

D3.1 Design of the CH Cloud and 4CH platform



Deliverable Report n. 3.1: final version, issue date on 28 July 2022

Grant Agreement number:	101004468
Project acronym:	4CH
Project title:	Competence Centre for the Conservation of Cultural Heritage
Funding Scheme:	H2020
Project coordinator:	Francesco Taccetti, INFN
Tel:	+39 3201806514
E-mail:	francesco.taccetti@fi.infn.it
Project website address:	www.4ch-project.eu

Title: Design of the CH Cloud and 4CH platform
Issue Date: 28 July 2022
Produced by: INCEPTION
Main authors: Marco Medici (INCEPTION), Alessandro Costantini (INFN), Franco Niccolucci (PIN, section 7)
Co-authors: Federico Ferrari (INCEPTION), Peter Bonsma (RDF), Alessandro Bombini (INFN), Francesco Giacomini (INFN), Laura Cappelli (INFN), Achille Felicetti (PIN), Maria Theodoridou (FORTH), Sorin Hermon (CYI)
Version: Final
Reviewed by: Francesco Taccetti (all), Franco Niccolucci (sections 1-6, 8)
Approved by: Francesco Taccetti
Quality Control: Paola Ronzino (PIN)
Dissemination: Public

Colophon

Copyright © 2021 by 4CH consortium

Distributed under the CC-BY-NC-SA 4.0 license



Use of any knowledge, information or data contained in this document shall be at the user's sole risk. Neither the 4CH Consortium nor any of its members, their officers, employees or agents accept shall be liable or responsible, in negligence or otherwise, for any loss, damage or expense whatever sustained by any person as a result of the use, in any manner or form, of any knowledge, information or data contained in this document, or due to any inaccuracy, omission or error therein contained. If you notice information in this publication that you believe should be corrected or updated, please contact us. We shall try to remedy the problem.

The authors intended not to use any copyrighted material for the publication or, if not possible, to indicate the copyright of the respective object. The copyright for any material created by the authors is reserved. Any duplication or use of objects such as diagrams, sounds or texts in other electronic or printed publications is not permitted without the author's agreement.

4CH is a Horizon 2020 project funded by the European Commission under Grant Agreement n.101004468 – 4CH.



Document History

- 01.05.2022: Draft version 0.1
- 20.06.2022: Draft version 0.2
- 30.06.2022: Draft version 1.1
- 07.06.2022 Draft version 1.2
- 15.07.2022: Draft version 1.3
- 21.07.2022 Draft version 1.4
- 27.07.2022 Quality Control
- 28.07.2022 Final version

List of acronyms and abbreviations

3M: Mapping Memory Manager
 AEC: Architecture Engineering and Construction
 API: Application Programming Interface
 AR: Augmented Reality
 ASCII: American Standard Code for Information Interchange
 B-Rep: Boundary Representation
 BCF: "BIM Collaboration Format" file format
 BIM: Building Information Modelling
 CAD: Computer-Aided Design
 CAM: Computer-Aided Manufacturing
 CC: Competence Center
 CDE: Common Data Environment
 CGI: Computer-Generated Imagery
 CH: Cultural Heritage
 CHNet: Cultural Heritage Network
 CIDOC-CRM: Center for Intercultural DOCumentation - Conceptual Reference Model
 COBie: "Construction Operations Building Information Exchange" file format
 CPU: Central Processing Unit
 CRMhs: CRM extension for heritage sciences
 CSS: Cascading Style Sheets
 DAE: "Digital Asset Exchange" or "Collada" file format
 DB: DataBase
 DCHE: Digital Cultural Heritage and Europeana
 DoA: Document of Action
 DT: Digital Twin
 DXF: "Design Web Format" file format
 E57: Lidar Point Cloud Data File
 Getty AAT: Getty Art&Architecture Thesaurus
 GIS: Geographic Information Systems
 glB: "GL Transmission Format Binary file" file format
 glTF: "Graphics Language Transmission Format" file format
 GPS: Global Positioning System
 GPU: Graphic Processing Unit
 HBIM: Historical/Heritage Building Information Modeling
 HDF5: Hierarchical Data Format version 5
 HDT: Heritage Digital Twin
 HPC: High Performance Computing
 HTML: HyperText Markup Language
 HTTP: HyperText Transfer Protocol
 HTTP/2: HyperText Transfer Protocol 2.0
 IaaS: Infrastructures as a Service
 IaC: Infrastructure-as-Code
 ICT: Information and Communication Technologies
 IFC: "Industry Foundation Classes" file format
 IFC 2X3: "Industry Foundation Classes 2x3" file format
 IFC 4: "Industry Foundation Classes 4" file format
 IGES: "Initial Graphics Exchange Specification" file format
 IAM: Identity and Access Management
 ISO: International Standards Organization
 JSON: JavaScript Object Notation
 JWT: JSON Web Token

KB: Knowledge Base
LAS: "LASer" file format
LAZ: "LASzip" file format
MEP: Mechanical Electrical Plumbing
MR: Mixed Reality
NURBS: Non-Uniform Rational Basis-Splines
OAuth2: Open Authorization 2.0
OBJ: "Object" or "Wavefront" file format
OIDC: OpenID Connect
OPA: Open Policy Agent
PaaS: Platform as a Service
PLY: "Polygon" or "Stanford Triangle" file format
PTS: Laser scan plain data format
PTX: Laser scan plain data format
RDF: Resource Description Framework
REST: Representational State Transfer
RGB: Red Green Blue
RKE: Rancher Kubernetes Engine
SaaS: Software as a Service
SfM: Structure from Motion
SSL: Secure Sockets Layer
STEP: "Standard for the Exchange of Product Data" file format
STL: "Standard Triangle Language" or "Standard Tessellation Language" file format
Sub-D: Subdivision Surface Modeling
THESPIAN: Tool for HERitage Science Processing, Integration and ANalysis
URL: Uniform Resource Locator
VM: Virtual Machine
VR: Virtual Reality
VRLM: "Virtual Reality Modeling Language" file format
WP: Work Package
WSGI: Web Server Gateway Interface
X3D: "Extensible 3D" file format
X3ML: code-name for an XML based language to describe schema mappings
XDC: eXtreme-DataCloud
XML: Extensible Markup Language
XR: eXtended Reality
XRF: X-Ray Fluorescence
XYZ: "XYZ" file format

List of figures

Figure 1: Schema of digital contents produced in the process of digitization of Cultural Heritage with specific reference to 3D models.	4
Figure 2: Parameters for the evaluation of point cloud models (reality captured).	10
Figure 3: Parameters for the evaluation of mesh models (reality captured).	11
Figure 4: Parameters for the evaluation of solid, surface and mesh models (digital born).	12
Figure 5: Cloud service layers	21
Figure 6: 4CH platform design and components	23
Figure 7: Longhorn service deployed on 4CH platform	24
Figure 8: 4CH Platform and ancillary services.....	25
Figure 9: Authorization and authentication workflow for applications within the 4CH platform.....	26
Figure 10: 4CH platform monitoring using Grafana and Prometheus.....	27
Figure 11: 4CH platform image repository using Harbor.	28
Figure 13: THESPIAN-Mask application integration in the 4CH Cloud Platform	32
Figure 14: THESPIAM-XRF application federation in the 4CH Cloud Platform.....	34
Figure 15: Example of Integrated and Federated applications in the 4CH Cloud Platform	35
Figure 16: Example of hybrid application deployment in the 4CH Cloud Platform	35

List of tables

Table 1: Main open, public and standards formats organized on 3 macro-categories of models: point cloud models, solid, surface or mesh models, and BIM - HBIM models	8
Table 2: Main features of open, public and standard formats for point cloud models.	9
Table 3: Main features of open, public and standard formats for solid, surface and mesh models.....	9
Table 4: Distribution model and license type for each analyzed solution.	17
Table 5: 4CH Platform services, related application versions and endpoints.....	24

Table of Contents

<i>Executive summary</i>	1
1. <i>Introduction</i>	2
1.1 Objectives and structure of the deliverable	2
2. <i>Analysis of 3D data generation and creation for Cultural Heritage</i>	3
2.1 3D data types	4
2.1.1 3D models	5
2.1.2 Images, videos and technical sheets	6
2.2 File formats	7
2.2.1 File formats for 3D point clouds.....	8
2.2.2 File formats for 3D models	9
2.3 Evaluation parameters by type	10
2.3.1 Reality captured 3D models	10
2.3.2 Digital born 3D models	12
2.4 Other 3D digitization activities	12
3. <i>Software tools</i>	14
3.1 Areas and main functions by software	14
3.2 Software by licence	17
4. <i>Towards the development of a 4CH platform</i>	19
5. <i>The architecture of the 4CH Cloud Platform</i>	21
5.1 The 4CH Cloud Platform: design and approach	22
5.2 4CH Platform services	24
5.2.1 Authentication / Authorization workflow.....	25
5.2.2 Monitoring.....	26
5.2.3 Image repository.....	27
5.2.4 Web portal homepage.....	28
5.2.5 Developing services	28
6. <i>4CH Platform policies and requirements: integration and federation</i>	29
6.1 Open-source advantages and possibilities	29
6.2 Standard, Best Practices and conventions	29
6.2.1 Implementation Standards	29

6.2.2	Documentation standards	29
6.2.3	Software components and classification	29
6.3	Requirements for service integration	30
6.3.1	The integration use case: THESPIAN – Mask	31
6.4	Requirements for service federation	33
6.4.1	The federated use case: THESPIAN-XRF	34
6.5	Mixing service Integration and Federation	35
7.	<i>The 4CH Knowledge Base</i>	36
7.1	Introduction to the Knowledge Base	36
7.2	Organizing the data: the HDT ontology	37
7.3	The aggregation pipeline	40
7.3.1	The mapping phase.....	40
7.3.2	The cleansing phase	41
7.3.3	The aggregation workflow	41
7.3.4	Implementing the 4CH KB.....	42
8.	<i>Conclusions</i>	43
9.	<i>References</i>	44
	<i>Appendix 1</i>	39

Executive summary

The Competence Centre for the Conservation of Cultural (4CH) project is a project approved in January 2021 within the DT-TRANSFORMATIONS-20-2020 call of the Horizon 2020 framework program of the European Community. Its goal is to design and prepare for a European Competence Centre (CC) on the Conservation of Cultural Heritage which will work proactively for the preservation and conservation of cultural heritage (CH).

The main high-level topics addressed by the project include the implementation of structure, organization and services of the CC which will operate as a virtual infrastructure providing expertise, advice and services using state-of-the-art ICT with a special focus on 3D technology.

In particular, the document provides the initial definition of the 4CH Cloud Platform as resulting from the first stage of the work of T3.1 - 4CH platform project and architecture and integration of existing and available tools and T3.3 - Cultural Heritage Cloud, including:

- the analysis of types of data generated in the 3D digitization process;
- the recognition of the software tools available and used by the Cultural Heritage community in order to design those that will be proposed them as services via the 4CH Platform;
- the recognition and adoption of cloud-oriented services aimed to build a suitable platform where to host the CH software tools.

Such information has been then collected and used by the task participants to design and implement the Cultural Heritage Cloud platform in distributed computing and its involvement in many e-infrastructure EU-funded projects. The functionalities of the CH Cloud, and the related services and tools adopted, will be tailored to the specific user needs by enabling a collaborative approach.

The document will also provide policies and requirements needed for the integration and federation of the CH tools and services deployed into the Platform and made available for the heritage community.

Moreover, a selection of the CH services hosted and made available through the 4CH Platform will be promoted via the EOSC marketplace in order to be used and adopted also from a wider community.

1. Introduction

1.1 Objectives and structure of the deliverable

The main objective of the Deliverable is to present an overview on available software tools used by the heritage community in the digital documentation of Cultural Heritage assets and describe the architecture of the 4CH platform as a cloud-based infrastructure providing services to the community.

The document is structured as follows.

In Section 2, the outputs of the 3D digitization process of Cultural Heritage are analysed in terms of data types, evidencing the characteristics of different 3D data formats, together with technical parameters related to their features and intended for data quality evaluation.

Section 3 is an overview on available software tools for creating, processing, and managing 3D data, including solutions both commercial and free. It has been organized to be offered to users for orienting them in a wide and complex scenario. The analysis collected more than 70 software tools, representing those mainly used by professionals working in the sector.

Section 4 bridges the gap between the latest efforts to standardize the 3D digitization process and the need for a seamless ecosystem providing services and tools to the whole Cultural Heritage community. In this sense, the section introduces the 4CH Cloud Platform in relation with two future actions such as the Data Space for Cultural Heritage and the European Collaborative Cloud for Cultural Heritage.

In Section 5, the 4CH Cloud Platform architecture is defined and explained. The 4CH Platform is a cloud-based infrastructure for hosting the 4CH services for cultural heritage and implemented as a pilot infrastructure hosted on the Openstack cloud at INFN. The adopted tools and solutions to provide distributed computing resources, network, and storage, are presented and discussed. Moreover, on top of the 4CH Platform, some ancillary services providing useful functionalities such as reverse proxy, image repository, monitoring and authorization and authentication, already deployed, are described.

In section 6, best practices, policy, and requirements put in place to integrate and federate CH applications made available for the project have been addressed and described. As a support, detailed examples describing both the integration and the federation of applications with the 4CH Cloud Platform are described and analysed.

Section 7 introduces the main aspects of the 4CH Knowledge Base, how it is organized and integrated in the 4CH platform.

Section 8 drafts the final conclusions and outlines future activities.

2. Analysis of 3D data generation and creation for Cultural Heritage

3D digitization of Cultural Heritage can be performed by means of several procedures and tools that could differ for scope and methods. The deliverable “D4.1 - Report on standards, procedures and protocols”, currently under development and foreseen for M21, will analyse existing procedures. Among others, it will normalize standard workflows, relating them to specific scopes as well as values and needs for digitization. However, based on the work carried out for “D1.1 – Initial survey of the experiences and technology state of the art” [1], it is already possible to note how all the 3D digitization processes and solutions have in common the generation of a specific range of digital data. That deliverable includes a list of digitization instruments, main data processing and valorisation techniques, and an initial list of data formats.

In the analysis presented in this deliverable, according with the DoA (Document of Action – annex 1 to the Grant Agreement), a particular focus has been put on the role of 3D data in the digitization of Cultural Heritage and how these data can offer an added value being a useful resource that can generate a variety of traditional outputs (images, videos, technical drawings, quantities, etc.). Thus, data types and formats have been arranged in a more structured way, aggregating data types by the overall type of the content, distinguishing, where relevant, between reality-captured and digital-born data.

3D digitization of an asset, and especially of Cultural Heritage ones, is generally intended as the production of a 3D digital asset (a 3D model) which represents a specific instance of reality, based on information directly collected on or gathered from original but somehow interpreted. The perfect and exact copy of a real item doesn't exist; it can be close to reality as much as possible, but it will always contain a deviation due to instrumental errors, approximations introduced by the operator in the phase of capturing or processing (resolutions, processing algorithms, decimations, etc.) or even interpretation performed by the user in the modelling phase. This gap between reality and the model can be effectively reduced by collecting metadata tracking all these aspects, helping in the correct interpretation of the created 3D model. The already mentioned D4.1 will explore these possibilities.

Anyway, the 3D models described above can be easily called reality-captured since the source data are directly coming from the original assets. But in the sector, we also face the use of 3D models where the interpretative part is pushed even further. In fact, 3D models may also contain or represent information extracted not directly from reality but from critical assumptions based on historical documents, exploring its configuration in different ages, only partially corresponding to the actual situation. In this case we speak of 3D reconstructions, a digital born model where 3D modelling technologies for creating such digital assets can significantly differ.

Schemas and tables presented in the following paragraphs are intended to differentiate data types in such direction, list the related file formats together with their features and eventually indicate possible evaluation parameters for data quality.

2.1 3D data types

3D data types are usually directly related to data formats since a content, being an image, a video, or a text, can be archived and stored in the proper file format. However, when it comes to 3D models, due to the complexity of the content and the features that it requires, the differentiation could become quite complex if we do not break down the general category of 3D models into sub-categories. The schema below (Figure 1) is intended to help in understanding such subdivision and to highlight the relation and dependencies between 3D models and other content types. Contents of different types (images, videos, technical sheets, etc.) are related to the 3D model as data that can contribute to 3D digitization (green) or possible extrapolation (light blue). In addition, some 3D models can generate the creation of other models of different types or with different interpretative purposes (yellow).

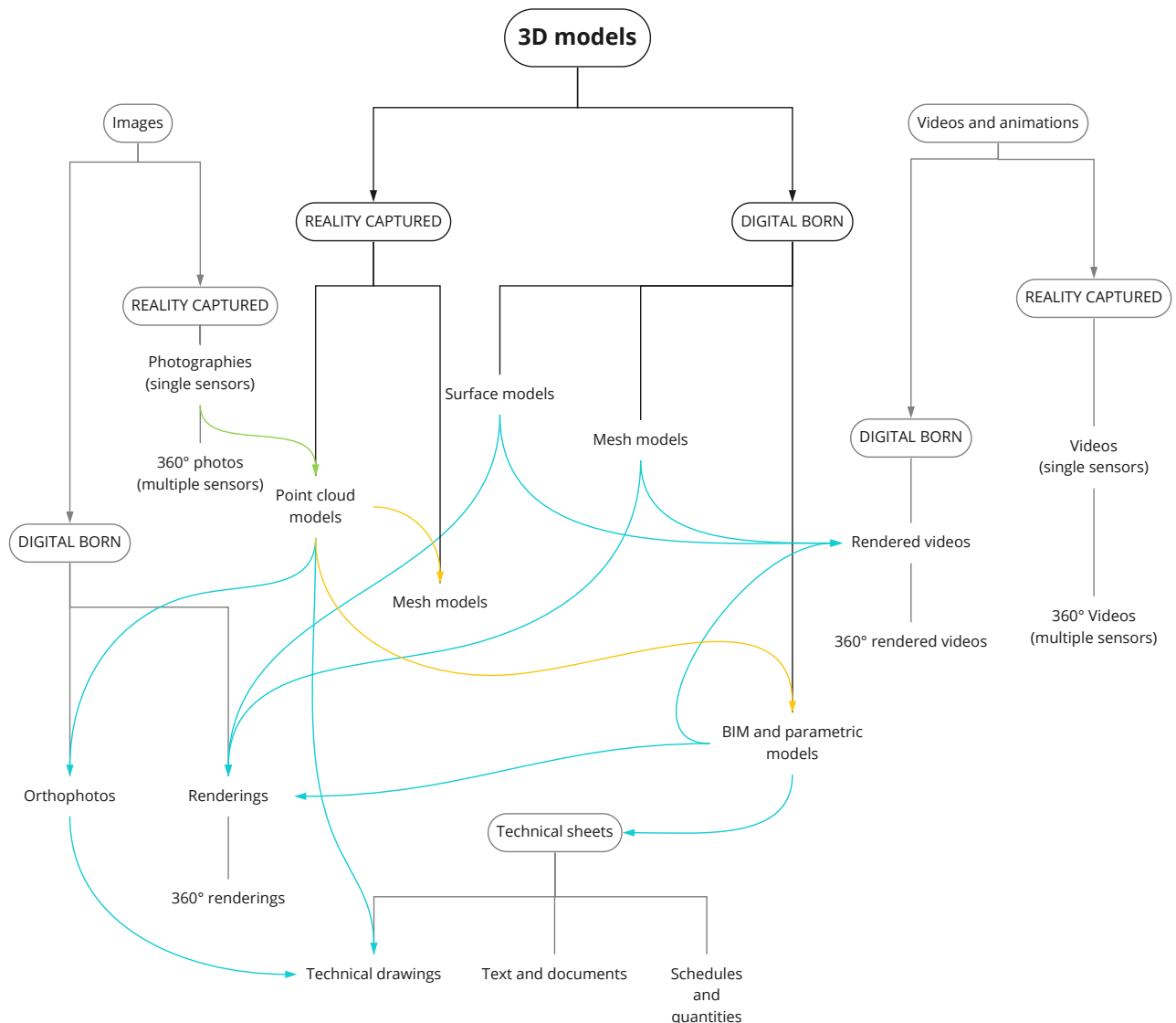


Figure 1: Schema of digital contents produced in the process of digitization of Cultural Heritage with specific reference to 3D models.

The following descriptions illustrate some details of the schema presented in Figure 1.

2.1.1 3D models

In the Figure 1, 3D models are broken down into 4 main typologies:

- Point cloud models coming directly from the digitization procedures;
- Mesh models deriving from them or totally digital born;
- Solid and surface models, digital born but ideated for mimicking reality;
- BIM and Parametric models, digital born and intended to be used for the management and interpretation of the reality.

A short description of these data types is reported below.

Point cloud models (reality captured)

A point cloud is a set of points in space. The points may represent a 3D shape or object. Each point position has its set of Cartesian coordinates (X, Y, Z). Point clouds are generally produced by:

- 3D scanners, which measure many points on the external surfaces of objects around them;
- photogrammetry software, which uses Structure from Motion (SfM)¹ algorithms for recreating the shape of an object or an environment from a bunch of photos from different angles. The green arrow of the schema sets a dependency from the traditional digitization performed by photos and these models: existing set of photos can be processed for creating new 3D models or photos captured for this purpose can be reused differently.

Further information on capturing technologies is already available in D1.1 [1]. It can be highlighted here that, as the main output of 3D data capturing processes, point clouds are used for many purposes and as the basis for creating other 3D models, such as mesh models (for analysis, visualization, and documentation purposes) or BIM and parametric models (for analysis and technical purposes). If coming from scanner, point clouds always have a metric value (accuracy may differ) while if generated by SfM the metric information is valuable only if specific procedures are applied (integration of reference measurements such as laser, GPS or topography for instance). In both cases point may also have colour value (RGB).

Mesh models (reality captured and digital born)

In 3D computer graphics and solid modelling, a polygon mesh is a collection of vertices, edges and faces that defines the shape of a polyhedral object. In the Cultural Heritage sector (as in many others), the faces consist of triangles (triangle mesh) even if technically they could be also quadrilaterals (quads), or other simple convex polygons (n-gons). The triangle mesh model is the easiest way to represent a surface or a solid (closed surface as boundary representation – B-Rep method²).

Mesh models can be reality captured if generated from point clouds through interpolation, creating flat surfaces that connect the detected points and approximate the surface of the original asset. This operation is performed via software based on approximation parameters (e.g., point decimation, offset angle, point normal, etc.) provided by the operator. Otherwise, mesh models can be generated by solid and surface modelling procedures, saving 3D

¹ https://en.wikipedia.org/wiki/Structure_from_motion (accessed on 18/07/2022)

² https://en.wikipedia.org/wiki/Boundary_representation (accessed on 18/07/2022)

modelling instances based on a more complex CGI techniques (primitives and booleans, NURBS, Sub-D, sculpting, etc.) into a simpler shape. In this case, most of the 3D models are digital born.

Mesh models may have coloured vertex and triangles, as well as textures projected on the surface. In both cases, the procedure requires one or multiple images (as external reference) and a project matrix saved within the file.

Solid and surface models (digital born)

Solid modelling is a consistent set of principles for mathematical and computer modelling of three-dimensional solids. Mesh models are actually a sub-category of such models, but they are treated separately here because of their use scope. Solid and surface models are a more generic category including a vast amount of CGI techniques such as primitive solids (e.g., cube, sphere, cone, etc.); Boolean operations (union, difference, intersection) on those shapes or more complex ones; 3D generation from 2D shapes (e.g., extrusion, revolution, sweep, loft, etc.); and Non-Uniform Rational Basis-Splines (NURBS) and Sub-D surfaces, sculpting, etc. All these methods have been developed as 3D modelling techniques intended to represent with physical fidelity, mimicking the reality but modelling it by observation rather than by measurements. This means that these models must rely on mathematical formulas for describing shapes and other characteristics, resulting in more complex file format. Furthermore, such models in general support also animations, simulations, annotations, texture, and material behaviours and they are optimized in terms of file size (storing only the generating maths). These models are oriented to the creation of digital assets offering a high-quality engaging experience (e.g., edutainment, games, immersive environments, etc.) of the content and are conceived as digital born.

BIM and parametric models (digital born)

Building Information Modelling (BIM) is a process supported by various tools and technologies involving the generation and management of digital representations of physical and functional characteristics of places. The concept of BIM, introduced in the 1970s but used without a common definition until the early 2000s, was originally developed for adopting 3D technologies in the Architecture Engineering and Construction (AEC) sector with particular reference to the design and construction of new buildings. In BIM models, the digital representation of geometries goes together with the virtualization of qualitative and quantitative parameters such as, for instance, those relating to construction materials, mechanical strength, transmittance, or cost.

As an extension of this technology, H-BIM (where H stands for Historical or Heritage) is becoming more common for management and conservation of historic buildings, and in some cases for archaeological monuments. BIM is a mature system developed over twenty years, which allows for great information management and collaboration between multi-disciplinary teams. For its nature, it can integrate a variety of data (including point clouds and reality captured mesh models). Its added value can also be exploited as a technical model for simulations and management of interventions, where the data need to be consistently interpreted in the generation of a digital born model. Recent research also includes using parametric modelling and algorithms that can explicitly represent the generation of the shape, offering still unexplored ways of understanding their significance.

2.1.2 Images, videos and technical sheets

Figure 1 also shows images, videos, and technical sheets as other outputs of the digitization effort. In the first two cases, images and videos can be captured from reality or generated by

computer models, while technical sheets (with particular reference to drawings) are nowadays generally created by extrapolation from 3D models.

It is worth to notice the increasing use of 360° reality captured photos and videos, which can offer to the final user an immersive and somehow 3D experience. They are actually little more than traditional photos and video captured by the use of multiple sensors, where the outputs are stitched together. Similarly, digital born ones exploit the generation from a 3D model but still result in regular images and videos. Quite interesting is the increasing use of 3D models for generating orthophotos: a corrected (orthorectified) photo where scale is uniform and that can be measured. Orthophotos, traditionally used in the CH sector, were once obtained by correcting the distortion while nowadays they are generated by 3D models (in particular from those created with SfM techniques) resulting in a more accurate as well as faster result. As stated before, also producing technical drawings, rendered images and videos representing simulations on the digitized object benefits of the possibility of being derived from 3D models. For this reason, a clear overview on 3D data can provide advantages to the overall digitization process.

2.2 File formats

Three-dimensional acquisition instruments and modelling software generate and work primarily on proprietary file formats. Proprietary formats are often optimized both in terms of performance and file size, but they force the user to remain loyal to one or more manufacturers or developers and thus, in fact, dependent on their commercial policies.

On the other hand, there are the public, open, standard file formats that can be used for modelling or generally manage 3D data. The advantage of such formats is the ability to preserve and reuse data over time and to view, manage, and modify it in all those tools, open or proprietary, that will have implemented such a standard. The adoption of open formats has already widely demonstrated a significant benefit in terms of use and reuse possibility.

In contrast, open formats are not always optimized often do not contain all the information that original, proprietary formats had.

However, adopting open and standard protocols for files can be considered overall an advantage as well as a best practice. Thus, in Table 1 below, the main open, public, and standard formats are listed, organized in 3 macro-categories of content types: point cloud models, surface or mesh models (which even if different for the purpose of application, often do not differ in file format – please see paragraph 3.2.1 for more details), and BIM (and HBIM) models. It should be noted how most of the formats can technically also be used in more than one category (e.g., the PLY format can contain both mesh and points, just as IFC models can contain both objects defined by parametric primitives and by mesh), but we prefer to emphasize their most common use scope.

Table 1: Main open, public and standards formats organized on 3 macro-categories of models: point cloud models, solid, surface or mesh models, and BIM - HBIM models

<i>Point cloud models</i>	<i>Solid, surface or mesh models</i>	<i>BIM – HBIM models</i>
<i>LAS/LAZ</i>	<i>DXF</i>	<i>IFC 2X3</i>
<i>PLY</i>	<i>OBJ</i>	<i>IFC 4</i>
<i>XYZ</i>	<i>DAE</i>	<i>COBie</i>
<i>PTS</i>	<i>PLY</i>	<i>BCF</i>
<i>PTX</i>	<i>STL</i>	
<i>E57</i>	<i>IGES</i>	
	<i>STEP</i>	
	<i>VRLM</i>	
	<i>X3D</i>	
	<i>glB/glTF</i>	

Next then, focusing particularly on the first two categories (point clouds and 3D models for surfaces or meshes), there is an attempt to make a comparison between the formats, highlighting the types of data that can be stored and the available features. BIM models are not yet analysed since the only open, public, and standard format for storing geometries together with information is the IFC (COBie and BCF are format storing only information). However, an in-depth analysis of BIM capabilities will be presented in the “D3.2 - Integration of the INCEPTION 3D and HBIM technologies”.

2.2.1 File formats for 3D point clouds

Point cloud models, as highlighted before, do not actually contain 3D geometries but rather a set of points. Thus, the file format needed for that is quite simple, requiring just 3 coordinates (X,Y,Z) to define a point. This is the only necessary and sufficient condition for dealing with a point cloud. However, a point cloud file format can (and should) contains more data. Among them, those analysed in the following comparisons are:

- RGB (Red, Green and Blue) value for identifying the colour of the point (values are extracted from photos matching points and pixels);
- Intensity value representing a function of the power of the received backscattered signal (which could be useful for data interpretation: different scanning technologies report different intensity values on different materials).

Furthermore, a point cloud model could be defined as structured or unstructured. Structured point clouds express points referred to the origin of scan's position and organized in order of acquisition.

Structured point clouds can be achieved only using scan-based solutions with fixed-origin or tripod scans and without merging multiple scan positions.

The table below also shows if the file format can be saved as ASCII or Binary.

Table 2: Main features of open, public and standard formats for point cloud models.

	XYZ	RGB	Intensity	Structured point clouds	Unstructured point clouds	ASCII	Binary
LAS/LAZ	X	X	X		X		X
PLY	X	X	X	X	X	X	X
XYZ	X	X	X		X	X	
PTS	X	X	X		X	X	
PTX	X	X	X	X		X	
E57	X	X	X	X	X	X	X

2.2.2 File formats for 3D models

In the case of complex 3D models, the number of entities to be stored is large. Therefore, it is important to evaluate the complexity of geometries that can be stored in open, public, and standard formats. In most of the cases, mesh surfaces can be saved since they represent the easiest way to describe a geometry. IGES and STEP, intended to be used for CAD/CAM applications, require instead geometries expressed in terms of their generative mathematical formulae. As mentioned before, sometimes also point clouds can be stored in such formats, even if they are not optimized for this purpose. All the features related to the model visualization are instead quite important here: textures and materials, lights and cameras, as well as animations, kinematics and physical effects to visualize movements and simulations.

Table 3: Main features of open, public and standard formats for solid, surface and mesh models.

	Point clouds	Mesh	Solid or Surfaces	Colour	Texture and materials	Audio	Lightning	Cameras	Animations	Kinematics	Physical effects
DXF	YES	YES	YES	YES	NO	NO	NO	NO	NO	NO	NO
OBJ	YES	YES	YES	YES	YES	NO	NO	NO	NO	NO	NO
DAE	YES	YES	NO	YES	YES	NO	YES	YES	YES	YES	YES
PLY	YES	YES	NO	YES	NO	NO	NO	NO	NO	NO	NO
STL	NO	YES	YES	YES	NO	NO	NO	NO	NO	NO	NO
IGES	YES	NO	YES	YES	NO	NO	NO	NO	NO	NO	NO
STEP	NO	NO	YES	YES	NO	NO	NO	NO	NO	NO	NO
VRML	YES	YES	NO	YES	YES	YES	NO	NO	NO	NO	NO
X3D	YES	YES	NO	YES	YES	YES	YES	YES	YES	YES	YES
glTF	NO	YES	NO	YES	YES	YES	YES	YES	YES	YES	YES

2.3 Evaluation parameters by type

The main purpose of analysing 3D data created and generated in the digitization of Cultural Heritage is creating a common basis for data use and a possible technical evaluation.

The analysis of data types and file format presented above allows:

- a more aware management of created and generated data;
- together with the “Digitizing Cultural Heritage - A Reconnaissance Investigation on the Data Infrastructure” [2], published by 4CH as an internal report and analysing data in terms of size, it constitutes a solid base for designing the 4CH cloud platform and services running on it.

However, in order to allow the services to run effectively, data should be complete and reliable. The following three schemas try to summarize what can be technically evaluated on 3D data, listing parameters referred to the most relevant data types.

2.3.1 Reality captured 3D models

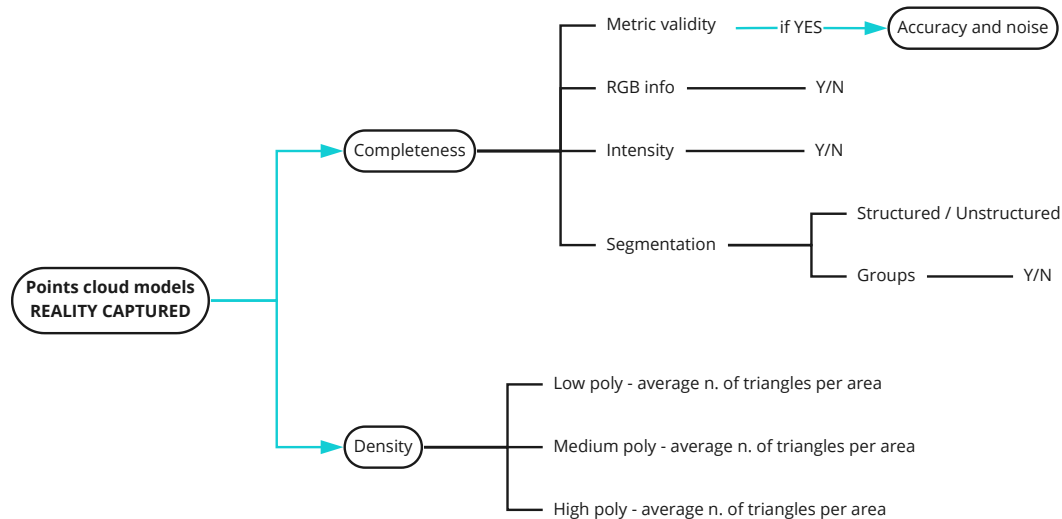


Figure 2: Parameters for the evaluation of point cloud models (reality captured).

For example, in archaeological sites it is frequent the case of scattered remains with few still standing structures. 3D cannot provide the information concerning the spatial and logical relationships between the different parts on which archaeologists rely to produce their interpretation. A significant step forward in computer-assisted archaeological research was the introduction (some 30 years ago) of Geographical Information Systems (GIS) as a tool to store and synthesize spatial relationships among parts of an archaeological site, often identified as such only by the objects found there. A recent publication [5] explores the advantages of combining GIS tools with 3D digitization.

Both historic gardens and fountains have an alive component, the vegetation in the former and water in the latter. Vegetation changes continuously and changes the appearance of the garden accordingly. Fontana di Trevi in Rome, famous also in the popular culture, would lose much of its appeal without its water jets. All these rapidly changing features are difficult to document in a relatively static 3D model and will need further investigation to be properly documented.

Finally, spaces that derive their heritage value from what happens, or happened, are very poorly represented by a 3D model. Examples from somehow opposite perspectives are a pilgrimage route as the Camino de Santiago and a battlefield such as Culloden in Scotland, the latter outside the EU but an exemplary case of heritage at risk. In such cases links to history, traditions and beliefs must accompany and integrate the visual documentation.

It is anticipated that these issues will be addressed in the planned Knowledge Base (see section 7 for details), in which the visual aspect is just one – although very important – component of heritage documentation. Understanding which additional visual documents, for example videos, are required will also allow to evaluate the related formats.

3. Software tools

The capturing and processing phases of the 3D digitization of Cultural Heritage strongly relies on software tools. An overview on capturing technologies and generally on software tools has been already provided by the already mentioned “D1.1 - Initial survey of the experiences and technology state of the art” [1] and the “Internal Technical Report - Digitization techniques in the field of Cultural Heritage”[6]. To understand how software tools answer the need of the CH professionals, understand where gaps are, and provide indications for future developments, an in-depth review of the current solution has been carried out.

More than 70 software tools and services have been compared and evaluated here. In order to provide a common ground and thanks to the activities developed together with the participants of the Task 4.2, the analysed solutions has been mapped within the 3D digitization workflows under development³, highlighting available features that could be referred to specific digitization phases. In the analysis, both commercial and free solutions have been taken into consideration and, for completeness, also licensing options, software distribution and usage policies have been annotated.

As a general consideration, must be noticed how the current solutions for the sector are still desktop-based in a world where everything is moving to the web, as demonstrated by the increasing need for visualizing solutions. Resource-consuming processes and the increasing size of data slowed down the transitions to web applications, but the overall trend is clearly delineated: that means we should have a clear picture of what is existing today, what is working and what is missing. The interpretation of the comparison presented below should provide the answers to those questions.

Furthermore, this analysis should be considered as a tool for orienting professionals in a wide and complex scenario, where it's quite rare to make a transparent choice. A comparison that is not based on generally available software features but that is structured on features needed for performing a 3D digitization process as a whole, it represents also a contribution to the knowledge of the sector. For this reason, the analysis can be considered also as one of the pieces of the Knowledge Base that will be designed and realized in the just started Task 3.4 - Cultural Heritage Knowledge Base.

3.1 Areas and main functions by software

The 75 solutions analysed and reported in the matrix in the Appendix 1 of this deliverable, cover the following areas of the 3D digitization.

Data capturing

- Laser scanning
 - Scan acquisition
 - Scan pre-registering

³ The list of features to related to specific digitization workflows were still preliminary, thus some minor discrepancies can be there. The analysis, however, will be kept up to date for the whole duration of the project and fine-tuned as soon as the D4.1 will be released.

- Scan importing
- Scan registering
- Scan merging
- Photogrammetry (SfM)
 - Photo alignment
 - Photo target matching
 - Topographic correction
 - Points cloud generation

Data processing

- 3D object modelling
 - 3D object importing
 - 3D object creation
 - 3D primitive modelling
 - 3D NURBS modelling
 - 3D mesh modelling
 - Point cloud meshing
 - 3D object managing
 - 3D object editing
 - Mesh decimating
 - Mesh repairing
 - Free form editing
 - Sculpting
 - Boolean editing
 - 3D object analysis
 - 3D properties analysis
 - 3D object comparison
 - Mesh doctor
- BIM modelling
 - BIM model importing
 - BIM model authoring
 - Architectural
 - Structural
 - MEP
 - Other
 - BIM model coordination
 - IFC data editing
 - IFC geometric data editing
 - IFC semantic data editing
 - IFC metadata editing
 - CDE repository
 - Model checking
 - Visual checking
 - Clash detection
 - Code checking
- Technical drawing
 - Layout organization
 - 2D extraction/generation
 - 2D drawing
 - 2D annotation
 - Sheet elaboration

- Texturing
 - Texture managing
 - Photo projection
 - Texture mapping
 - Texture baking
 - Mosaicking
- Rendering
 - Rendering managing
 - 3D city map importing
 - Weather simulation
 - Material characterization
 - Cameras creating
 - Lightings creating
 - Rendering
- Animating
 - Animation managing
 - Rigging/Kinematic chain
 - Keyframing
 - Physical effects

Data use and visualization

- Viewer platforms
 - Repository
 - Geometric viewer
 - Point cloud model viewer
 - Surface model viewer
 - Mesh model viewer
 - BIM model viewer
 - Extended reality viewer
 - AR
 - MR
 - VR
 - XR
 - Time machine navigator
 - Interactive tools
 - Annotating
 - Measuring
 - Slicing tools
 - Download

From the matrix emerges that tools intended for the first processing of captured 3D data perfectly manages point cloud models and initial meshing operations but lacks all the other features beside some exporting for technical drawings.

Modelling tools are nowadays the solutions covering most of the workflow phases, from some basic function on point clouds models to technical drawing, animation, and rendering. Sometimes they also already offer basic BIM features.

BIM tools are instead quite specialized. They can exploit the point cloud but they are specifically addressed to the creation of architectural models and technical drawings generation (even if the IFC format can be actually exploited in several ways for storing information together with geometries).

Visualization solutions are also quite specific. From one side, there are software addressed to the generation of rendered images and videos and, on the other, web-based services and platforms for visualizing data.

In the end, the table in Appendix 1 shows how professionals still need to integrate the use of several software tools when processing data across the whole 3D digitisation, being requested to install and manage different solutions. The transition to a seamless web solution will for sure provide a huge benefit in terms of always up to date tools, working despite the platform and accessing to computing resources only when needed.

On the other side, CH professionals need also to use, most of the time, commercial solution and thus get the correct license.

3.2 Software by licence

This paragraph is simply intended to report on the distribution model and licence type for each analysed solution. Different colours in rows indicate the intended use of the software solutions: blue for data capturing, yellow for 3D modelling, green for animations and renderings, red for BIM and light blue for visualization.

Table 4: Distribution model and license type for each analyzed solution.

License typology	Open source	Free software	Proprietary software (freeware)	Proprietary software (shareware)	Proprietary software (periodic licence fee - annual/monthly)	Proprietary software (perpetual licence)	SaaS - Software as a Service
CAM2 Suite				X		X	
Leica Cyclone			X	X	X	X	
Trimble Suite							
AVEVA LFM							
Stonex				X	X	X	
Suite X-PAD							
Artec Studio 16				X	X	X	
Recap PRO			X	X	X		
Metashape				X	X	X	
3DF Zephyr				X	X	X	
Reality Capture				X	X	X	
Cloud Compare	X	X					
Meshroom	X						
PIX4D Suite				X	X	X	
Regard3D	X	X					
OpenDroneMap	X						
Geomagic Design X				X	X		
Geomagic Wrap				X	X		
DJI Terra Pro				X	X	X	
Meshlab	X						
Sketchup			X	X	X		
AutoCAD				X	X		
Rhinoceros 3D				X		X	
Blender	X	X					
Maya 3D				X	X	X	
Modo				X	X		
Zbrush				X	X	X	
Unity			X	X	X		
Vuforia			X	X	X		
NetFabb				X	X		
Magics				X	X	X	
Twinmotion				X		X	
Unreal Engine							
3DS Studio Max				X	X		

License typology	Open source	Free software	Proprietary software (freeware)	Proprietary software (shareware)	Proprietary software (periodic licence fee - annual/monthly)	Proprietary software (perpetual licence)	SaaS - Software as a Service
Cinema 4D				X	X	X	
Lumion				X		X	
Enscape				X	X		
Keyshot				X	X		
Maxwell Render Bundle				X	X		
Revit				X	X		
Civil 3D				X	X		
Naviswork				X	X		
ArchiCAD			X	X	X		
Edificius				X	X		
usBIM.checker				X	X		
usBIM.clash				X	X		
usBIM.editor				X	X		
usBIM.viewer+			X				
usBIM.platform				X	X		
AllPlan				X	X		
Vectorworks				X	X	X	
Solibri				X	X	X	
Tekla Structures			X	X	X	X	
Trimble Connect				X	X		
Bentley Micostation				X	X		
Bricsys BricsCAD BIM				X	X	X	
SimpleBIM				X	X	X	
3D Hop	X						
Hexalab	X						
Inception							X
Plas.io	X						
Potree	X						
Sketchfab							X
Smithsonian Museum X3D visualizer	X						
Pano2vr				X	X		
X3DOM	X						
Flyvast				X	X		
Cintoo				X	X		
3dusernet				X	X		
Euclidean			X	X	X		
Cesium	X						
Web BIM3D							X
3D warehouse							X
itwin platform	X						
Modelo.io				X			X

As it can be noted, more technical solutions such those intended for BIM or advanced visualization are available as proprietary software. Open-source software can compete in the pure 3D modelling and somehow in the first processing of captured data (even if here the dependency instrument – software is quite high). Web-based visualization solution are instead equally spread among licence types, even if solutions are most of the time deeply different in what they offer.

4. Towards the development of a 4CH platform

Achieving a rich and comprehensive digital documentation of the Cultural Heritage has been among the main European challenges since long. In 2018 the European Year of Cultural Heritage was celebrated to recognize the existing interest of Europe on Cultural Heritage. The entire Horizon 2020 framework program invested more than €500 million, in order to start to capitalize on the many research results achieved over the years. After that, with the “Declaration of Cooperation on Advancing Digitisation of Cultural Heritage”⁴, signed by 24 EU Member States in April 2019, the joint efforts in European initiatives were even more consistent and articulated on three pillars: the 3D digitization of Cultural Heritage (pillar I); the re-use of digitized cultural resources to foster citizen engagement and innovative use in other sectors (pillar II); enhanced cross-sector, cross-border cooperation and capacity building (pillar III).

Following such intention, several activities have been developed both to facilitate the 3D digitization and to offer a common environment for storing, accessing, and using the collected data.

The European Commission tasked the DCHE (Digital Cultural Heritage and Europeana) Expert Group to the development of guidelines on 3D cultural heritage assets. Thus, the Expert Group elaborated a list of 10 basic principles and a number of related tips for each of them geared toward cultural heritage professionals, institutions and regional authorities in charge of Europe's precious cultural heritage. The principles were published in August 2020 [7]. This topic was further investigated in the “Study on quality in 3D digitisation of tangible cultural heritage” [4] that demonstrates how complexity and quality are fundamental considerations in determining the necessary effort for a 3D digitisation project to achieve the required value of the output. In short, the study highlights how there are no internationally recognized standards or guidelines for planning, organising, setting up and implementing a 3D data acquisition project for any size of asset; quality parameters refer to different stages of the 3D digitisation process; and thus, there is a pressing and urgent need for a technical specifications and guidelines. The results presented above aim to contribute to that by organizing types, formats, and technical evaluation parameters. However, this is just a starting point and a more details will be provided in the “D4.1 - Report on standards, procedures and protocols” as already mentioned at the beginning of Section 2.

On the other hand, specific actions pursued the achievement of an enhanced valorisation of 3D digital data. Between 2019 and 2020, Europeana established the Task Force “3D Content in Europeana”. The Task Force analysed valuable content on 3D digitization of cultural artefacts at large, in the perspective of their integration in Europeana, also discussing a number of related issues such as, for example, data formats, standards and storage of the models [8]. The effort precluded the explicit need made clear by the DIGITAL EUROPE Work Programme 2021-2022 [9] foreseeing the creation of Data Space for Cultural Heritage: an European common data space to provide support to the digital transformation of Europe's cultural sector, and foster the creation and reuse of content in cultural and creative sectors. It will build on the current Europeana platform, vastly expanding the current 3D functionalities. It must be noticed that improvement in 3D data management will be not only qualitative but also quantitative,

⁴ <https://digital-strategy.ec.europa.eu/en/news/eu-member-states-sign-cooperate-digitising-cultural-heritage> (accessed on 21 June 2022)

since the “Recommendation on a common European data space for cultural heritage” [10] published by the Commission in November 2021 suggest the 3D digitization of all monuments and sites at risk, and at least 50% of the most visited ones, by 2030.

Simultaneously, the “Ex - ante impact assessment: Report on a European collaborative cloud for cultural heritage” [11], published in May 2022, lay the groundwork for the preparation of the Cluster 2 WP Calls to support the European Collaborative Cloud for Cultural Heritage for the period of 2023-2025. The scope of the basic European Collaborative Cloud for Cultural Heritage platform will be therefore to offer possible higher-level tools and instruments for storing, accessing, using, and documenting digital twins, for supporting the activities of the digital continuum and for connecting all actors in the digital ecosystem. As stated in the report: *“The EC-funded competence centre 4CH could play an important role in the design of a European Collaborative Cloud for Cultural Heritage, since it could contribute to the identification of competences, requirements, and potential integration of existing services; it could also help in selecting experts who might work on the European Collaborative Cloud for Cultural Heritage assessment of tools and services. Therefore, a collaboration with 4CH would be opportune.”*

The 4CH project wants indeed to act in that direction and contribute as much as possible: as foreseen in the DoA, the 4CH project is going also to design and implement a cloud-based infrastructure for services that could be managed by the future Competence Centre on Cultural Heritage, enabling the CH community (including software providers) to share competences and information, and providing a distributed platform to run applications. Most of these needs of the CH community are currently partially satisfied by the various solutions analysed above. In such respect, one of the objectives of the present deliverable is indeed to analyse several tools and solutions as a good candidate for the design and the definition of the 4CH Cloud Platform. The next sections are going in this direction by providing details on the future architecture of the 4CH Cloud Platform pilot and related services adopted, and on the other hand, policies, best practices, and requirements to integrate and federate CH tools and applications.

5. The architecture of the 4CH Cloud Platform

The initial and fundamental activity has been focused on the design and the definition of the 4CH project platform to host the services suitable for cultural heritage study and provided by the 4CH project members.

The challenge has been to design an architecture, which contains all the elements needed to provide users with the capability of using different software services made available via a platform in order to let them be able to use resource infrastructures in a seamless and transparent way (see Figure 5 for an overview of the different cloud layers). The current technology based on lightweight containers and related virtualization developments make it possible to design such Platform as a Service (PaaS) layer in a relatively straightforward way. As a first step, the activities were focused on providing compute, storage, and networking resources (IaaS) in a standardized, reliable, and performing way to remote customers.

For such reason, a cloud-based infrastructure has been chosen and made available by one of the project partners (INFN) which has a long-term experience in distributed computing, and it has been involved in many e-infrastructure EU-funded projects such as (i) the DEEP Hybrid DataCloud (DEEP-HDC) project [12], which aimed to bridge together cloud and intensive computing resources in order to explore different datasets for artificial intelligence and deep and machine learning, (ii) the eXtreme-DataCloud (XDC) project [13], which focused on developing scalable technologies aimed at federating storage resources and managing data in highly distributed computing environments, (iii) the INDIGO-DataCloud project [14], which centred on developing a computing platform that can be deployed on different hardware and provisioned over hybrid infrastructures.

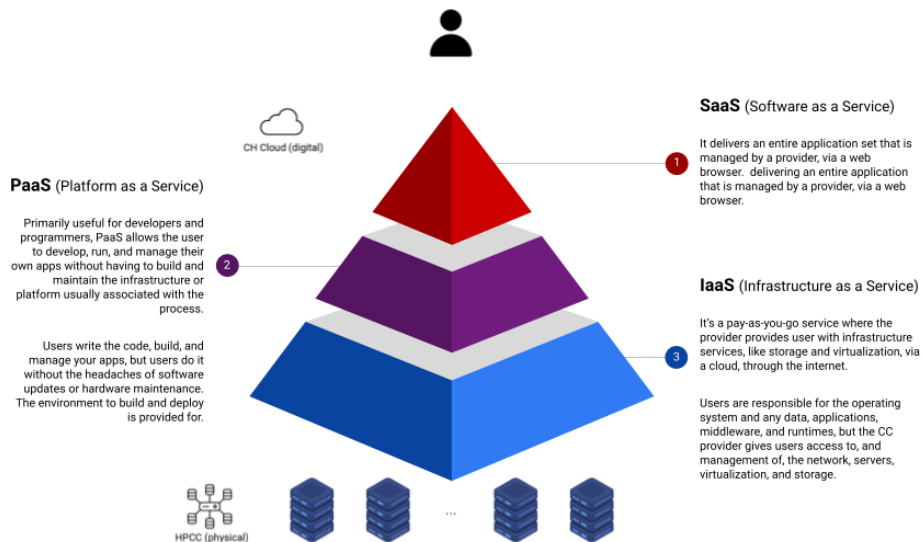


Figure 5: Cloud service layers

5.1 The 4CH Cloud Platform: design and approach

The 4CH project infrastructure is a cloud-based infrastructure. The resources used to host the 4CH Cloud Platform, in fact, are provided by the CLOUD@CNAF infrastructure, based on Openstack⁵, present at INFN-CNAF, the main ICT centre of the National Institute for Nuclear Physics in Italy⁶.

In particular, to support the development of project activities, some IaaS cloud resources, based on OpenStack, have been made available.

Those resources will be increased in the deployment phase of the 4CH Cloud Platform (indicated in the text also as 4CH Platform). Nevertheless, the 4CH platform development approach adopted the so-called Infrastructure-as-Code⁷ (IaC); IaC is the process of managing and provisioning computing resources through machine-readable definition files, rather than physical hardware configuration or interactive configuration tools; this allows for a faster, easier deployment of the CC cloud platform once the final High Performance Computing (HPC) resources will be made available for the production phase.

Following the IaC approach, the Kubernetes solution has been adopted and implemented. Kubernetes⁸ is an open-source container orchestration system for automating software deployment, scaling, and management.

By defining a set of building blocks ("primitives"), Kubernetes can collectively provide mechanisms that deploy, maintain, and scale applications based on CPU, memory, or custom metrics. Thanks to its loosely coupled nature, Kubernetes is extensible enough to meet different workloads. The internal components as well as extensions and containers that run on Kubernetes rely on the Kubernetes API. The platform exerts its control over compute and storage resources by defining resources as Objects, which can then be managed as such.

In such respect, to enhance the use of resources and orchestration of different services, a Kubernetes cluster has been created. The cluster is composed of a master and several worker nodes suitable to host the 4CH services. To expose services and applications to the users, an ingress object has been also made available.

The platform was deployed using the Rancher Kubernetes Engine (RKE)⁹, version 1.3.11 which is based on Kubernetes version 1.23.6.

On top of the Kubernetes cluster, some components have been deployed to improve the cluster functionalities. From one hand, network service relies on Calico¹⁰, from the other storage service rely on Longhorn¹¹.

A graphic schema of the 4CH Cloud Platform is shown in Figure 6.

⁵ Openstack, <https://www.openstack.org/>

⁶ INFN CNA, <https://www.cnaf.infn.it/>

⁷ IaC, https://en.wikipedia.org/wiki/Infrastructure_as_code

⁸ Kubernetes, <https://kubernetes.io>

⁹ Rancher Kubernetes Engine (RKE), <https://rancher.com/docs/rke/latest/en/>

¹⁰ Project Calico, <https://www.tigera.io/project-calico/>

¹¹ Longhorn, <https://longhorn.io/>

Calico is an open-source networking and network security solution for containers, virtual machines, and native host-based workloads which provides the network connectivity internal and external to the Kubernetes cluster.

Calico extends the built-in API of Kubernetes and is responsible for adding or deleting pods to/from the Kubernetes pod network, including creating/deleting each pod's network interface and connecting/disconnecting it to the rest of the network implementation.

Longhorn is an open-source software that implements distributed block storage using containers and microservices. Longhorn creates a dedicated storage controller for each block device volume and synchronously replicates the volume across multiple replicas stored on multiple nodes. The storage controller and replicas are themselves orchestrated using Kubernetes. Moreover, to increase availability Longhorn creates replicas of each volume. Replicas contain a chain of snapshots of the volume, with each snapshot storing the change from a previous snapshot. A graphical overview of Longhorn is given in Figure 7.

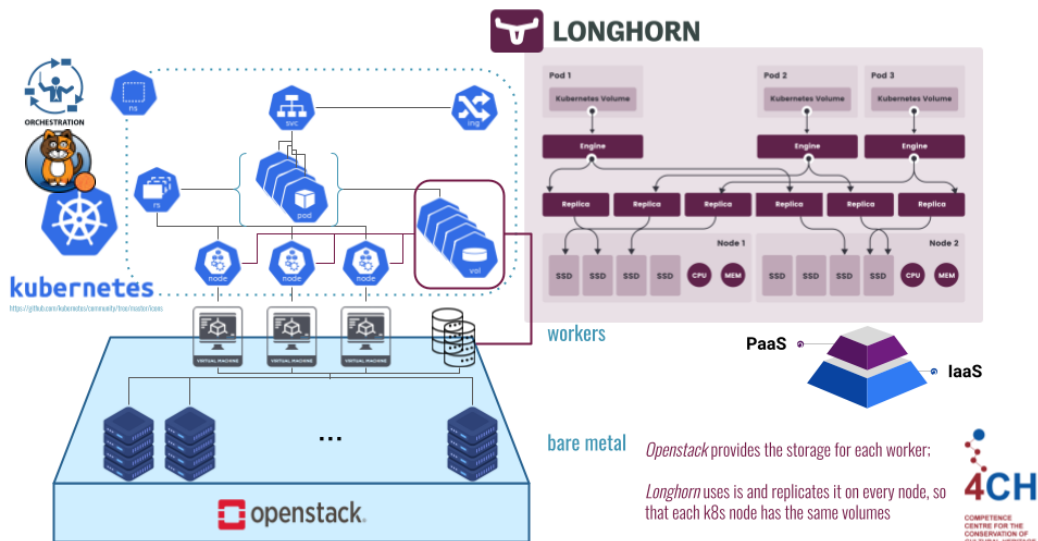


Figure 6: 4CH platform design and components

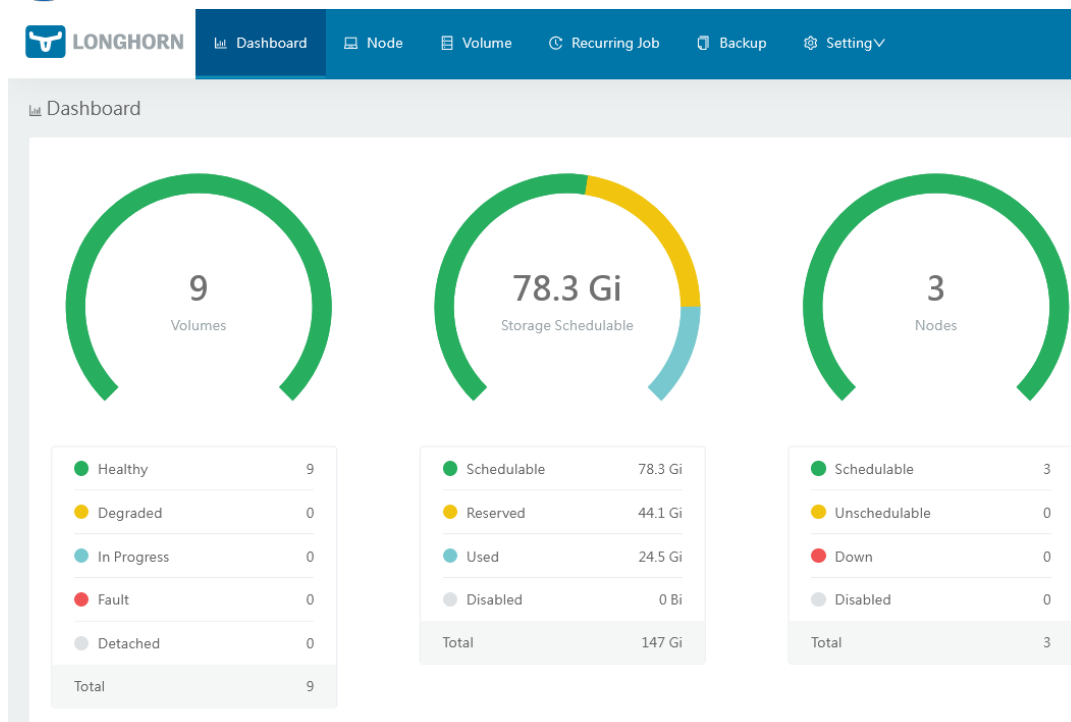


Figure 7: Longhorn service deployed on 4CH platform

5.2 4CH Platform services

The 4CH Platform comes with some services needed to maintain the functionalities of the platform itself. All the internal services hosted by the 4CH Platform are provided as container and are orchestrated via Kubernetes (see Figure 8 for details).

Hereafter, a description of the services deployed is provided and a complete list of endpoints related to the 4CH platform is made available on Table 5.

Table 5: 4CH Platform services, related application versions and endpoints

4CH Service	Application	Version	Endpoint
Authentication/Authorization	INDIGO-IAM	1.7.2	https://chnet-iam.cloud.cnaf.infn.it
Platform Monitoring	Grafana	8.5.3	https://ch.cloud.cnaf.infn.it/grafana
Image repository	Harbor	2.5.1	https://ch.cloud.cnaf.infn.it:32254/harbor
Longhorn	Longhorn	1.2.4	https://ch.cloud.cnaf.infn.it/longhorn/
CHNet web Portal	Nginx	openresty 1.21.4.1	https://ch.cloud.cnaf.infn.it/
Assisted Metadata generation and data retrieval	THESPIAN-Mask	0.2	https://ch.cloud.cnaf.infn.it/metadatamask/
Online visualizer for X-ray fluorescence raw data	THESPIAN-XRF	0.1	https://ch.cloud.cnaf.infn.it/XRF/

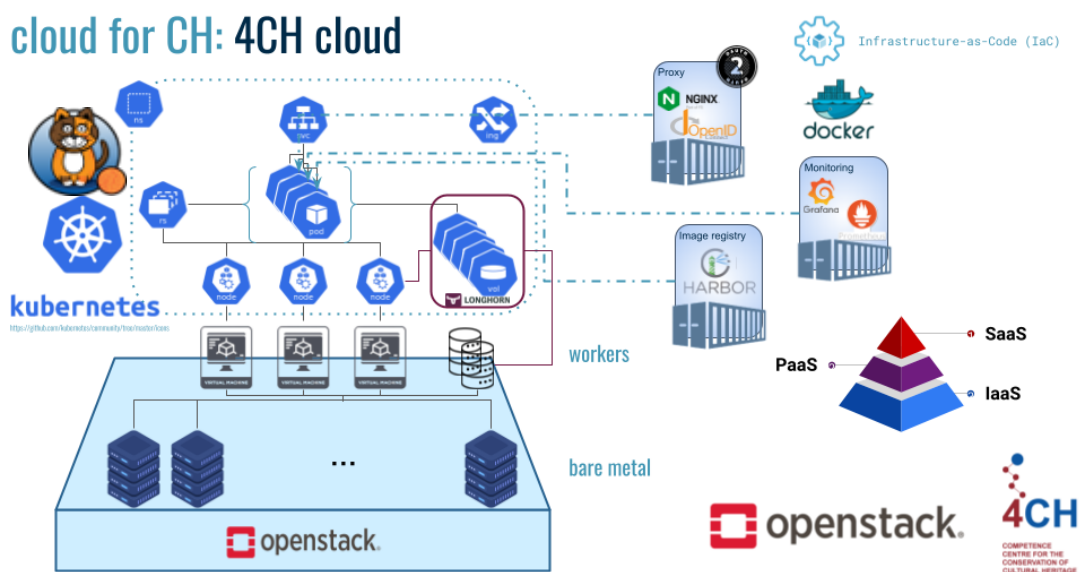


Figure 8: 4CH Platform and ancillary services

5.2.1 Authentication / Authorization workflow

The 4CH Cloud Platform rely on INDIGO-IAM¹² service which provides the authentication and authorization mechanisms and deployed for the Cultural Heritage Network (CHNet) community. INDIGO-IAM implements (i) the OAuth2.0¹³ standard authorization framework with the OpenID Connect¹⁴ (OIDC) layer, (ii) the User-Group model to manage the authorization procedure.

Moreover, the 4CH platform relies on Nginx-proxy¹⁵ that acts as reverse proxy which intercepts users' requests and redirects them to the correct service. Nginx-proxy takes care also of the authentication mechanism provided via INDIGO-IAM.

Nginx-proxy includes a reverse proxy exposed both to general internet and to the deployed applications. The reverse proxy intercepts users' requests and redirects them to the correct service. Currently, the reverse proxy is Nginx-based but we are also evaluating other solutions.

¹² INDIGO-IAM, <https://github.com/indigo-iam/iam>

¹³ Oauth2.0, <https://oauth.net/2/>

¹⁴ OIDC, <https://openid.net/connect/>

¹⁵ Nginx-proxy, <https://docs.nginx.com/nginx/admin-guide/web-server/reverse-proxy/>

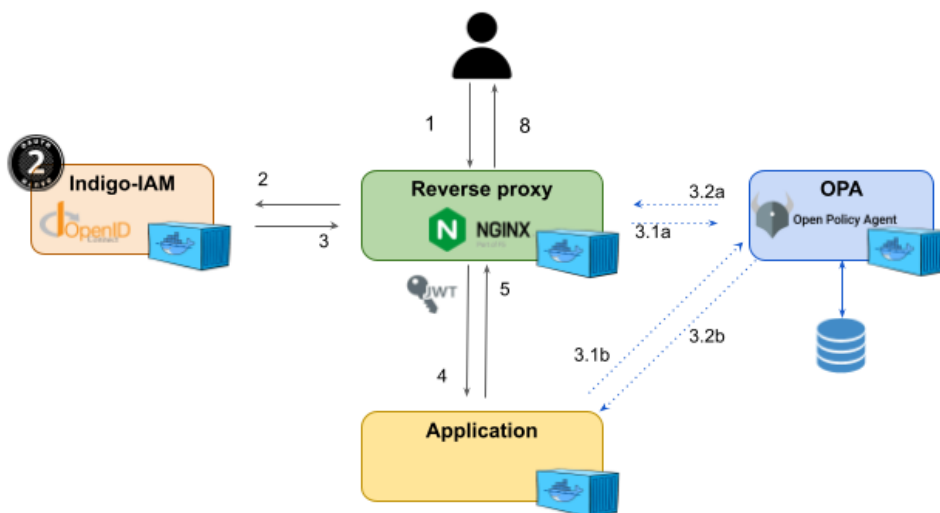


Figure 9: Authorization and authentication workflow for applications within the 4CH platform.

Users can authenticate to a proper application, exposed by the reverse proxy, only if they have an account on the INDIGO-IAM project instance.

As from Figure 9, users can access the web portal exposed by NGINX (1) where they can login. The authentication is managed by INDIGO-IAM (2) instance deployed for CHNet community. INDIGO IAM implements the OAuth2.0 standard authorization framework with the OpenID Connect (OIDC) layer. If the authentication succeeds, an ID token and an access token are provided (3). The ID token, which is in JWT¹⁶ format, according to the standard, can be used by the application to obtain identity and other information about users' membership. The access token, again in JWT format, it is forwarded by the reverse proxy to the upstream application (4) as an access key to the resources. If the token is valid, the application sends back the resource to the reverse proxy (5) which forwards the response back to the users (6). As a further development, the introduction in the workflow of the Open Policy Agent (OPA)¹⁷ is foreseen. The OPA service will provide the needed authorization rules, which are expressed in the Rego language. The OPA service can be queried either by the Nginx reverse proxy (3.1a, 3.2a) or directly by the application (3.1b, 3.2b), presenting the access token.

5.2.2 Monitoring

The 4CH Platform provides a service to monitor the status of the platform.

Grafana¹⁸ and **Prometheus**¹⁹ have been deployed to monitor both the Kubernetes cluster and the different services running on top of it. Grafana provides a graphical dashboard where both capacities and usage of virtual resources and services are shown. Grafana is a multi-platform open-source analytics and interactive visualization web application. It provides charts, graphs, and alerts for the web when connected to supported data sources. Prometheus is a free

¹⁶ JSON Web Tokens, <https://jwt.io/>

¹⁷ Open Policy Agent, <https://www.openpolicyagent.org/>

¹⁸ Grafana, <https://grafana.com/>

¹⁹ Prometheus, <https://prometheus.io/>

software application used for event monitoring and alerting. It records real-time metrics in a time series database (allowing for high dimensionality) built using a HTTP pull model, with flexible queries and real-time alerting.

In the present deployment, Prometheus is collecting metrics from the 4CH platform and made them available to be visualized with Grafana. In Figure 10 there is a graphical example of the metrics collected by Prometheus and visualized by Grafana.



Figure 10: 4CH platform monitoring using Grafana and Prometheus.

5.2.3 Image repository

Another, important service provided in the 4CH platform is the image repository based on **Harbor**²⁰ (see Figure 11). Harbor, in fact, is an open-source software aimed at storing the images related to the different services running in the 4CH Platform for an easy redeployment and to have up-to-date images available. Harbor has been integrated with INDIGO-IAM, so specific members of the project registered in IAM can be authenticated also to push images to Harbor, in particular in the implementation of Harbor for the 4CH project. Moreover, Harbor is also providing statistics and logs that can be used for maintenance and accounting scopes.

²⁰ Harbor image repository, <https://goharbor.io/>

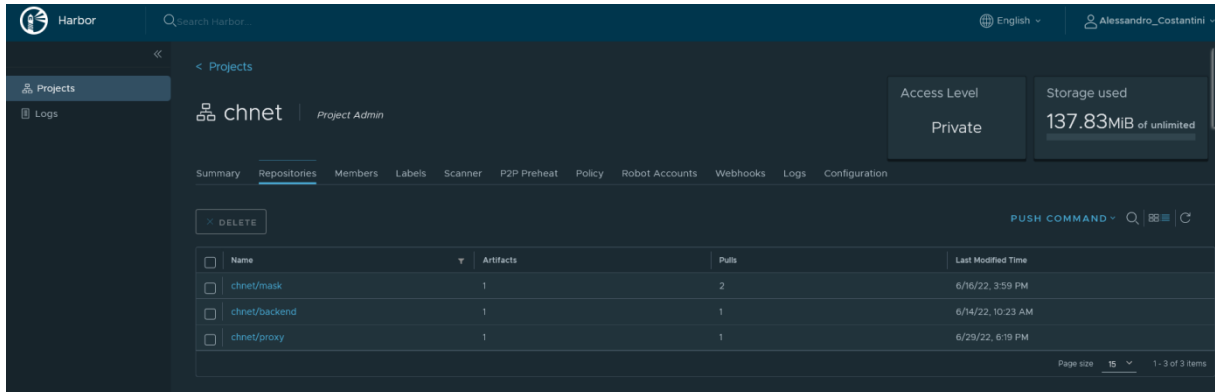


Figure 11: 4CH platform image repository using Harbor.

5.2.4 Web portal homepage

To access the Cultural Heritage applications made available by the 4CH project partners, a web portal will be developed and integrated with the INDIGO-IAM service to provide authentication and authorization mechanisms. At present, a previous web page is used for testing the access system. It will be replaced by a proper 4CH web page at the time of launching the 4CH cloud and its services.

The Web portal will provide the information required to access the services and the KB, together with the available services and relevant links to the project web site.

5.2.5 Developing services

As a further development, the following implementations and services are foreseen:

- High Availability of the platform services (Kubernetes controller, ETCD, Ingress controller) in order to improve resilience of the deployed services;
- Sync&Share service offered to the 4CH users to easily share documentation and user data.

6.4CH Platform policies and requirements: integration and federation

6.1 Open-source advantages and possibilities

To foster the guiding principles expressed by the Tallin Declaration on e-Government [15,16], and since the 4CH Cloud Platform is currently hosted by a public research institution (INFN, which is also a partner of the 4CH project), a set of regulations related to the deployment, adoption and use of commercial software should be applied to the whole lifecycle of the services that have to be integrated in the 4CH Cloud Platform. In particular, the integrated applications should embrace the open-source, open-access policies.

The platform described in Section 5 satisfies the open-source open-access principles, since the IaC was developed employing open-source software, frameworks, and tools. Additionally, all the web applications already integrated, offering both the Platform-as-a-Service (PaaS) and Software-as-a-Service (SaaS) stack of the 4CH Platform, were compliant with these principles and policies.

Those principles and policies have been considered to define the requirements related to the application deployment, method and resource exposed and authorization and authentication mechanisms.

Such requirements are presented hereafter, where different application integration use cases are provided and described.

6.2 Standard, Best Practices and conventions

6.2.1 Implementation Standards

The software products coming from the CH partner that is hosted to the 4CH Cloud must follow the Configuration and Integration Policy to build and package the products. In particular, the project will adopt the well-known packaging guidelines and policies defined by Docker for containers [17].

6.2.2 Documentation standards

Technical Documentation as any document providing information on the behaviour and interfaces of the 4CH products algorithms, APIs, containers, VM images, etc. will be provided by the project. The Technical Documentation includes requirements description and analysis, architectural and design documents and specifications of interfaces and APIs. The Technical Documentation will be maintained and kept updated by the project partners.

6.2.3 Software components and classification

A general architectural description of all 4CH services, their interfaces and protocols and their interactions and dependencies is available in Section 5.

4CH classification for software components can be classified as:

- **Service:** provides a well-defined set of features (behaviours or capabilities) through a published interface. Both interfaces and behaviour must be publicly documented, supported and maintained and are subject to public transition and lifecycle processes.
- **Client:** provides capabilities to interact with the services through their published interface and a separate user interface accessible with command-line or graphical commands. The user interface must be publicly documented, supported and maintained and are subject to public transition and lifecycle processes.
- **Library:** provides capabilities to interact with the services through their published interface or implements a set of tools and utilities. It provides a programmatic interface with bindings for one or more programming languages. The programmatic interface must be publicly documented, supported and maintained and is subject to public transition and lifecycle processes.
- **Internal component:** a sub-element of a service, client or library that does not expose its interfaces to users or external programs. Its interfaces must be documented but are not subject to public transition or life-cycle processes and can be changed at any time provided the change does not introduce changes in the published interface or behaviour of services, clients, and libraries using them.
- **Virtual images**²¹ (container images²² and VM images²³): Containers use Operating-system level virtualization a server-virtualization method where the kernel of an operating system allows for multiple isolated user-space instances, instead of just one. A container encapsulates everything that is needed for some code to be run, such as the code itself, its dependencies, system libraries, etc. VM images, are the conventional virtual machine images that contain a complete operating system, plus applications and services.

6.3 Requirements for service integration

Applications providing services for Cultural Heritage applications can be integrated in the 4CH Platform and made available to the community via the 4CH dashboard as from Figure 12.

To be integrated in the 4CH platform, the application (or a set of applications) should agree with the following requirements:

- 1) The application must expose a REST API [18] or a reachable endpoint and must be configured to be proxied to a proper location (e.g., `https://<application_domain>:<application_port>/<application_prefix>`, or using proper variables such as, `HOST=https://<application_domain>`,

²¹ A virtual image is the abstraction of a computer disk or media to be used in a Virtual Machine as well as in a Container able to provide a fully configured/configurable virtual environment.

²² Docker overview, <https://docs.docker.com/get-started/overview/>

²³ Virtual Machine image, https://en.wikipedia.org/wiki/Virtual_machine

BASE_URL=/`<application_prefix>`/, PORT=`<application_port>`). All the application resources (URL) must have the same path URL.

- 2) The resources exposed by the application must be labelled as free-access or authentication protected. In the first case the application exposes free access content. In the latter, the content must be accessed after authentication, which is managed by INDIGO-IAM. The access token in JWT format is forwarded by the reverse proxy to the application as an access key to the resources. Without the JWT access token, the access must be rejected.
- 3) The application (or the applications) must be containerized and the image, based on standard images available on the common image repositories, must be made available.
- 4) The image must respect the common security aspects and must be compliant with the security best practices. Moreover, the image must be kept updated, in respect of the Common Vulnerabilities and Exposures²⁴ which identify, define, and catalogue publicly disclosed cybersecurity vulnerabilities.
- 5) Authentication mechanism:
 - a) Could be demanded to the 4CH reverse proxy. In this case all the communications must be performed from and to the reverse proxy itself.
- 6) Authorization mechanism
 - a) Could be demanded to the application: in this case, the application must use the JWT, by extracting from it the needed user information to be used for authorization proposes.
 - b) If the point above (3.a) is not applicable, authorization could be managed by the Open Policy Agent (OPA) internal service that will be integrated soon (see Section 5.2.1). In this case, a set of authorization rules must be provided.

Moreover, to be integrated in the 4CH platform, the application (or a set of applications) should agree with the following policy:

- 7) The Operating System and related components within the container should be compliant with the national and EU regulations and recommendations related to the adoption and use of Open Source and commercial software.

6.3.1 The integration use case: THESPIAN – Mask

To illustrate the integration and federation approaches we refer to a set of tools already created to manage the results of scientific analyses of heritage assets, called THESPIAN Mask, and to visualize the results of an XRF analysis, called THESPIAN-XRF.

The Tool for HEritage Science Processing, Integration and Analysis (THESPIAN) Mask [19] is a service for assisted standardised metadata generation and extraction. It is based on an ad-hoc devised ontology, CRMhs [20], extension of the CIDOC-CRM ontology [21].

²⁴ CVE® Program Mission, <https://www.cve.org/>

THESPIAN-Mask consists of the following services²⁵:

1. A frontend, hosting an endpoint furnishing the HTML²⁶, CSS²⁷, and JavaScript²⁸ bundle to be rendered in the client browser **{Req 1}**;
2. A backend, offering a set of RESTful APIs contactable via HTTP/2²⁹ protocol, for metadata/data upload and retrieval, and offering a set of tools, based on third parties API, for metadata standardisation **{Req 1}**;
3. A Database (DB), based on MongoDB³⁰ No-SQL³¹ database, for storing the metadata entries.

A visual representation of the service is available in Figure 13.

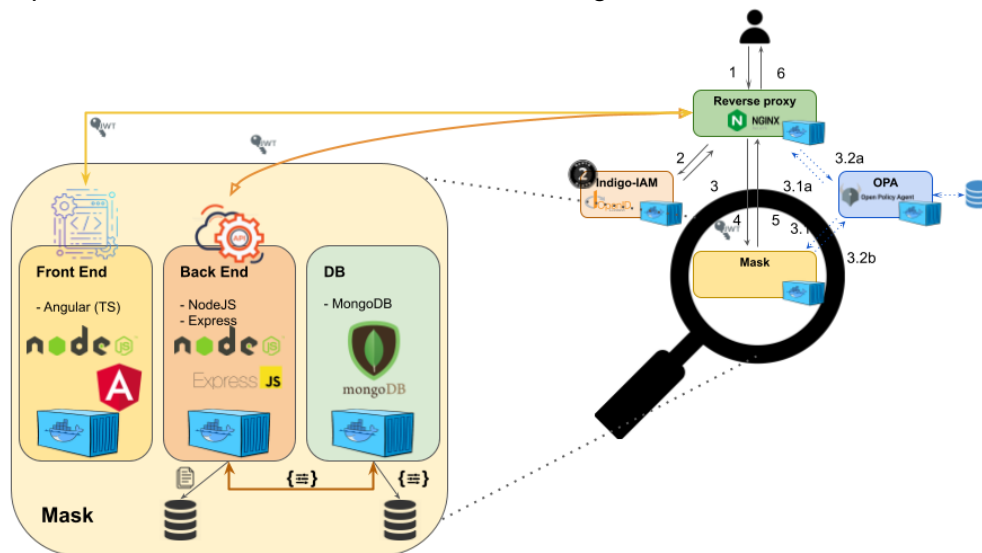


Figure 13: THESPIAN-Mask application integration in the 4CH Cloud Platform

Each application adopts open-source software **{Pol 1}** and is containerised by using an official image **{Req 3}** as following:

- Frontend and backend are based on the official alpine image³² based Node.js³³.
- Database is based on the MongoDB official image³⁴.

²⁵ Hereafter, the requirements listed in Section 6.2 that are satisfied by the application are highlighted with the curly brackets, e.g. **{Req 1, Pol 1}**.

²⁶ HyperText Markup Language (HTML), <https://html.spec.whatwg.org/multipage/>

²⁷ Cascading Style Sheets (CSS), <https://www.w3.org/Style/CSS/>

²⁸ JavaScript, <https://www.javascript.com/>

²⁹ HyperText Transfer Protocol 2 (HTTP/2), <https://datatracker.ietf.org/doc/html/rfc7540>

³⁰ MongoDB, <https://www.mongodb.com>

³¹ No-SQL Database, <http://nosql-database.org/>

³² Node.js official image on Docker Hub: https://hub.docker.com/_/node

³³ Node.js, <https://nodejs.org>

³⁴ MongoDB official image on Docker Hub, https://hub.docker.com/_/mongo

The official Node.js image has been enriched by allowing TypeScript compilation for the Angular³⁵ framework, following the Docker best practises **{Req 4}**.

All images adopt environment variables that can be passed to the containers where they are running. Moreover, the used and enriched images have been tagged and pushed on the 4CH image repository service based on Harbor and made available by the 4CH Cloud Platform. Backend and database use persistent volumes and the related volume location is set by an appropriate environment variable.

To ensure files preservation, the volumes are not shared among the applications.

The connections between frontend and backend are mediated by the reverse proxy **{Req 5, Req 6}**, this means that only authorized and authenticated users can access the services. By adopting a common best practice in containerization, the DB can only be reached from the backend, via an internal (virtual) network.

The third-party APIs are contacted by backend. The request formulated by the frontend uses HTTP/2 protocol and is sent to the backend via the reverse proxy through an established connection. The backend contacts the 3rd-party APIs, get the replies and replies via the reverse proxy to the frontend.

6.4 Requirements for service federation

Applications that cannot be integrated in the 4CH Cloud Platform can be hosted outside and can be made available to the community via the web portal.

By adopting this approach, Cultural Heritage applications that have specific needs (e.g. application cannot be virtualized, in house application with particular dependences, commercial software) can in any case be made available to the final user via the 4CH web portal adopting the authentication and authorization mechanism provided by the 4CH Platform.

For such reason, the application (or a set of applications) must respect the following requirements:

- 1) The application must expose a REST API or a reachable endpoint using SSL. The application must be configured to be proxied to a proper location (e.g., `https://<application_domain>:<application_port>/<application_prefix>`, or using proper variables such as, `HOST=https://<application_domain>`, `BASE_URL=/<application_prefix>/`, `PORT=<application_port>`). All the application resources (URL) must have the same path URL.
- 2) The resources exposed by the application must be labelled as free-access or authentication protected. In the first case the application exposes free access content. In the latter, the content must be accessed after authentication, which is managed by INDIGO-IAM. The access token in JWT format is forwarded by the reverse proxy to the application as an access key to the resources. Without the JWT access token, the access must be rejected.
- 3) Authentication mechanism:

³⁵ Angular framework, <https://angular.io/>

- a) Could be demanded to the 4CH reverse proxy. In this case all the communications must be performed from and to the reverse proxy itself.
- 4) Authorization mechanism
- a) Could be demanded to the application: in this case, the application must use the JWT, by extracting from it the needed user information to be used for authorization proposes.
 - b) If the point above (3.a) is not applicable, authorization could be managed by the Open Policy Agent (OPA) internal service that will be integrated soon (see Section 5.2.1). In this case, a set of authorization rules must be provided.

6.4.1 The federated use case: THESPIAN-XRF

THESPIAN-XRF is a web service for analysis and visualisation of X-Ray Fluorescence data, stored as standard ISO HDF5 format³⁶.

The application is hosted on the LABEC node of INFN-CHNet and can be accessed only through the 4CH Cloud Platform by a connection handled by the reverse proxy {Req 3, Req 4}.

Gunicorn³⁷ and WSGI³⁸ server offer a dynamic content with responsive endpoint {Req 1}, by using HTTP2 request to handle the computations. Also in this case, the application is containerised (even if it is not mandatory and is not listed in the requirements) to improve both development and operations.

A schematic representation of the service, and its federation with the 4CH Cloud Platform, is shown in Figure 14.

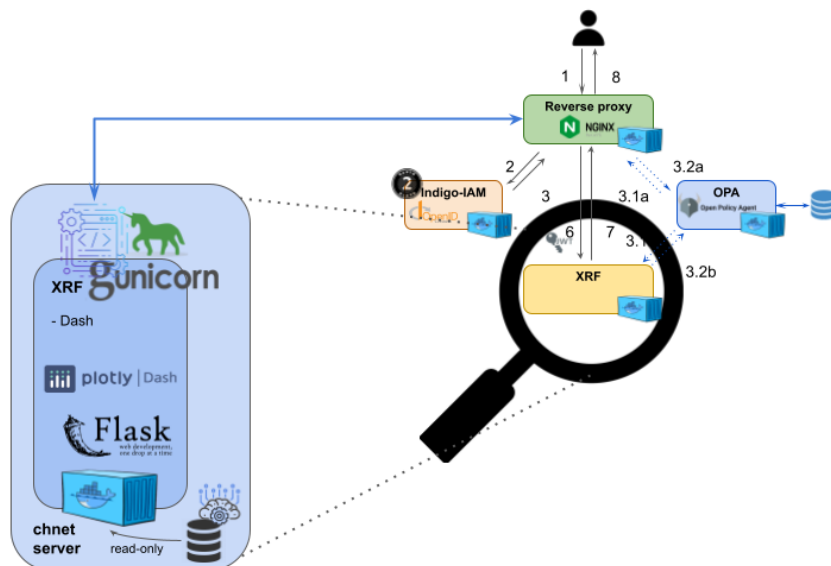


Figure 14: THESPIAN-XRF application federation in the 4CH Cloud Platform

³⁶ ISO HDF5 format, <https://www.loc.gov/preservation/digital/formats/fdd/fdd000229.shtml>

³⁷ Gunicorn, <https://gunicorn.org/>

³⁸ Web Server Gateway Interface (WSGI), <https://wsgi.readthedocs.io/en/latest/>

6.5 Mixing service Integration and Federation

On the sections above, examples of Integrated (Section 6.2) and Federated (Section 6.3) deployments of applications in the 4CH Cloud Platform have been shown (see Figure 15 for details).

Note that a single application, comprising multiple independent services, may be deployed in a mixture of federation and integration deployments. As an example, we report another way the THESPIAN-Mask service could have been deployed.

In this scenario, the THESPIAN-Mask Frontend is deployed as an integrated service within the 4CH Cloud Platform, while the Backend and the database are deployed on a remote server and federated (see Figure 16 for a visual depiction). Having mixed deployments arises from different needs: starting from (some of the) services which benefits of specialized hardware (i.e., HPC, GPU, etc.) that can be made available both in the 4CH Cloud Platform as from outside the platform, to the fact that (some other) services cannot fulfil the requirements to be integrated in the 4CH Cloud Platform but can still be federated and made accessible to the community.

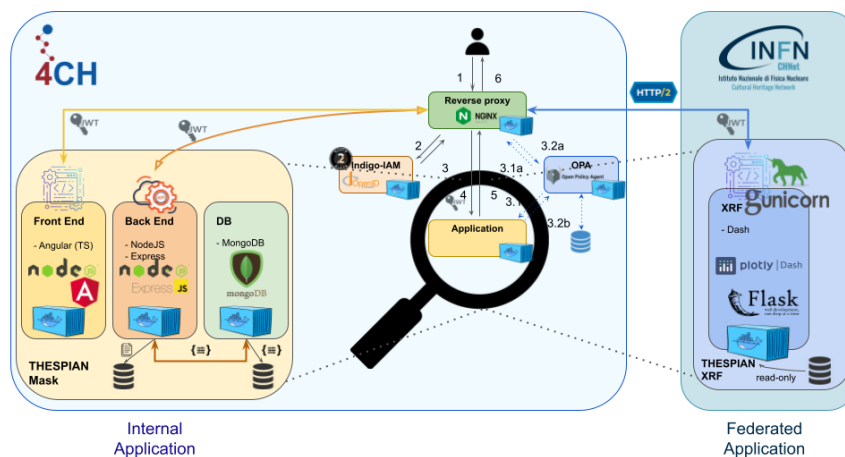


Figure 15: Example of Integrated and Federated applications in the 4CH Cloud Platform

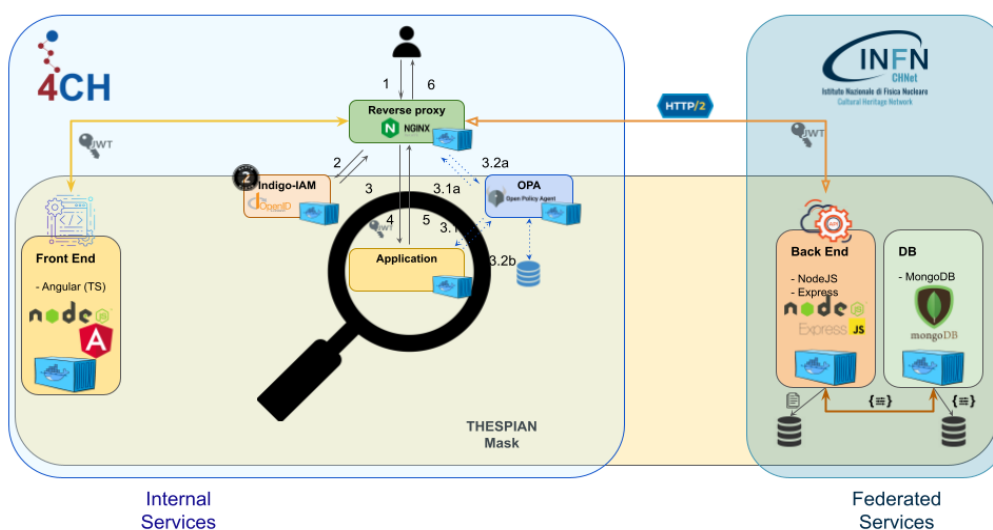


Figure 16: Example of hybrid application deployment in the 4CH Cloud Platform

7. The 4CH Knowledge Base

7.1 Introduction to the Knowledge Base

The 4CH Knowledge Base (KB) will collect and organize all the information pertaining to cultural heritage as far as the 4CH scope is concerned. Of course, it is not a mammoth repository where all such data are stored: it is instead an intelligent catalogue of such information, organizing into it the metadata of datasets stored in different systems and relevant for the 4CH topical scope, which at present form a fragmented global universe of silos disconnected from each other. Just to give an example, such datasets are currently stored into local repositories of research centres working on heritage sciences; into large repositories managed by heritage agencies for management purposes, such as the system managed by heritage ministries or agencies in different countries, as well as into smaller ones created by private organizations (for example, Churches) owning and managing heritage assets; and so on. It is anticipated that by progressively integrating such metadata a global network of heritage-related information will be set up. In parallel, such system will be integrated with global technical information concerning heritage management and conservation, for example restoration techniques and material reference collections, and linked to Big Data sources, as the ones concerning risks like earthquakes, floods, atmospheric events, pollution and so on. The creation of the KB will progressively proceed according to dataset availability and their owner's willingness to collaborate.

A similar concept is present in Europeana and will be hopefully exploited in the future Data Space for Cultural Heritage. However, the purpose of this system is rather different from the 4CH one under several aspects. Firstly, the metadata they currently collect are not suitable to an advanced use such as the one envisaged by 4CH. Secondly, the quality of the data they organise is dubious, what may prevent researchers to use Europeana for their work. According to the Annex to the Recommendation issued by the European Commission on cultural heritage digitization [10], only about 60% of Europeana data are currently "High Value Datasets" and even for them it is not known how much they pertain to the 4CH scope – most likely not too many. It is hoped that the new undertaking will improve these indicators, but until then most researchers and heritage managers will probably refrain to use Europeana data for their work. Relatively to the archaeological domain, the ARIADNE series of projects have created an aggregation system of about 3 million datasets. While in this case almost all of them are qualified as high-quality research data, there are sectors in which the data availability is less abundant than in others. This is for example the case of scientific analyses and of 3D data. The latter are much underrepresented compared to the production of 3D models in research because of the difficulty of storing them and making such repositories available for consultation.

To set up such a system, it is necessary to create a theoretical background in which data are organized, the 4CH ontology, described in section 7.2; to set up a procedure for aggregation, outlined in section 7.3; and to manage a cloud-based information retrieval system based on such ontology, summarily sketched in section 7.4, which also introduces a general framework for virtual research environments and virtual management environments, to address both the research and the management needs of the heritage community as far as the 4CH topical scope is concerned.

As concerns potential data providers, in France *Huma-Num* [22] is a very important framework for humanities-related data. It includes sections for archaeology (MASA) and tangible heritage assets. The large-scale *Espadon* French project [23] aims at setting up a framework for heritage sciences data. Belgium has created at KIK-IRPA (Royal Institute for Cultural Heritage) the *HESCIDA* project [24], a national database for heritage sciences. The Netherlands National

Academy of Sciences (KNAW) manages *DANS* [25], a Dutch national repository for research data with sections for heritage, archaeology in particular. Such national repositories are just examples of potential contributors to the 4CH KB.

7.2 Organizing the data: the HDT ontology

It is apparent that powerful semantic tools are required to deal with heritage data in the way the 4CH goal – and hence its KB – requires.

For this reason, a task force from Task T3.4 is working to the definition of a novel ontology based on the *digital twin* concept. This ontology is named the ***Heritage Digital Twin (HDT) ontology***, from the name of its key concept, the ***Heritage Digital Twin***³⁹.

The concept of digital twin is not new. A digital twin can be broadly defined as a *virtual representation that serves as the real-time digital counterpart of a physical object*. Digital twins have been applied in many industrial fields where the idea was born to test components, devices and, later on, to simulate the real behaviour of complex appliances in a digital way [27]. Then, digital twins made their way in machinery control applications, by using sensors surveying the behaviour of a device and sending an alert when an anomalous value is measured, or directly activating specific components to return to a normal condition. This kind of applications required the use of simulation processes within the model, which were eventually incorporated in the concept of digital twin. Thus, for industry digital twins, the data component has a relatively simple schema, while the process part is more complex. The stress is on how the system behaves rather than on how the information about the system is structured. The European Commission has recently proposed to create a digital twin of the Earth [28] to evaluate complex environmental processes and their impact on the whole system, and to forecast the effects of mitigating measures against, for instance, global warming.

Digital Twins (DTs) are nowadays extensively used among others in mechanical engineering [27, 29], architecture and especially in the building industry, where they belong to the approach called BIM, already introduced in the previous sections, in particular to the HBIM one. From a semantic point of view, such a system is rather flat, with a root class in the drawings of the building and of its components, to which all the extra information is appended. A recent important project combining BIM and digital twins is the UK Gemini project, proposing to use digital twins nationwide for town planning [30]. It seems that at present it is very difficult to incorporate a much wider set of concepts – including some that are of a non-physical nature, the so-called intangible heritage – and relations among them.

The 4CH approach as envisaged so far goes beyond the above-mentioned limitations. From an ontological point of view, we define a ***Heritage Digital Twin*** as the complex of information about Cultural Heritage objects, embedded in a sophisticated digital system, operated within a virtual environment or digital platform, on which digital information concerning the replicated entities of the real world (e.g., 2D and 3D models, structured and unstructured documentation etc.) is stored, endowed with powerful analysis tools, capable of connecting and interacting with all the data and the models, of acquiring and processing external information, to provide real-time indications, simulate and predict future situations and conditions in the real world. The full technical definition of the related HDT ontology is still in preparation, but we can describe some features of its most important classes as follows. To recognise them, HDT class and property names are *in italics*, and in capital letters for classes. It must however be noted

³⁹ The remainder of this section is derived from a paper accepted for publication on DATA [26] and from the draft of another one in preparation for ORE.

that the names quoted here must be considered as interim ones, subject to change as the precision of the ontology description improves.

All HDT classes are considered as subclasses of an overarching one, *Heritage Entity*, which has no instances. The pivot concepts are the class *Heritage Asset*, corresponding to actual heritage assets of any nature (physical, both movable or immovable, immaterial or born digital), and the class *Heritage Digital Twin*, to indicate the whole of the digital information pertaining to Heritage Asset. The *Heritage Digital Twin* is related to the *Heritage Asset* via the property *is digital twin of (has digital twin)*.⁴⁰

Any instance of the *Heritage Asset* class may be related to other instances of the same class, for example being a part of them, as the tower of a castle, or forming with them a collection or a more general asset, as the paintings of an art gallery. Another recent example is the new UNESCO World Heritage site “The Great Spa towns of Europe” [31], which is made of eleven towns, each one developed around a spa, far from each other but belonging to a common cultural framework. Each one will have its own Heritage Asset instance, and there will also be a collective Heritage Asset instance representing the whole UNESCO World Heritage site.

Information about the parts forming an asset is relevant also to the whole, as it describes important details as, for example, the material of which the part is made. Thus, parts of a heritage asset may be considered as assets as well, and correspondingly the HDT of the whole will incorporate the HDTs of all the parts of the entire heritage asset which are identified as heritage asset on their own. For example, a church, i.e. the ‘main’ heritage asset, has many parts as chapels, paintings, architectural components, furniture, and so on, each one a heritage asset on its own and with an own heritage digital twin. The heritage digital twin of the church is then the assemblage of all these heritage digital twins, plus features concerning the whole, for example the style, the cult, the architect and so on. A more generic connection than part-whole, represented by a property *is part of (has part)*, is described by the property *is related to (relates)*, relating heritage assets to other assets.

The instances of all classes have an *Identifier*, i.e. a code attached to them, for example an inventory number; and may have one (or more) *Title*, often according to language, e.g. “Mona Lisa” (English), “Monna Lisa” or “La Gioconda” (Italian), “La Joconde” (French) etc., all referring to the same famous Leonardo’s painting. Thus, the properties *is identified by (identifies)* and *is titled (titles)* have *Heritage Entity* as domain, since identifiers may be defined also for digital twins, e.g. their URL, as well as names if one likes to do so for digital artefacts.

The *Heritage Document* class includes all the documentation items pertaining to a *Heritage Asset*. The documentation may consist in digital or analogic objects as printed or handwritten documents, old analogic photos, drawings, and so on; or digital, either born digital like digital photos or digitised from analogic one.

The property linking *Heritage Asset* to *Heritage Document* is the property *is documented in (documents)*, which may apply to the whole asset or to specific parts of it.

Among others, we may distinguish among the following ones:

- *3D Model*, resulting from any of the various techniques available as 3D scanning, wireframe modelling and so on. Practice will assess if it is more convenient to distinguish among them introducing different types, e.g. 3D point-cloud model, 3D CAD model and so on.
- *Imagery*, as photos and videos, but also special imagery as X-ray images, spectra of chemical and physical analyses, and so on. Among them, particularly relevant are the Virtual reality (VR) and Augmented Reality (AR) models, other types of visual digital artefacts pertaining to Heritage Asset. Both VR and AR models rely on 3D models of

⁴⁰ The text that follows is taken from the already mentioned DATA publication.

the related heritage asset, but may add or remove parts of it, or require further digital input as in AR, so they should be catalogued separately from 3D models. 3D models may correspond to actual objects – artefacts or built structures – or to conceptual ones, often representing the reconstruction of what is presumed to be the original configuration (and often, the use) of the reconstructed object. Similar models are named Virtual Reconstructions and are a commonplace in archaeology to communicate to researchers or, more frequently, to the public, the interpretation of past appearance. In general, Virtual Reality (VR) models enable virtual visits, while models incorporating the present appearance of heritage assets. i.e. Augmented Reality (AR) models can be viewed only on site, as they need the real time acquisition of the current asset appearance.

More types/subclasses may be introduced according to needs. For example, it may be worth considering *Conservation Document*, data about the asset conservation both in terms of the documentation of past interventions, the materials, and the analyses carried out on it. It consists in text files, numeric files (e.g. the results of analyses), images, videos, and special 3D objects, for example the results of tomography.

In general, different types of models may generate a subclass if they have special properties that apply only to them. Otherwise, it is simpler to consider their *Type* only, associated via the property *has type (is type of)*.

A *Heritage Asset* pertaining to a tangible asset *is located in* (is location of) a *Place*. The location may have different levels of precision, as a determinate position or area, a town, a region, and so on, for example “Room 1 of the Uffizi Museum”, “Athens”, “Cyprus”. It may also vary in time, if the related asset is moved elsewhere. A concept similar to location may be considered also for some intangible assets, as is done for some members of the UNESCO *Intangible Heritage Representative List* [32]. For others, it is intrinsic to the intangible asset to be on the move or to have no location, for example music. Thus, the (important) location property domain does not coincide with all the *Heritage Asset* class and is somehow different between tangible and intangible heritage: for the former it defines where the asset is located, while for the latter it indicates where the asset manifests itself. Therefore, the *Heritage Asset* class must split into the two subclasses *Tangible Asset* and *Intangible Asset*. For the former, *is located in (is location of)* specifies where the asset is placed. For the latter, the sister property is *happened at location (was location for)*. Both have *Heritage Location* as the range. Defining the location of a tangible asset may have various degrees of difficulty: if it is easy for most of them, in some cases it relies on research as there is no physical evidence confirming the location. This is the case of battlefields, for instance, when no finds or traces exists. An example is the Cannae site, where the battle (216 BC) between Romans and Carthaginians took place during the Second Punic War. Historical reports by Polybius and Titus Livius are available and enable a trustworthy identification of the location.

A *Heritage Asset* has many a *Heritage Story* that are associated to it. A *Heritage Story* includes any kind of witness related to the asset: it can be a narrative, a historical source, a popular attribution, co-created content and so on. A very special case concerns a literary itinerary as Leopold Bloom’s route through Dublin in James Joyce’s *Ulysses*, or – in a light-hearted perspective – real or almost imaginary places featured in popular novels or TV series, for example crime series as George Simenon’s *Commissaire Maigret*, in Paris; or Andrea Camilleri’s *Inspector Montalbano*, based in southern Sicily; both generating a substantial flow of ‘cultural’ tourism. In such cases, it is the *Heritage Story* that creates the *Heritage Asset*.

In general, a *Heritage Story* relates tangible heritage assets to their intangible components and to their reference communities. They are therefore of paramount importance also for the asset physical conservation and the safeguard of its intangible value. Before starting a conservation intervention or the evaluation of an activity from the heritage conservation perspective, it is necessary to consider the impact on the intangible component especially

when the impact on the tangible one is irrelevant. For example, fast food shops are often banned from historic centres, although their visual impact may be negligible; there are, instead, provisions to preserve the permanence of historic shops. Locating a MacDonalD's in the basement of a historic palace is unacceptable even if the building statics is unaffected. Ignoring these intangible aspects as it happens in BIM and HBIM models is a serious shortcoming for any heritage application of these approaches and one of the main reasons to propose a generalisation as the HDT ontology for the broader 4CH KB. A very nice example concerns Orthodox sacred icons. The devotion to a particular icon is manifested by lighting candles in front of it, which in time causes the blackening of the painting. Thus, the icon blackening level is the evidence of the believers' devotion and a confirmation of its religious value. Therefore, cleaning it would damage this intangible value: a common and much valued restoration practice for paintings would have an unexpected adverse effect if this intangible component is not considered. There are of course stories about intangible heritage as well, often in a much greater amount than for the tangible one.

The activity of producing a *Heritage Document* is called *Heritage Documenting*. It is characterised by the intentionality of creating new knowledge. We deliberately avoid using the term "documentation" in the class names as it has an ambiguous meaning, the activity and its result: "Documentation (i.e. the activity) produces documentation (the outcomes)". The term can be disambiguated only by the context, obviously unavailable in an isolated name.

On the other hand, a *Heritage Story* is the account of facts about a *Heritage Asset*, including but not limited to descriptions based on documents and on the interpretation of these documents. They are usually formulated in an attractive and accessible way to facilitate visitors' understanding and stimulate their curiosity. Frequently they avail of communication techniques as drawings, physical or digital reconstructions, and increasingly use VR and AR technology. If need be, it is very easy to extend the level of detail of classes to incorporate concepts with a very fine conceptual resolution. The previous discussion on 3D modelling, for example, may require a very detailed classification of the so far generic class *3D Model*, as already mentioned above.

A final, but important, note, about the HDT ontology is the consideration that it is fully compliant with the standard documentation system for cultural heritage, the CIDOC CRM model. Actually, the HDT ontology may be considered a compatible extension of the CRM. This guarantees interoperability with all the semantic infrastructures based on it, in practice all the cultural heritage repositories.

7.3 The aggregation pipeline

7.3.1 The mapping phase

To summarily describe the aggregation pipeline of datasets into the 4CH KB, it is possible to rely on the ARIADNE model which has proved to be efficient and successful in several years of use and after successfully aggregating 3 millions of disparate archaeological datasets.

The aggregation process starts with the mapping of the original ontology of the datasets onto the HDT. This part is developed through a collaboration between the dataset owner and the 4CH team. The tool used to describe the mapping and to implement it is X3ML [33], an open-source mapping tool developed and made publicly available at FORTH. According to FORTH, "following the X3ML definition, we developed and used an efficient transformation algorithm

which processes the declarative X3ML statements and produces equivalent RDF⁴¹ statements. Different cases of semantic heterogeneity that were encountered in different applications are covered [...] The aggregation of heterogeneous data from different institutions has the potential to create rich data resources useful for a range of different purposes, from research to education and public interests. Along this line, CCI has developed the X3ML Toolkit, a set of small, open source, microservices [...] designed with open interfaces and that can be easily customized and adapted to complex environments. The X3ML Toolkit consists of a set of software components that assist the data provisioning process for information integration. The key components of the toolkit are: (a) Mapping Memory Manager, (b) 3M Editor, and (c) X3ML Engine.” In the ARIADNE experience, X3ML has proved to be usable also by moderately computer-literate heritage professionals.

7.3.2 The cleansing phase

To be able to aggregate data from different sources, a first step is referring the content to a common set of vocabularies, gazetteers, and period gazetteers. To carry out this task, the Vocabulary Matching Tool developed in ARIADNE will be used. Such tool has the ability to define Exact, Close, or Broad matches between the Getty AAT [34] and target concepts in native vocabularies deployed by partners. A similar task will concern the mapping of named time periods to time intervals according to places: the term “Renaissance” corresponds to different time spans in Italy and Germany, for example. For this, the PeriodO [35] methodology will be used. ARIADNE has already developed several mappings for many institutions and languages that may be reused in 4CH.

7.3.3 The aggregation workflow

The aggregation workflow⁴² schematically consists of the following stages:

- 1) *Ingestion of XML records of the provider.* Records are collected, transformed into HDT by applying a dedicated 3M mapping, and fed into the knowledge graph.
- 2) *Enrichment with Getty AAT subjects.* The subject mapping to Getty AAT is transformed into RDF by applying the 3M mapping and fed into the knowledge graph. As a result, the knowledge graph will contain the correspondences between native subjects and terms of the Getty AAT vocabulary (or of the other chosen vocabularies).
- 3) *Enrichment with PeriodO.* The PeriodO datasets defined by the provider is ingested into the knowledge graph and used to generate explicit period properties.
- 4) *Feed the “staging” knowledge graph.* In order to support the providers at checking the content before it is made publicly available, the push on the knowledge graph initially targets a “staging” instance.
- 5) *Feed the staging 4CH KB portal.* For each set of records ingested with flow 1, a procedure reads the corresponding triples and builds JSON records to feed the index server that serves the 4CH KB portal. The goal of this flow is initially to generate records that can be used by the 4CH KB portal.
- 6) *Feed the public knowledge graph.* This flow is executed if the provider successfully completed the content checking on the staging knowledge graph and portal.

⁴¹ Throughout this document RDF means, as is well known, *Resource Description Format*. It must not be confused with the company named RDF, a 4CH partner to which one of the co-authors belongs, mentioned only in the list of co-authors.

⁴² This description is derived from the procedure successfully adopted by ARIADNE and published in its deliverable D5.2 [36].

7) *Feed the public 4CH KB portal.* The flow applies the same procedure of 5, using the public instances of the knowledge graph and portal instead of the staging instances.

7.3.4 Implementing the 4CH KB

To implement the 4CH KB a decision must be taken about the package to be used to manage the knowledge graph. So far, the INFN cloud adopts MongoDB while ARIADNE successfully used GraphDB, both open-source software. A comparative evaluation⁴³ of the different requirements, features and performances will be carried out as a desk analysis and possibly on a test dataset if necessary.

⁴³ See for example this comparison: <https://db-engines.com/en/system/GraphDB%3BMongoDB>

8. Conclusions

The present document described the activities carried out by the active tasks involved in WP3 and the connections with the other WPs in the project. Some of those activities will continue and progress in the next months as further improvements.

Some examples of continuing activities are:

- the HDT ontology will be completely described in semantic terms in Task 3.4 - Cultural Heritage Knowledge Base;
- the tools to be used to manage the creation and functionality of the KB will be analysed and implemented in joint work of Task 3.3 – Cultural Heritage Cloud and Task 3.4;
- the analysis on types of data and of the feature available with the software tools will be kept up to date and will contribute to the development of the Task 3.4 as a tool for orienting the CH community in the current wide and complex scenario;
- collaboration with the ongoing Task 4.2 - Standard and guidelines to CH digitization to assess controlled vocabularies for use in the 4CH KB;
- collaboration with the ongoing Task 4.2 to detail, among others, the organization of the 3D digitization into workflows;
- the ongoing Task 3.2 - Integration of 3D H-BIM technologies from INCEPTION (and its first deliverable D3.2) will also build on the current analysis for expanding the part related to BIM and H-BIM models;
- generally, a critical analysis of the currently available software tools and standard file formats will also allow an aware design of future services running on the 4CH Cloud Platform (WP3) and data management (WP5).

Furthermore, some already developed services, in particular those related to the operativity of the 4CH Cloud Platform, shall be maintained and may be kept updated during the project activities. The actual implementation details together with any adjustments, including the related documentation, will be made available at every stage. Moreover, the WP3 services and processes are part of a lively ecosystem of activities that are expected to evolve within the 4CH project to accommodate further needs and requirements that may not be raised so far. The interaction and feedback collected from service providers and user communities will also contribute to improve the project developments as well as to refine the WP3 services offered to the CH community.

9. References

- 1 4CH project - D1.1 - *Initial survey of the experiences and technology state of the art* - <https://doi.org/10.5281/zenodo.6698707>
- 2 4CH internal technical report - *Digitizing Cultural Heritage - A Reconnaissance Investigation on the Data Infrastructure* - <https://doi.org/10.5281/zenodo.6370361>
- 3 Lerones, P. M., Vélez, D. O., Rojo, F. G., Gómez-García-Bermejo, J., & Casanova, E. Z. (2016). "Moisture detection in heritage buildings by 3D laser scanning" *Studies in conservation*, 61(sup1), 46-54. - <https://doi.org/10.1179/2047058415Y.0000000017>
- 4 European Commission DG Connect Unit G2- *Study on quality in 3D digitisation of tangible cultural heritage: mapping parameters, formats, standards, benchmarks, methodologies, and guidelines* - *VIGIE 2020/654*: <https://digital-strategy.ec.europa.eu/en/library/study-quality-3d-digitisation-tangible-cultural-heritage>
- 5 Dell'Unto, N. & Landeschi, G. (2022) *Archaeological 3D GIS*. Taylor & Francis, London-New York
- 6 4CH internal technical report - *Digitization techniques in the field of Cultural Heritage* - <https://doi.org/10.5281/zenodo.6529064>
- 7 European Commission DG Connect Unit G2 - *Report by the Expert Group on Digital Cultural Heritage and Europeana: Basic principles and tips for 3D digitisation of tangible cultural heritage for cultural heritage professionals and institutions and other custodians of cultural heritage* - <https://digital-strategy.ec.europa.eu/en/library/basic-principles-and-tips-3d-digitisation-cultural-heritage>
- 8 Europeana Network Association Members Council - *Task force report "3D content in Europeana"*: <https://pro.europeana.eu/project/3d-content-in-europeana>
- 9 European Commission – *ANNEX to the Commission Implementing Decision on the financing of the Digital Europe Programme and the adoption of the multiannual work programme for 2021 – 2022*: https://ec.europa.eu/newsroom/repository/document/2021-46/C_2021_7914_1_EN_annexe_acte_autonome_cp_part1_v3_x3qnsqH6g4B4JabSGBY9UatCRc8_81099.pdf
- 10 European Commission *Recommendation of 10.11.2021 on a common European data space for cultural heritage*: <https://digital-strategy.ec.europa.eu/en/news/commission-proposes-common-european-data-space-cultural-heritage>
- 11 European Commission - Independent Expert Report on a European Collaborative Cloud for Cultural Heritage - Ex-Ante Impact Assessment: <https://op.europa.eu/en/publication-detail/-/publication/90f1ee85-ca88-11ec-b6f4-01aa75ed71a1/language-en>
- 12 DEEP-Hybrid-DataCloud (DEEP-HDC) EU Funded Project. Available online: <https://deep-hybrid-datacloud.eu/>
- 13 eXtreme-DataCloud (XDC) EU Funded Project. Available online: <http://www.extreme-datacloud.eu/>
- 14 INDIGO-DataCloud EU Funded Project, Web Site. Available online: <https://repo.indigo-datacloud.eu/>
- 15 The Tallinn Declaration: <https://digital-strategy.ec.europa.eu/en/news/ministerial-declaration-egovernment-tallinn-declaration>
- 16 OPEN SOURCE SOFTWARE STRATEGY 2020 – 2023: https://ec.europa.eu/info/sites/default/files/en_ec_open_source_strategy_2020-2023.pdf
- 17 Docker containerization guidelines: https://docs.docker.com/engine/userguide/eng-image/dockerfile_best-practices/
- 18 Fielding, Roy Thomas (2000). "Chapter 5: Representational State Transfer (REST)". *Architectural Styles and the Design of Network-based Software Architectures* (Ph.D.). University of California, Irvine.
- 19 Castelli, L., Felicetti, A. & Proietti, F. (2021) "Heritage Science and Cultural Heritage: standards and tools for establishing cross-domain data interoperability" *Int J Digit Libr* 22, 279–287.
- 20 CRMhs: F. Niccolucci and A. Felicetti (2018) "A CIDOC CRM-based Model for the Documentation of Heritage Sciences," *2018 3rd Digital Heritage International Congress (DigitalHERITAGE (VSMM 2018))*, 2018, pp. 1-6, doi: 10.1109/DigitalHeritage.2018.8810109.
- 21 CIDOC-CRM: "CIDOC CRM. International Committee for Documentation (CIDOC) of the International Council of Museums (ICOM). Version 7.1.1. <http://www.cidoc-crm.org/version/version-7.1.1>
- 22 Humanum: <https://www.huma-num.fr/#>
- 23 Espadon: <http://www.sciences-patrimoine.org/2020/12/selection-espadon/>
- 24 Hescida: <https://hescida.kikirpa.be/>
- 25 DANS: <https://dans.knaw.nl/en/>
- 26 Niccolucci, F., Felicetti, A. & Hermon, S. (2022) "Populating the Data Space for Cultural Heritage with Heritage Digital Twins", *Data, Special Issue on Special Issue "A European Approach to the Establishment of Data Spaces"*, in press.
- 27 Semeraro, C., Lezoche, M., Panetto, H. & Dassisti, M. (2021). "Digital twin paradigm: A systematic literature review". *Computers in Industry*, 130, 103469

- 28 European Commission (2022). *Destination Earth*. <https://digital-strategy.ec.europa.eu/en/policies/destination-earth>
- 29 Costantini, A et al. (2022) “IoTwin: Toward Implementation of Distributed Digital Twins in Industry 4.0 Settings” *Computers* 2022, 11, 67. <https://doi.org/10.3390/computers11050067>
- 30 Bolton, A.; Lorraine, B.; Dabson, I.; Enzer, M.; Evans, M; Fenemore, T.; Harradence, F.; Keaney, E.; Kemp, A.; Luck, A.; Pawsey, N.; Saville, S.; Schooling, J.; Sharp, M.; Smith, T.; Tennison, J.; Whyte, J.; Wilson, A.; Makri, C. (2017). *The Gemini Principles: Guiding values for the national digital twin and information management framework*. <https://doi.org/10.17863/CAM32260>
- 31 *The Great Spa towns of Europ*. <https://whc.unesco.org/en/list/1613/>
- 32 UNESCO. *Intangible Cultural Heritage List*. Available at <https://ich.unesco.org/en/lists>
- 33 <https://www.ics.forth.gr/isl/x3ml-toolkit>
- 34 <https://www.getty.edu/research/tools/vocabularies/aat/>
- 35 <https://perio.do/en/>
- 36 ARIADNEplus - D5.2: *Data Infrastructure update and extension*. <https://zenodo.org/record/4922749#.Yt2zVOzOM0R>

All web references visited on 30/7/2022.

