitter regarding specifics hashtags or topics.
- *User Monitoring,* identification and monitoring of most popular users who comment about specific hashtags or topics.

Additional analysis such as social media-based Location Surveillance or Topic-related expressions identification (identification of new words or expression which might have hidden relationships with known ones) can be also objective of T4.4 task.

This component will follow the same approach as the sentiment analysis component using Apache NiFi to create a pipeline in a modular way to achieve the described objectives.

7.6.9 Sentiment Analysis

This component will perform a sentiment analysis based on the input received from the pilots about their policies. This input could come from what the citizens say in social media channels, from platforms owned by the pilot (getting feedback on various subjects), or other channels that will be discussed through the duration of the project. Having this input as also additional information extracted about a specific topic (such as which entities are involved), a sentiment will be assigned (Positive, Negative, or Neutral). To achieve this, it is needed to train the sentiment models with different types of data from different scenarios in order to receive the best accuracy possible.

The development of this component will take advantage of powerful tools such as Apache NiFi, in order to create pipelines in an easy and modular way to be adapted to vary situations without the necessity of repetitive working. It will have a common NLP part to analyse the text arriving as an input from different

sources (social media, text files, or others). The sentiment value assignment for each text will be stored in the database provided by PolicyCLOUD to be used by other components.

7.6.10 Social Dynamics

The Social Dynamics component will consist of a concurrent, web-based environment for social simulation. The environment will allow the user to create graph-based population models online. These models will satisfy various parameters set by the user in terms of size, individual characteristics affecting social behavior, link characteristics, individual and connection dynamics. In addition, it will be able to upload appropriately structured population data from databases conforming with these parameters. Individual characteristics will consist of sets of variables that capture the relevant attributes for each individual in the model. Link characteristics will specify a set of variables used for the creation of weighted links between individuals. Individual dynamics will consist of a set of rules describing the conditions under which individual characteristics can change and the ways these changes can affect individual characteristics. In an analogous way, connection dynamics will consist of a set of rules describing the conditions under which link weights can change and the ways these changes can affect link characteristics. A special-purpose modelling language will be developed that will allow users to specify all these parameters online in the simulation environment. Based on these specifications, the environment will be able to simulate in real-time the dynamics of such populations and store the results in a database for further processing by interested parties. The environment will exploit opportunities for the breakdown of the tasks in each simulation into concurrent units that will allow the simulator to optimize its use of computational resources.

7.6.11 Data Acquisition and Analytics Integration

7.6.11.1 ARCHITECTURE INTEGRATION

A testbed has been installed for the containerized environment at <https://indigo-paas.cloud.ba.infn.it>. The testbed consists of a Kubernetes cluster with master and two nodes, 8GB of RAM and 4 vcpus each.

The following components have been installed as containers on the testbed:

- OpenWhisk
- Kafka
- Leanxcale datastore
- A Cloud Gateway

For the testbed the initial gitlab instance at <https://registry.grid.ece.ntua.gr/> was used (before it was also moved to the cloud environment) providing repositories for the project code and containers.

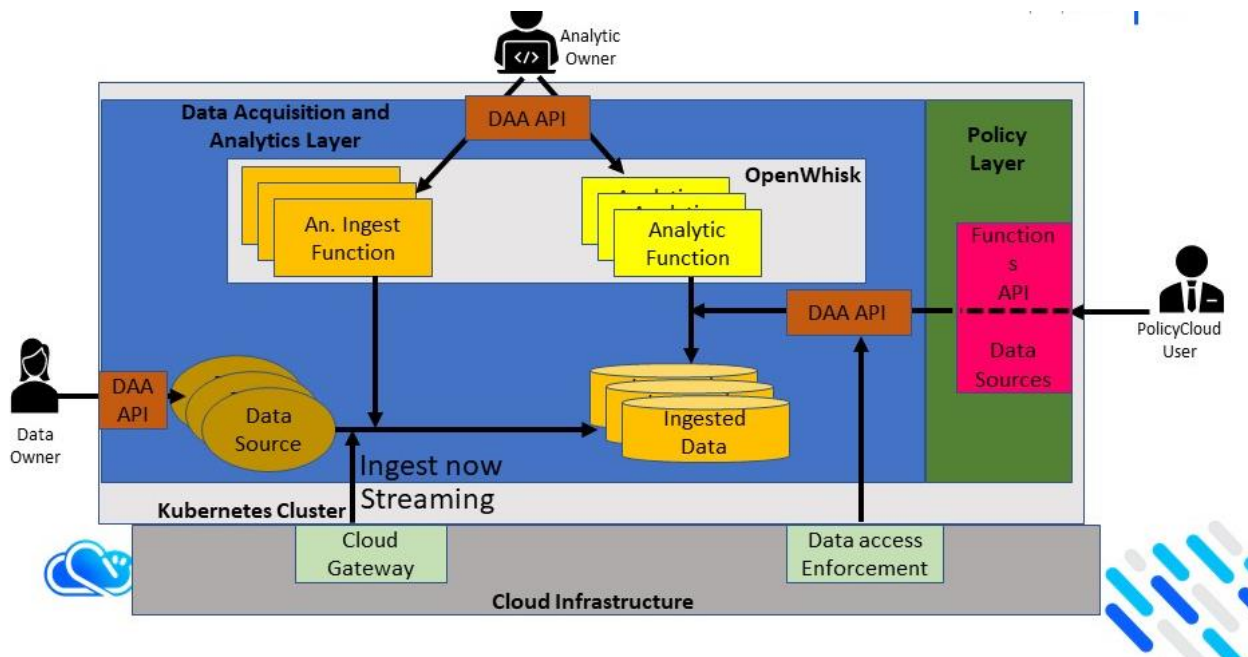


FIGURE 10 - WP4 CONSTITUENT ARCHITECTURE

The architecture implemented shown in Figure 10 demonstrates the following capabilities:

1. Registration of analytic ingestion functions which are implemented as OpenWhisk serverless functions. The ingest functions can be written in any OpenWhisk supported programming languages (e.g. Java, Python, node.js).
2. Registration of data sources (stream/ingest-now). From these data sources, data is imported to PolicyCLOUD's backend datastore (the LXS database) while being transformed by the ingest functions described in the previous paragraph. In the process, Kafka is used for buffering.
3. Analytic functions can be registered and then run on ingested data. The data is read from the data store (the LXS database), then processed while the output is presented to the PolicyCLOUD user who has initiated the function invocation.

7.6.11.2 INTEGRATION WITH THE KUBERNETES CLUSTER

Integration of OpenWhisk, Kafka and LeanXcale database with Kubernetes has been successfully achieved based on extensive test during the trial period described in previous sections. This integration has been fundamental for the successful implementation of the PolicyCLOUD environment.

7.6.11.3 INTEGRATION OF THE SOCIAL DYNAMICS COMPONENT

The integration of the social dynamics component with the overall PolicyCLOUD environment is progressing. Currently, social dynamics runs as a stand-alone web-based component. This component is

an interactive meta-simulation tool as it provides a modelling environment for developing social simulations, it automatically executes simulations and it reasons about the suitability of different policy alternatives. To this end, it accepts/generates a variety of inputs (policy models, high volumes of population data, simulation models, evaluation criteria for policy alternatives), it generates high data volumes as outputs, it requires significant computational resources (memory/cores) for simulation execution, and, because of its interactive nature, it assumes not a one-shot but a constant interaction with its user to develop, run and evaluate alternatives. We are currently analysing the pros and cons of integrating such a component in a serverless paradigm versus using a virtual server for this component from a cloud provider that can communicate via messages with the data store and the rest of the analytics components. Our analysis is performed in terms of the resource, scheduling and reliability requirements, along with the response times that will be feasible under a constant user interaction scenario in both integration alternatives.

7.7 Layer 3 – Policy Management Framework

Components: Policies Modelling (T5.2), Policies Implementation (T5.1), Policies Clusters (T5.4), Policies Experimentation (T5.5), Policies Evaluation (T5.6)

7.7.1 Policy Modelling & KPIs Identification

The Policy Model Editor (PME) is the component that supports and guides the policy maker (PM) to effectively model policies by selecting a data schema, applying relevant Key Performance Indicators (KPIs) or setting new ones with simple linear functions, and creating a set of rules (criteria). As for the existing policies, the PM shall name a description with a set of rules (criteria) which applies the values of a specific data schema and KPIs.

7.7.2 Middleware for Policies

A middleware based on .NET Core has been designed and implemented as the adapter pattern to retrieve data from the policy datastore. At the other end of the adapter lies a REST API as a mechanism that allows policies to be modelled and designed based on specific structural representations.

For more details, please refer to Deliverable D5.2 "Cross-sector Policy Lifecycle Management: Design and Open Specification 1".

7.8 Layer 4 - Policy Development Toolkit

Components: Policy Development Toolkit (T5.3), Data Visualization (T5.3)

7.8.1 Policy Development Toolkit and Data Visualization

The Policy Development Toolkit (PDT), along with the Policy Model Editor (PME), constitute the Front-End of the PolicyCLOUD platform. They integrate several sub-components to enable policy makers (PMs) to create, update and validate their policies. The PM will trigger the underlying analytics mechanisms to provide the corresponding quantitative information, while integrating the visualization component to ensure that the results are presented in a meaningful way. It includes mechanisms to explore and incorporate available analytics into new or existing policy models. The PM will set Key Performance Indicators (KPIs) that support the policy in focus, which will be calculated through the triggering of selected suitable analytics along with the provision of the respective parameters regarding datasets, temporal or spatial constraints, population filtering etc.

For the visualization of analytical tools results, the PolicyCLOUD platform provides a reporting tool that enables to build visual analytical reports. The reporting is produced from analytical queries and includes summary tables as well as graphical charts resulted from the analytics. The dashboard will be adaptable, since it will enable to include different charts with the KPIs chosen by the PM and a set of transformation operators that can aggregate and correlate the received policies KPIs.

The PDT directly interacts with the Data Acquisition and Analytics (DAA) Layer, the datastore and the integrated visualization as presented in the next section.

7.8.2 PDT Architecture

The present section describes the functional architecture of the Policy Development Toolkit (PDT). As a single page web application, PDT hides the complexity of the system dataflow to provide to policy makers (PMs) an integrated Decision Support System (DSS) towards the application of evidence-based Public Policies (PPs).

The general interconnection of the PDT with the other PolicyCLOUD components is illustrated in Figure 11. PDT may be considered as the point of integration and interaction of the platform with the PMs. Through the PDT, the PM will be able to question the platform data and exploit the analytics tools to perform policy modelling and evaluation.

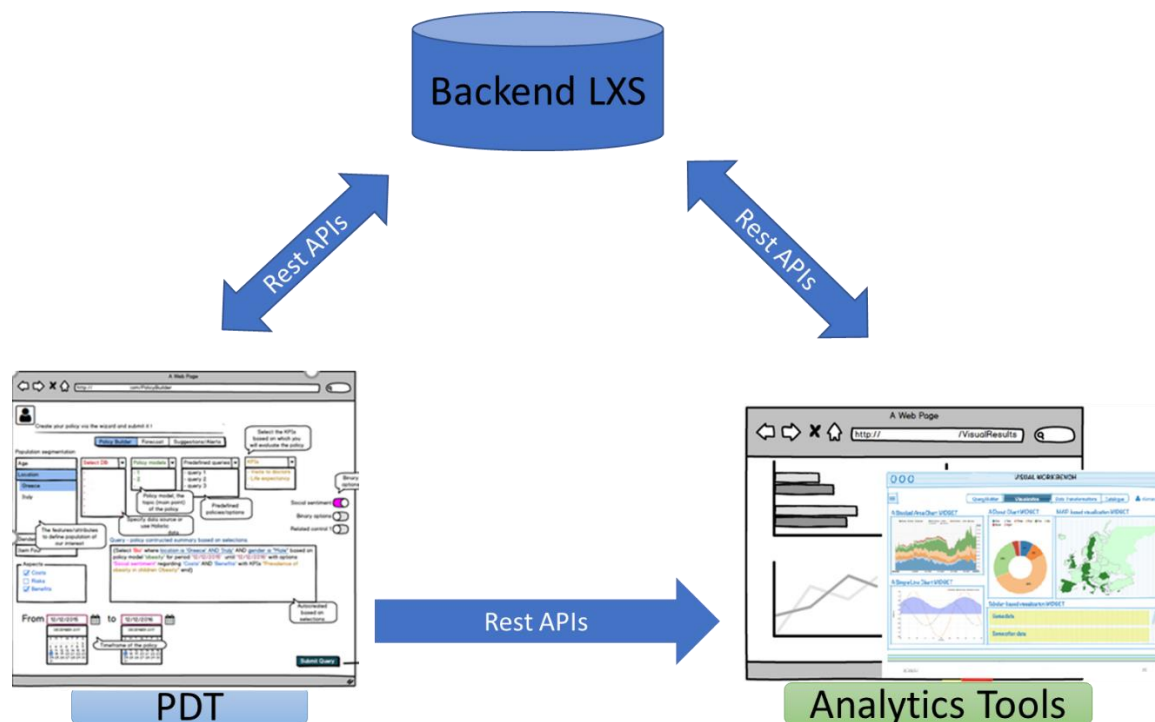


FIGURE 11 – POLICY DEVELOPMENT TOOLKIT COMMUNICATION COMPONENTS

Figure 11 shows the two main components with which PDT will communicate: Backend / Data Repository and the various Analytics Tools.

Both components will expose API Interfaces so that PDT - as the front-end UI - receives the policy model related data from datastore along with the list of registered policy-related data sources and analytic functions. It then activates the selected analytic function on a predefined data source with the parameters specified by the PM. The arrows in Figure 11 depict the communication between the components through REST APIs. The Analytics Tools become available to the PDT once they are registered to the platform. The Analytics Tools registration sequence is provided in Section 7.6

The Policies will be serialized in a predefined format following common syntax (in JSON) into the datastore. The PDT will translate/deserialize the policy objects retrieved from the datastore into UI objects to provide the visual environment for the policymaker actions.

The arrow between PDT and datastore also encompasses the process of semantic or rule-based reasoning and querying. Based on the process set out in T5.2, the semantic processing of emerging policies for lifecycle policy modeling is intervened, which enables the validation of the policy structure in terms of their proper construction. They also guide policymakers to choose KPIs, avoid dysfunctional policies, and provide cross-sectional policy optimization information.

In the architecture proposed in Figure 10 each component is decoupled from the others. The modular structure allows versatility and extensibility, regarding analytics tools providers, analytics frameworks, cloud providers and deployment patterns. The -also- modular UI intentionally hides the big complexity for the users, as each component is decoupled and focused on their properties and functions. So, a Policy Model is composed and supported by related KPIs, which in turn are composed of related Analytics Tools that provide their visualization graphs. The Service-Oriented Architecture (SOA) pattern is followed by requiring the components to adhere to a common communication protocol, and by exposing consistent RESTful APIs.

7.8.3 PDT Architecture Integration

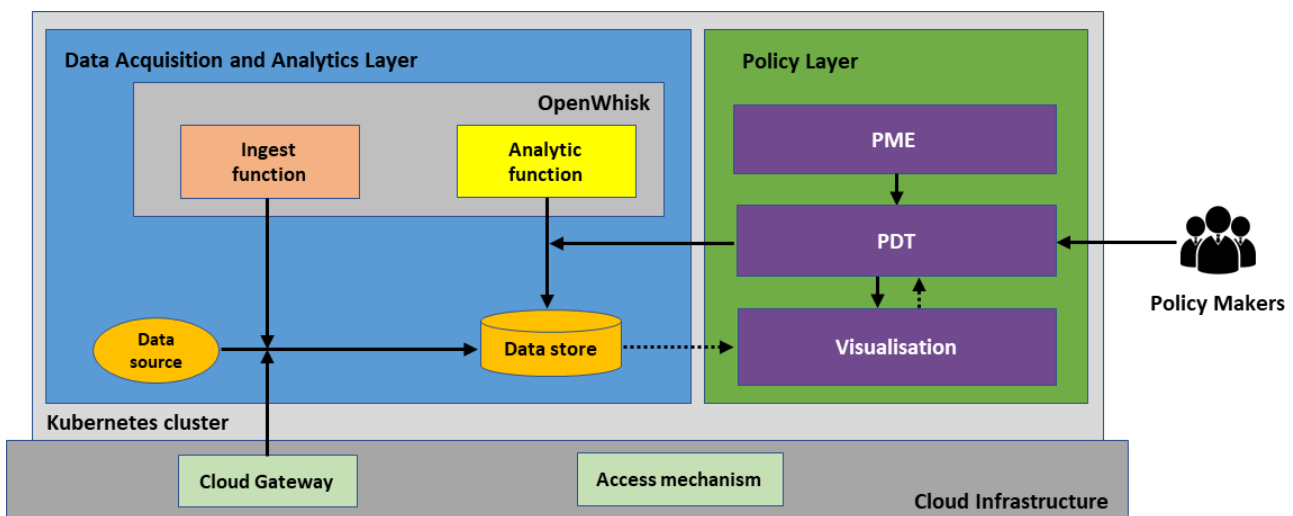


FIGURE 12 – POLICY MODEL DEVELOPMENT INTEGRATION SCHEMA

PME and PDT act in unison as the User Interface for the creation and evaluation of Policy Models (Figure 12). PME guides the user into the creation of proper policy models, while through PDT the user can call for the evaluation of policy KPIs. The analytics results are shown into the same UI depicting KPI values calculation / trends by the integrated visualization component. PME, PDT and the Visualization component share the same source code base, run as a Single Page Application (SPA), hosted on the same Web Server under the same Virtual Machine in the Cloud. All the three integrated components communicate through calls to the Rest APIs offered from the other PolicyCLOUD platform subsystems: PDT Backend, DataStore (WP4), Analytics Tools (WP4) and KeyCloak User Authentication (WP3).

7.9 Layer 5

7.9.1 Data Marketplace

In the context of PolicyCLOUD, the Data Marketplace is a public, unified and standalone platform with many different APIs, able to store several types of assets (solutions). The offered assets vary depending on the needs of the project's stakeholders. They may derive/result from the separate procedures and mechanisms that are implemented in the PolicyCLOUD platform or in general, may be outcomes of the project (policies, templates, tutorials, and others).

From its architecture perspective, the Data Marketplace is structured around two core services, the back-end and the front-end. Generally, the marketplace supports access to its offerings to both end-users and other services (through the respective interfaces). In this context, the end-users are able to interact with the market platform through the front-end that reflects a user-friendly platform (providing the UI), while other additional services (e.g. project's services, 3rd parties) may interact directly with the back-end. This separation contributes towards the platform's enhancements in terms of functionality as well as provides additional information and capabilities.

The back-end side of the marketplace is a RESTful API and receives HTTP requests that trigger the platform's implemented functionalities. As depicted in the following figure (Figure 13), the back-end includes three layers (i.e. Assets Storage Layer, Assets Management Layer, Interaction Layer), while the front-end includes the fourth layer of the Data Marketplace (i.e. the Presentation Layer) that in full consists of all of these 4 different layers. Their capabilities are shortly described below: The "Assets Storage Layer" is the layer in which the platform's offered assets are stored.

- The "Assets Management Layer" delivers all the needed principles and techniques for the management of the marketplace's assets.
- The "Interaction Layer" supports the communication between the marketplace and the end-users (i.e. human end-users, machine end-users), by providing discrete APIs for exploiting each different type of asset.
- The "Presentation Layer" (i.e. the front-end) provides the User Interface (UI) towards the different types of end-users that are willing to use the platform.

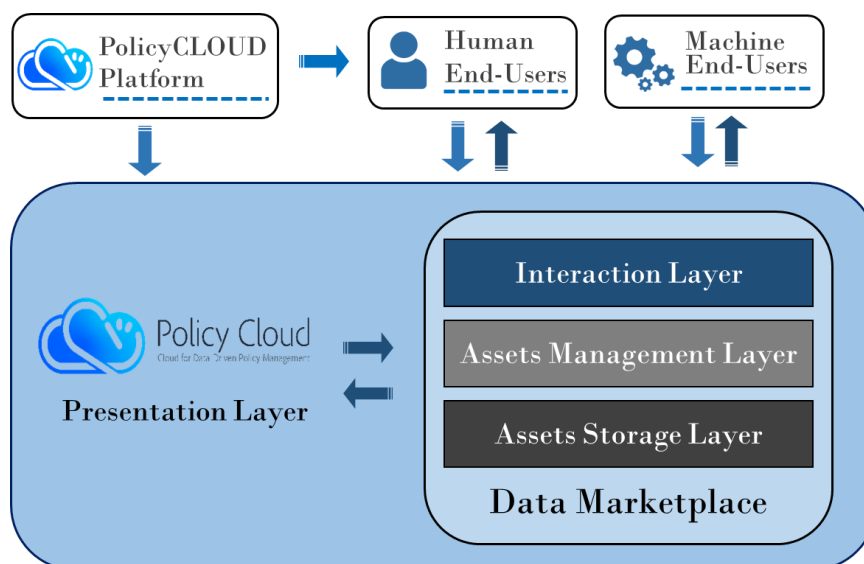


FIGURE 13 – DATA MARKETPLACE ARCHITECTURE

7.10 Data Governance Model, Protection and Privacy Enforcement

Components: Access Mechanisms (T3.6)

7.10.1 Data Governance Model, Protection and Privacy Enforcement

The data governance model and the tools for protection and privacy enforcement are used to protect data and ensure decisions across the complete path following specific guidelines and legislations. Data Governance Model and Privacy Enforcement mechanism is depicted vertically in the right part of the Overall Architecture in Figure 4. This includes three different parts, a) the access policy editor, b) the model and model editor and c) the ABAC authorization engine. The access policy editor will provide the user with the ability to define and store policies based on the ABAC scheme according to the XACML standard. The data governance model of PolicyCLOUD will be used for the definition of these policies, and also for the actual enforcement of the policies by the authorization engine that will be able to evaluate the policies and the attributes, thus enforcing protection and privacy-preserving policies.

In addition, as depicted in Figure 4 and presented for convenience in Figures 14 (A), (B) and (C), the components developed in the scope of T3.6 regarding the protection of data and privacy enforcement, will be used in three separate parts of the overall architecture envisioned for the PolicyCLOUD. The first - Figure 14 (A) - is to provide an access control mechanism for the inclusion and usage of data sources that are being part of PolicyCLOUD. The second - Figure 14 (B) - is the access control being also provided at the level of data visualization, thus allowing or denying access to specific data analytics. The third - Figure 14 (C) - is the usage of the access control mechanisms for managing the control between the PolicyCLOUD datastore and any additional private data store that may be used.

Finally, for the whole mechanism to work properly it has to be mentioned that the authorization engine will need to have access to the attribute values regarding the data, the data sources/origins, the phase of the data lifecycle (e.g. stored data or analysed data) and the phase of the policy lifecycle (e.g. modelling or experimentation process); these can be provided by external components acting as adapters, and can be developed per use case.

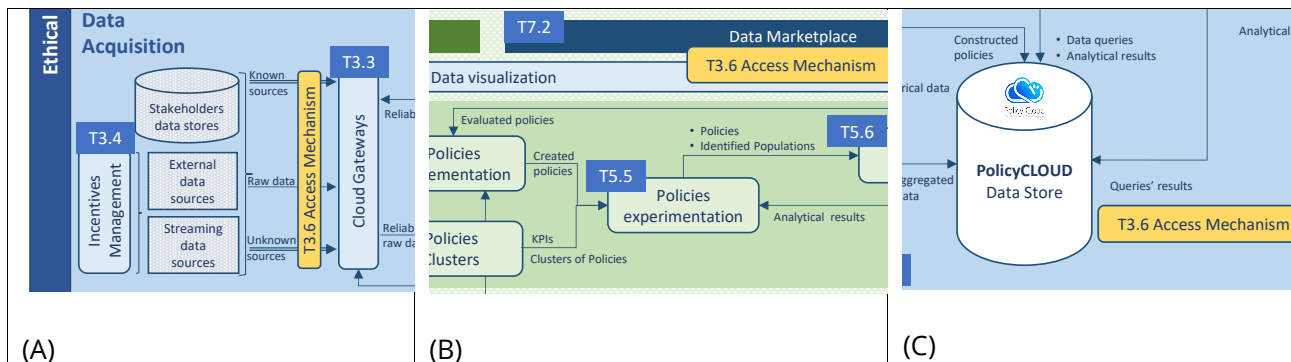
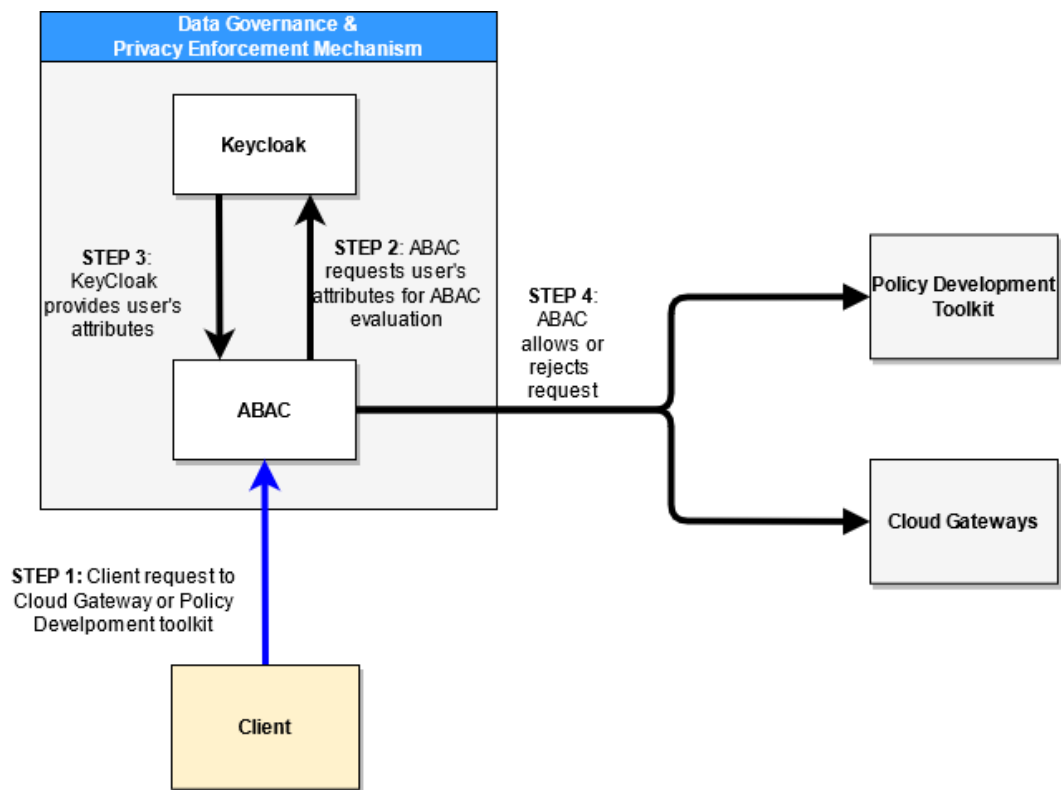


FIGURE 14 – DATA GOVERNANCE MODEL, PROTECTION AND PRIVACY ENFORCEMENT MECHANISMS – EXTRACTED VIEWS (A), (B) AND (C) FROM THE DIAGRAM OF POLICYCLOUD OVERALL ARCHITECTURE.

7.10.2 Data Governance model, protection and privacy enforcement mechanisms Integration

7.10.2.1 ARCHITECTURE INTEGRATION

The integration of the Data Governance and Privacy Enforcement Mechanism is achieved via two components, the Keycloak and the ABAC servers that are connected to each other. As presented in Figure 15, the pair of Keycloak and ABAC can intercept requests to both the Policy Development Toolkit and the Cloud Gateways, ensuring the privacy enforcement for both. A client makes a request to either of those services (STEP 1) and is immediately intercepted by the ABAC Engine. In order to make a decision, ABAC queries the Keycloak server regarding the attributes of the user making the request (STEP 2). Keycloak provides the requested user attributes (STEP 3) and depending on whether they satisfy the current implemented ABAC policy the request is allowed or denied to go through to its original destination (STEP 4).



ABAC Request Interception Example

FIGURE 15 - DATA GOVERNANCE & PRIVACY ENFORCEMENT MECHANISM INTEGRATION FLOW

7.10.2.2 INTEGRATION WITH THE KUBERNETES CLUSTER

The components of the architecture described in the previous sections have been deployed through EGI provisioned infrastructure and have thus been integrated to the Kubernetes cluster. More specifically, an instance of the Keycloak server, along with a connected instance of the ABAC Server have been made available.

8 Use Case examples for end-to-end data path analysis

The first two scenarios from two different Use Cases are summarized based on the information included in Deliverable D6.10 “Use Case Scenarios Definition & Design”, in order to serve as end-to-end examples, demonstrating the data ingest flow and data exploitation while analysing the processing and data transformations along the complete data path.

In the following sections a short description of the Use Case scenarios which includes problem statement, main objectives, Key Performance Indicators and data sources to be used, is provided (sections 8.1 and 8.2).

From the technical perspective an end-to-end data path analysis is provided through the integration of two subpaths: (i) the path from the External Data Sources-Cloud Gateway to the LXS database, which is analyzed in section 8.3 and (ii) the path from the LXS database-PDT backend to the Visualization-PDT interface with which the Policy Maker (as end-user) interacts, analyzed in section 8.4. The two subpaths constitute a complete end-to-end data path from an external data source to a semantically meaningful result to be presented to the end user.

8.1 Use Case 1: Participatory Policies Against Radicalization

8.1.1 Scenario A: Radicalization incidents

Description:

Monitor the occurrence of radicalization incidents in the geographic proximity of a region. Data coming from the GTD and RAND will be used. The Policy Maker can select the area of his/her interest and consult the different incidents that have taken place in a given period.

8.1.2 Main Objective

Validate existing policies and investigate if there is a need to update them or create new ones based on the retrieved information.

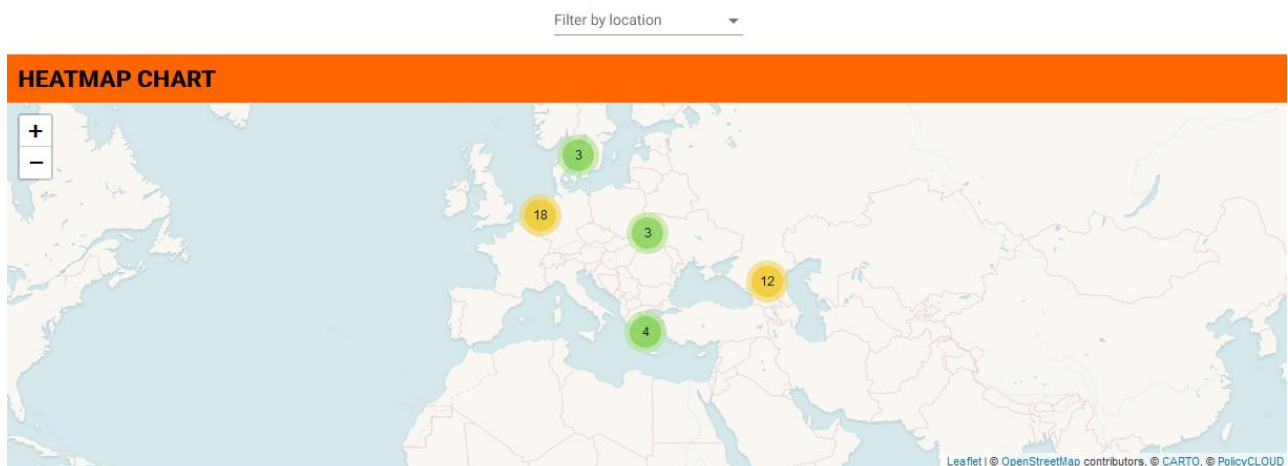


FIGURE 16 - VISUALIZATION ON POLICYCLOUD OF THE RESULT OF SCENARIO A: RADICALIZATION INCIDENTS OF USE CASE 1

8.1.3 Key Performance Indicators

Section	Description
ID	MAG-KPI9
Title	Number of identified occurrences of radicalization incidents in a given area
Priority	High
Reference Use Case	UC#1
Success Criteria	≥ 0

TABLE 1 – UC1 BUSINESS KPI9

8.1.4 Data Sources

Use Case	Scenario #	Data Source Description	Link(s)
Participatory Policies Against Radicalization	Scenario A	Managed by the National Consortium for the Study of Terrorism and Responses to Terrorism (START), the Global Terrorism Database includes more than 200,000 terrorist attacks dating back to 1970.	https://www.start.umd.edu/gtd/access/

TABLE 2 – DATA SOURCES LIST FOR SCENARIO A OF THE PARTICIPATORY POLICIES AGAINST RADICALIZATION USE CASE

8.2 Use Case 2: Intelligent policies for the development of the agrifood industry

8.2.1 Scenario B: Visualization of negative and positive opinions on social networks for different products

Description:

Visualize the negative and positive opinions on social networks of the different products analysed allowing an automatic and immediate response to the end user.

8.2.2 Main Objective

Create an immediate communication with the end user, knowing their impressions, both positive and negative, that will allow us to interact with the end customer more directly.

8.2.3 Key Performance Indicators

Section	Description
ID	SAR-KPI4
Title	Provide real-time calculation capacity
Priority	High
Reference Use Case	UC#2
Success Criteria	>20% of the data

TABLE 3 – UC2 TECHNICAL KPI4

Section	Description
ID	SAR-KPI6
Title	Increase process speed
Priority	High
Reference Use Case	UC#2
Success Criteria	>30% Reduce time

TABLE 4 – UC2 TECHNICAL KPI6

Section	Description
ID	SAR-KPI8
Title	Total number occurrences
Priority	High
Reference Use Case	UC#2
Success Criteria	>50%

TABLE 5 – UC2 BUSINESS KPI8

Section	Description
ID	SAR-KPI9
Title	Relative Total n° occurrences %
Priority	High
Reference Use Case	UC#2
Success Criteria	>10%

TABLE 6 – UC2 BUSINESS KPI9

Section	Description
ID	SAR-KPI10
Title	Opinion (-1 (negative) to 1 (positive)) impact
Priority	High
Reference Use Case	UC#2
Success Criteria	Average positive

TABLE 7 – UC2 BUSINESS KPI10

Section	Description
ID	SAR-KPI11
Title	Increment of the impact in the last month
Priority	High
Reference Use Case	UC#2
Success Criteria	>15%

TABLE 8 – UC2 BUSINESS KPI11

8.2.4 Data Sources

Link #	Link
1	https://opendata.aragon.es/datos/catalogo?texto=pac
2	https://www.aragon.es/en/-/vitivinicultura.-registro-viticola
3	https://www.aragon.es/en/temas/medio-rural-agricultura-ganaderia/agricultura/vinedos-vinos-bebidas-alcoholicas
4	https://opendata.aragon.es/datos/catalogo/busqueda/siu?tema=vinedos-vinos-bebidas-alcoholicas
5	https://opendata.aragon.es/servicios/open-social-data/#/main

TABLE 9 – LINKS TO ARAGON USE CASE DATA STORES

8.3 Data Path Analysis (From Cloud Gateways to LXS Database) based on the implemented Use Case scenarios

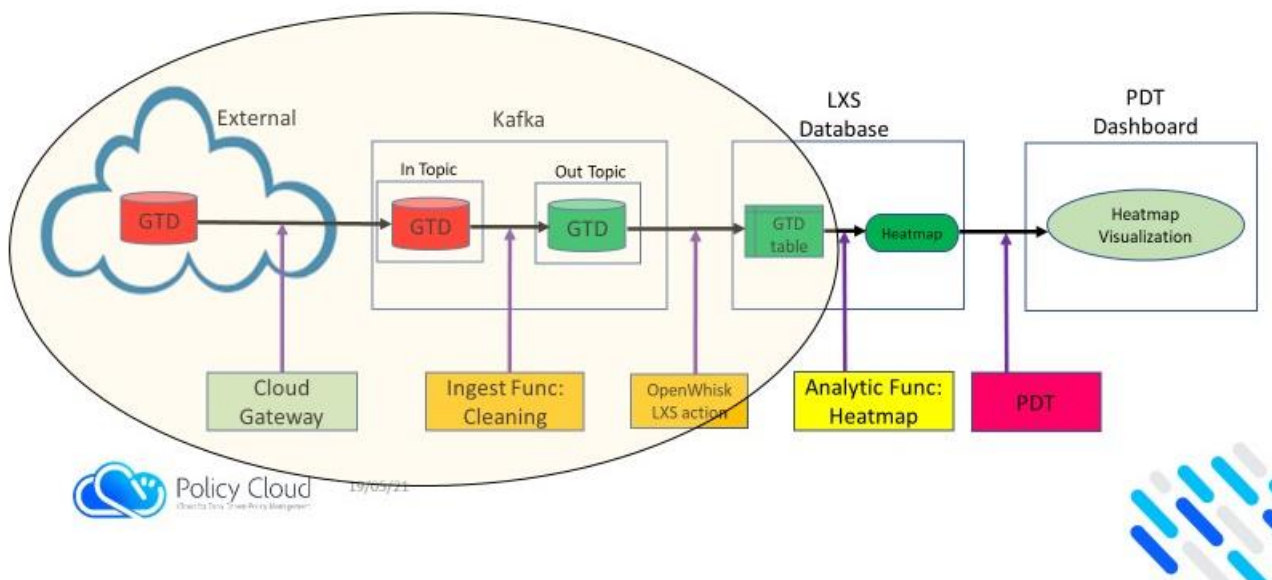


FIGURE 17 - DATA PATH ANALYSIS

The first scenario (scenario A) for Use Case 1 “Participatory policies against Radicalization” provides to a Policy Maker a visualization of a heatmap showing the frequency of occurrence of radicalization incidents in the geographic proximity of a region. Data coming from the GTD is used. Figure 17 demonstrates the data path for this end-to-end example, for the function implementing the heat map computation which is invoked when the heatmap visualisation is called.

The same architecture shown in Figure 17, is also used for the first scenario (scenario A) of Use Case 2 “Intelligent policies for the development of agrifood industry” which provides to a Policy Maker a visualization of the ARAGON wine sentiments with data received from Social Media.

It could be remarked that the SocialMedia (ARAGON wine sentiments) results can also be presented in the same manner as the results from the GTD.

8.4 Data Path Analysis (from the LXS database backend to visualization of result)

Data Path Analysis highlights the integration among the various components that consist the PDT on the one hand, and the corresponding building blocks of the overall PolicyCLOUD architecture on the other hand, Figure 18 provides the sequence diagram of all interactions that take place when a PM invokes an analytical function and receives the results in a visualized graph.

As depicted in the sequence diagram, the end user of the PDT, the Policy Maker, retrieves all existing policies stored in the platform, according to some filter criteria. The PDT returns these policies and visualizes them in its graphical user interface. Then, the Policy Maker wants to verify a specific policy, by making an analysis over the available data. Using the graphical user interface of the PDT, it selects and clicks on the KPI that the Policy Maker that wants to verify. The GUI invokes the corresponding REST web method of the backend to execute the relevant analytical function.

Subsequently, the backend **interacts with the data acquisition and analytics layer** of the PolicyCLOUD (sections 7.6.11.1 and 8.3). It contains all the information regarding registered analytical functions, their required input parameters, the type of their output etc. As a result, it collects all data received by the invocation of its REST web method, and further requests from the DAA layer to execute the function.

As explained in the previous subsection, the DAA layer incorporates the OpenWhisk serverless platform, and relies on the latter to deploy the requested function. Openwhisk takes the responsibility to create, via Kubernetes, the corresponding infrastructure resources, deploy its runtime execution environment there and finally execute the requested function.

The function on the other hand, receives its required input parameters that were passed to it via the PDT, the backend and the DAA layer, along with other meta-information (such as the connection url of the data management layer to retrieve data). Through this process it receives all required information to open a database connection with the datastore, and execute its relevant query, thus pushing a pre-processing down to the storage layer, in order to retrieve only the amount of data that is needed to run its AI algorithm.

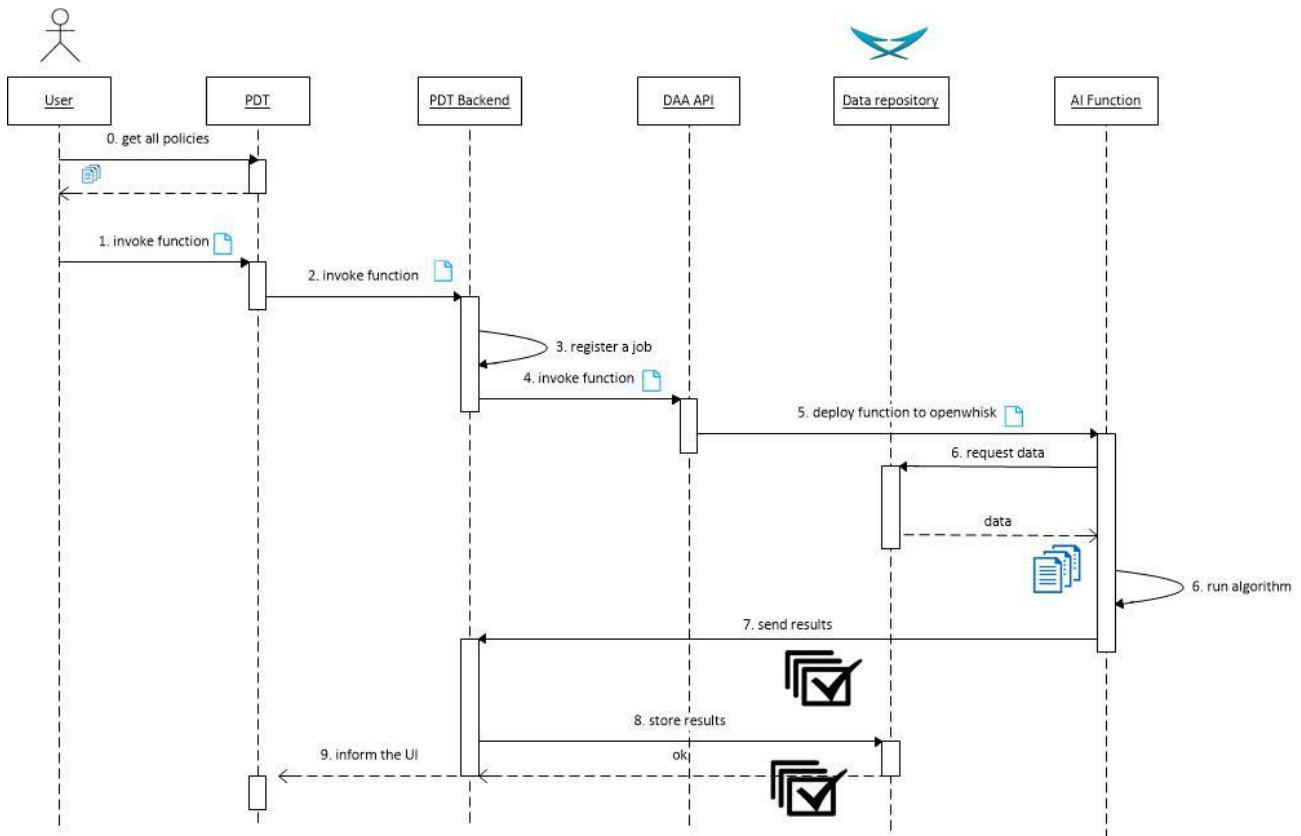


FIGURE 18 - SEQUENCE DIAGRAM FOR PDT-DAA INTERACTION

Among other meta-information parameters received, one important parameter is the URL of a REST web service that the function needs to communicate to persistently store the results of the analysis. In fact, when the AI algorithm produces results, and before the function completes, it firstly sends the results in this URL provided by the PDT backend, so that the latter can store this information and make it available to the end-user, the Policy Maker. After sending the results, the function returns, so that the Openwhisk can shut down the relevant run-time environment and release the resources used. When the function has been properly shut down, it informs the DAA layer, and the result can be further communicated to the PDT backend. Finally, when the REST web service of the backend is invoked in order to store the results, the backend persistently stores them into the data repository, by adding relevant meta-information. At the end of this process, it informs the PDT via web sockets that the results are now available and can be retrieved. The PDT sends a pop-up notification in the graphical user interface, and the Policy Maker can now click and see the results. The PDT will get the results upon request from the backend, and it will activate and make visible the corresponding type of visualization graph to show these results.

Our design allows all involved components to be generic enough and not locked into a specific deployment or implementation. For instance, the analytical functions do not need to know in advance where to connect to retrieve data, or what is the schema of the underlying data, or where to store results. They can be generic and receive this information during runtime. In the same sense, the DAA layer has

www.policycloud.eu

not been developed specifically for the needs of the PolicyCLOUD platform, rather than it is generic and can be used by any application that needs to access programmatically a serverless platform in order to administrate, deploy and execute functions over this environment. In the same direction, the PDT backend provides an interface for different connectors. An implementation of these connectors has been developed to allow the integration of the PDT with the DAA layer. In that sense, the PDT is not locked-in to a specific platform, but it can use any other type of environment by implementing the relevant connector.

Since the involved components are generic, individual components can increase their sustainability, by being exploitable in other deployments or other integrated solutions. The analytical provider does not need to implement its function specifically for the PolicyCLOUD environment only, rather can he or she only focus on the AI algorithm, which can be used in different solutions. In the same manner, the DAA layer does not provide capabilities specifically for the PolicyCLOUD environment only, but it can be exploited in any integrated solution that requires to programmatically administrate the serverless platform. Finally, the PDT can be deployed in other environments that might not allow for dynamic deployments. However, by putting all these layers together integrated into the PolicyCLOUD, the platform can benefit from all the advancements provided by the individual components, and this is what makes the overall integrated platform so innovative.

9 Conclusion

The second version of the PolicyCLOUD Conceptual Model & Reference Architecture (originally submitted as Deliverable D2.2) is presented in this document. The final update of the deliverable will be published in M30.

The architecture consists of the following five layers: Cloud Based Environment (Layer 1a), Data Management – Data Stores (Layer 1b), Data Acquisition and Analytics (Layer 2), Policies Management Framework (Layer 3), Policy Development Toolkit (Layer 4) and Data Marketplace (Layer 5). The architecture also includes the Ethical Framework and the Data Governance Model, Protection and Privacy Enforcement.

Within the updates of this document **special emphasis has been given on the Integration in PolicyCLOUD** which follows three directions: (i) architecture integration, (ii) integration with the cloud infrastructure and (iii) integration with Use Case scenarios through the implementation of end-to-end scenarios.

Two Use Case scenarios are used for end-to-end data path analysis: (i) the scenario of Use Case 1: “Radicalization incidents” and the scenario of Use Case 2: “Visualization of negative and positive opinions on social networks for different products”. The data path analysis consists of two subpaths: (i) the subpath from the Cloud Gateways to LXS database and (ii) the subpath from the LXS database backend to the visualization of result.

Additional integration activities took place along the two frameworks of PolicyCLOUD, (a) the Data Governance model, protection and privacy enforcement mechanism and (b) the Ethical and Legal Compliance framework.

References

- [1] Apache Mesos, "<http://mesos.apache.org/>"
- [2] IBM Cloud Object Storage, <https://www.ibm.com/cloud/object-storage>
- [3] PolicyCLOUD. D3.1 - Cloud Infrastructure Incentives Management and Data Governance Design and Open Specification 1. Ledakis, Giannis. 2020.
- [4] Load balancing, Molecular "<https://molecular.services/docs/0.14/balancing.html>"
- [5] Fault tolerance, Molecular "<https://molecular.services/docs/0.14/fault-tolerance.html>"
- [6] J. Kreps, N. Narkhede, and J. Rao, "Kafka: A distributed messaging system for log processing", NetDB, 2011.
- [7] Malone, Thomas W., Robert Laubacher, and Chrysanthos Dellarocas. "The collective intelligence genome." MIT Sloan Management Review 51, no. 3 (2010): 21.
- [8] PolicyCLOUD. D3.4 - Cloud Infrastructure Incentives Management and Data Governance Design and Open Specification 2. Ledakis, Giannis. 2021.
- [9] PolicyCLOUD, D3.3 PolicyCLOUD's Societal and Ethical Requirements & Guidelines, Audino Alice, 2020.
- [10] Gehra B., Leiendecker J. & Lienke G. (2017), White Paper. Compliance by Design: Banking's Unmissable Opportunity, "https://image-src.bcg.com/Images/Compliance-by-Design-Dec2017_tcm9-198779.pdf" retrieved 2020-11-28.
- [11] Balboni P., Taborda Barata M., Botsi A. & Francis K. (2019), Accountability and Enforcement Aspects of the EU General Data Protection Regulation – Methodology for the Creation of an Effective Compliance Framework and a Review of Recent Case Law, Indian Journal of Law and Technology, 15(1), 102-259.
- [12] JSON, "<http://json-ld.org/>"
- [13] Yamada, I., & Shindo, H. (2019). Neural attentive bag-of-entities model for text classification. *arXiv preprint arXiv:1909.01259*.
- [14] Attardi, G., Buzzelli, A., & Sartiano, D. (2013). Machine Translation for Entity Recognition across Languages in Biomedical Documents. In *CLEF (Working Notes)*.
- [15] Olszewski, Robert. (2001). Generalized feature extraction for structural pattern recognition in time-series data.