

## Data Pitch

H2020-ICT-2016-1

**Project number: 732506**

### **D3.6 Data catalogue V2**

**Coordinators:** Jérémy Decis and Anthoine Dusselier (Dawex)

**Contributor:** Stefano Modafferi (university of Southampton)

**Quality reviewer:** Elena Simperl (University of Southampton)

Deliverable nature:	Other
Dissemination level: (Confidentiality)	Public
Work package	WP3
Contractual delivery date:	30 June 2018
Actual delivery date:	29/06/2018
Version:	V2.0
Keywords:	Data catalogue, data-sharing platform

# Table of Contents

<b>Abstract</b>	<b>3</b>
<b>2. Executive summary</b>	<b>4</b>
<b>3. Introduction</b>	<b>5</b>
4.2 Data catalogue for the 2018 call	7
4.2.1 Stage 1: Contact with providers and contract signing	7
During Stage 1 of the project, the consortium will present the functionalities of the platform to the data providers and onboard them once the contracts are signed.	7
4.2.2 Stage 2: Challenge definition	7
4.2.3 Stage 3: call opens and applications	9
4.2.4 Stage 4: Negotiations	11
4.2.5 Stage 5: accelerator	11
<b>5. Conclusion</b>	<b>13</b>
<b>Tables</b>	
Table 1: Metadata description template	7
Table 2: Catalogue function by stage	17
<b>Figures</b>	
Figure 1: Screenshots of the current version of data catalogue	9
Figure 2: Screenshot of current version of the data catalogue (provider page)	10
Figure 3: Screenshots of the future data catalogue	15

## **1. ABSTRACT**

The DNA of Data Pitch is to allow experimenters and providers to extract value from closed data. This document describes the how the Data Pitch catalogue is utilised at each of the five stages of the process, from initial discussions with data providers to the creation of innovation data solutions.

## **2. EXECUTIVE SUMMARY**

This deliverable describes the part of the Data Pitch sharing platform related to the data catalogue, and the solution chosen by the Consortium for the second call.

This solution is an ad-hoc version of the existing Dawex Marketplace platform. This is the main technical asset brought by Dawex to the consortium and it includes lots of features and functionality that are required by Data Pitch.

### 3. INTRODUCTION

Data Pitch provides a bridge for experimenters and data providers to identify new instruments for extracting value from private datasets. The support the consortium offers to the applicants is global, covering all the legal, technical and business model aspects.

The Data Pitch goal is fulfilled via a process divided in stages.

Stage 1 of the process includes the discussions with data providers and the signing of all the legal contracts for including their data in one of the Data Pitch challenges. Stage 1 is concluded by the publication of relevant information and possibly a sample of datasets related to an identified challenge (i.e. the creation of a data catalogue).

Stage 2 includes the publication of the challenges for the benefit of the applicants and the collection of their applications.

Stage 3 includes the evaluation of the received proposals.

Stage 4 includes the negotiation phase where the fine detail of the experiment is defined in conjunction with the successful applicant(s) and the data provider.

Stage 5 is the actual experimentation period where data provider and the successful applicant work together to solve a relevant business problem.

Within the general expression of 'data sharing platform', several instruments are identified and available to support the different stages:

- The legal toolkit. The purpose of the legal toolkit is to make the data providers aware of the implication of the legal aspects of the data sharing (Cf. Del. D3.1). This instrument is relevant in Stage 1 and Stage 5.
- The application manager. This is an off the shelf service called F6S that is well-known in the startup world and provides a means for managing programme application and review. This instrument is relevant in Stage 2<sup>1</sup> and Stage 3<sup>2</sup> (<https://www.f6s.com/datapitchaccelerator/apply>).
- The data catalogue. This instrument is an ad-hoc version of the Dawex marketplace and it is relevant in Stage 1 to publish the dataset description, Stage 2 to provide access for the applicants to the dataset description, Stage 4 when contracts will be again one of the main topics and possibly Stage 5 when the access to the dataset is made via API.
- The experimental platform (cf D2.1). This instrument is used when either the data providers or the successful applicant requires an infrastructural support. It is relevant in Stage 5.

This deliverable focuses on the presentation of the data catalogue instrument.

Within the document, the acceleration stage is also called the experimentation stage and the successful applicants are also referred to as challenge winners or experimenters.

---

<sup>1</sup> See page 9

<sup>2</sup> See page 11

## 4.2 DATA CATALOGUE FOR THE 2018 CALL

As explained in the introduction the data catalogue will be used in Stage 1 to publish the dataset description, in Stage 2 to provide access for the applicants to the dataset description, in Stage 4 to onboard and vet the experimenters on the platform, and possibly in Stage 5 where the Dawex platform will support the access to data for the experimentation when the data are streamed or accessed via API. In the following paragraphs the catalogue functionalities and the Data Pitch requirements will be presented according to the different stages.

### 4.2.1 Stage 1: Contact with providers and contract signing

During Stage 1 of the project, the consortium will present the functionalities of the platform to the data providers and onboard them once the contracts are signed.

The early outcome of this stage is the definition of the challenge and the signing of the contract.

### 4.2.2 Stage 2: Challenge definition

Stage 2 includes the publication of the information related to the datasets. Ideally this step should be done by the data provider, who registers and joins the catalogue autonomously. To speed up the process an “acting on behalf” function will allow selected members of the consortium to input the data themselves if they are received in a different way. In the case where the data providers give us a link to an existing documentation, a member of the consortium will integrate it manually in the catalogue.

Once registered, each provider has the choice to create an “offering” or a “theme”, based on their data sets and depending on the level of sensitivity of their data.

Metadata both from offerings and thematics will thereafter be used in the public catalogue to describe the data provided.

#### **4.2.2.1 Offering:**

An “offering” is for a data provider who already has the data and is willing to upload it on the platform. By creating an offering, a data provider will enable the Data Pitch platform to create a sample. This sample will be generated with random algorithms. The sample will be available on the public catalogue during the Stage 3.

This feature will be used when the data are not sensitive. The offering supports data by APIs

#### **4.2.2.2 Metadata collected for an offering for data files (csv, json, pdf, txt, xml, xls, geojson):**

- Title of the data set
- Detailed description of the data set
  - o free text description of the data set
  - o free text description of each header present in the data
- Territories covered by the data
  - Business sectors covered by the data
  - Type(s) of data available
  - The data produced contains...

- o Anonymized data derived from personal data
- o Personal data
- o No personal data
- Keywords that describe the data set
- History of your data offering
- Level of data structure for this theme
  - o Not structured (never applied the resources to structure it)
  - o Not structured (by nature)
  - o Partially structured
  - o Highly structured
  - o I don't know how to evaluate it

#### 4.2.2.3 Metadata collected from Data provided via API's<sup>3</sup>:

One of the most challenging and also interesting case is the execution of analytic functions on data streams. The platform has the ability to manage the streaming live of data. The platform works as a proxy : it will play the role of a bridge between the data provider, and the data user, who wants to acquire the data. In the case of Data Pitch, the source will be the data provider and the destination will be the experimenter's infrastructure.

Using this functionality of the Dawex marketplace, we will ensure the consistency of the data and allows for API from IoT, Business Intelligence software, DPM and app.

Name of the data set

- Summary (250 characters)
- Data (The types of data making up your offering. Several choices are possible.)
  - Financial and administrative data
  - Sales and marketing data
  - Production and technical data
  - Etc.
- Indicate whether or not personal data is present
  - Anonymized data derived from personal data
  - Personal data
  - No personal data

Business sector (The business sectors covered by the data of your offering. )

- Geographical areas covered
- Production period (A date or period during which data is produced.)
- Description of the API
  - New URL (Fill in all the information needed to describe the URL for your API.)
    - Name
    - URL
    - Request (GET/POST)
- Output data format

---

<sup>3</sup> <https://www.dawex.com/en/product/#api>

- Error management protocol
- Does this URL contain parameters? (Yes/No)
- Details of your offering
  - Description – Keywords

#### 4.2.2.4 Theme:

If the data provider does not have the data ready, or if they do not want to upload it on the platform, they will create a “theme”. A “theme” is used only to describe the data that the provider will provide later to the experimenters. It is not possible to create a sample in this case. The metadata collected during the creation process will be used in the public catalogue during the Stage 3.

This feature will be used when the data are sensitive.

#### 4.2.2.5 Metadata collected for a theme

- Title of the theme
- Detailed description of the theme and of the data (This includes the free text description of each header present in the data)
- Territories covered by the data
- Business sectors covered by this theme’s data
- Type(s) of data available for this theme
- The data produced for this theme contains
  - Anonymized data derived from personal data
  - Personal data
  - No personal data
- Keywords that describe this theme
- History of the theme’s data
- Level of data structure for this theme
  - Not structured (never applied the resources to structure it)
  - Not structured (by nature)
  - Partially structured
  - Highly structured
  - I don’t how to evaluate it

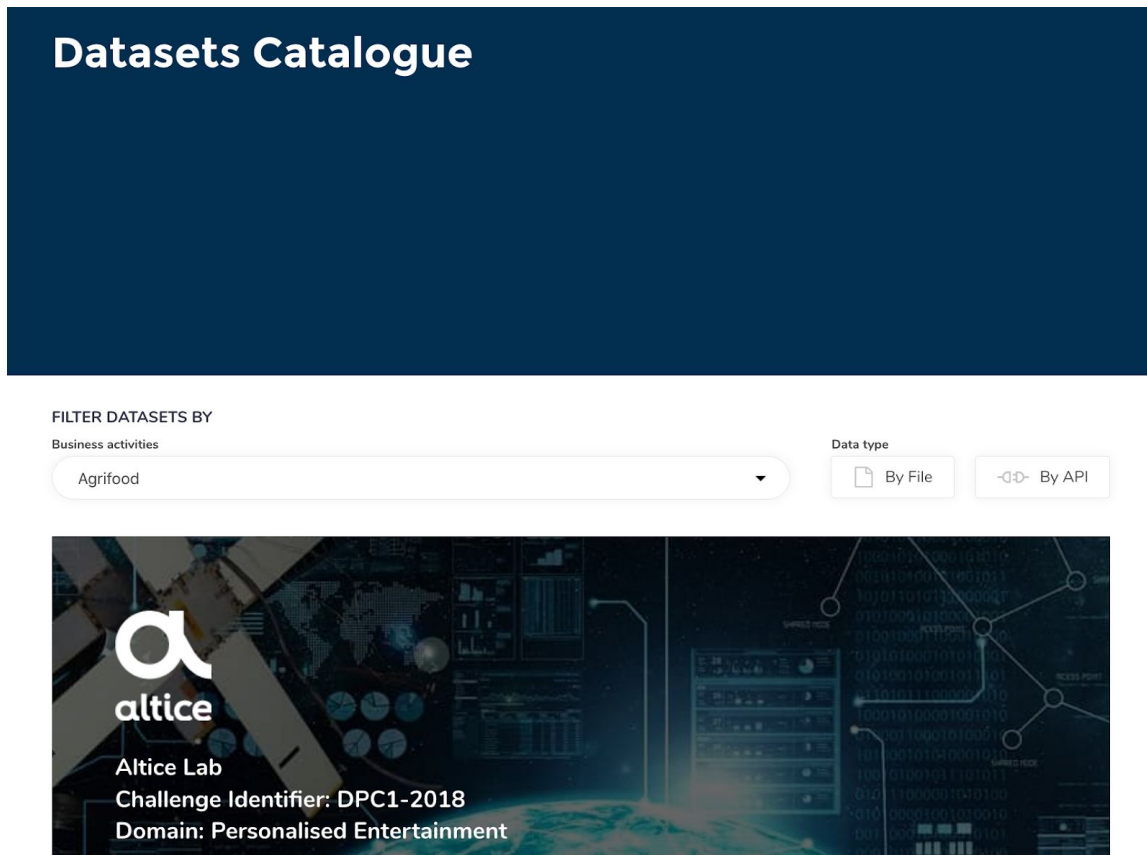
### 4.2.3 Stage 3: call opens and applications

After the call opens, a public catalogue has to be available to the applicants, to our associates and partners. This catalogue do not include the data itself, but either a sample generated from the data and a description of the data set, or only a description, allowing the applicants to build innovative solutions and to submit their proposals.

For these purposes, a landing page is created, gathering the samples generated from the offerings created on the platform, and the metadata from the themes. The catalogue will have a search and filter feature, allowing applicants to easily navigate through the datasets descriptions.

The data catalogue can be viewed here : <https://www.dawex.com/en/lp/datapitch-datasets/2018/>



**Figure 3: screenshots of the data catalogue**

#### 4.2.4 Stage 4: Negotiations<sup>4</sup>

After the end of the call, the applicant proposals will be evaluated by Data Pitch and the winners will be selected for the acceleration stage.

There will be 3 steps during this stage:

1. On boarding the experimenters on its platform.
2. Vetting the experimenters and due diligence on the status of the company<sup>5</sup> to enhance the security level of the platform and to reassure the data providers. This is a critical point in order for data providers agree to join Data Pitch.
3. When the two processes are finalized, the experimenters have access to the Datapitch platform. Which part of the Data Pitch sharing platform is going to be involved in the experimental phase is going to be part of the negotiation and will result from the analysis of the experiment requirements.

<sup>4</sup> Stage 3 is about the evaluation of the proposal and the data catalogue and contract repository does not play a relevant role.

<sup>5</sup> [https://drive.google.com/file/d/0B9IIZV\\_CjqlcOTVIYkh6V0xGNE0/view?pli=1](https://drive.google.com/file/d/0B9IIZV_CjqlcOTVIYkh6V0xGNE0/view?pli=1)

As in Stage 1 contracts and legal aspects in general will play a fundamental role in this stage. The repository part of the Data Pitch sharing platform will offer as many instruments as possible to facilitate and smooth the management of all the legal aspects.

### 4.2.5 Stage 5: accelerator

In this stage the data catalogue and document repository play a secondary role. Nevertheless other functionalities from the platform are going to be used in Data Pitch as already explained in Deliverable 2.1.

These functions will be activated when possible and convenient. This is most likely when a data stream is going to be involved or the data providers offer an API for accessing her datasets.

When data are going to be accessed via these functionalities, the dataset will remain completely under the control of the data provider that will have their own credentials to perform the task.

#### **4.2.5.1 API technical details**

The platform acts as an intermediate server (“proxy”) between the experimenter’s server and the provider’s server (which is providing the API)

- The experimenter makes a request to the server with its identifiers and a reference to the desired API
- The server checks if the subscription is valid, as well as the API reference and query volumes restrictions (daily and per second)
- If the request is valid, the server queries the provider’s server to retrieve the data
- The data is then returned to the experimenter via the server.

The server also performs a sampling based on vendor data at regular intervals to ensure that it is always available.

#### **4.2.5.2 Activity reports**

At any stage of the process, the data catalogue and the functionality for data access will regularly produce reports from the activities on the platform of both the data providers and of the experimenters.

These reports will include information about:

- searches on the platform made thanks to the search engine
- the pages and the data that the experimenters consulted
- contact requests
- downloads

These reports will allow to monitor the interactions between the providers and the experimenters, in order to evaluate their behaviours, to better understand their needs, but also to ensure that the rules of the project are respected.

## 5. CONCLUSION

This table summarizes all the functionalities provided by the Data Pitch sharing platform depending the stage of the project.

		Stages				
		Stage1: challenge definition	Stage2: applications	Stage3: Selection process	Stage4: negotiations	Stage5: acceleration of the experimenters
<b>F u n c t i o n a l i t i e s</b>	<b>Offering creation</b>	Collect metadata and generate sample	Provide metadata and sample to the applicants	N/A	N/A	Data access
	<b>Theme creation</b>	Collect metadata	Provide metadata to the applicants	N/A	N/A	N/A
	<b>Data from API</b>	Collect metadata	Provide metadata and description to the applicants	N/A	N/A	Data access
	<b>Sample Production</b>	N/A	Assess data quality and structure	N/A	N/A	N/A
	<b>Landing page</b>	N/A	Gather the metadata and the samples on a single place for the applicants	N/A	N/A	N/A
	<b>Vetting process</b>	N/A	N/A	N/A	Reassure the data providers and enhance the security of the platform.	N/A
	<b>Personalized data management and access</b>	Provide a dedicated environment to the Data providers on the platform.	N/A	N/A	N/A	Provide a dedicated environment both to the data providers and to the experimenters on the platform.
	<b>Activity reports</b>	N/A	Monitor the activity of the applicants on the catalogue	N/A	N/A	Monitor the activity of the data providers and the experimenters.

Table 2: Catalogue function by stage

This document explains the choice made by Data Pitch to provide data to the applicants and, later, to the experiments, by using a data sharing platform allowing to create the data catalogue, as well to share data between the providers and the experimenters in some cases (API transfer).