# SHARING, PRESERVING AND EXPLOITING DIGITAL COLLECTIONS AT THE VATICAN LIBRARY

*Paola Manoni, Mauro Mantovani**
Vatican Library, Vatican City, Italy

## 1. Introduction

For centuries, the Vatican Apostolic Library (BAV) has collected, preserved, studied and made available, manuscripts and archival documents, ancient and modern printed books, coins and medals, prints, drawings, photographs, molds and objects of the highest value, true treasures of science, spirituality, literature, music, and art [1-6]. Its main goal is to enable the progress of research and the development of culture. While it is a highly specialized library whose collections require specialized knowledge and skills in subjects such as linguistics, history, codicology, paleography and others, to correctly deal with the BAV collection, the Vatican Library has from its origins always been an open library - an open library endowed with a humanistic and universal vision, dedicated to promoting Gospel values and dialogue between persons, peoples and generations through culture, but first and foremost through its qualified service to scholars who seek and spread the truth in their very diverse fields of investigation [7].

Over the last twelve years, the Vatican Library has promoted a digitization project [8-9] aiming at the digital acquisition of its entire collection of manuscripts, which is composed, excluding archival units, of 80,000 codices mostly from the Middle Ages and humanistic period. The implementation of an interoperable digital library and the exploitation of its data is a challenge of our times and ideally relies on the same original purpose that expressed the humanistic idealism of the age of Popes Nicholas V (1447-1455), Sixtus IV (1471-1484) and Sixtus V (1585-1590), in the 15th-16th centuries, with the creation of a universal library, open to the public, "for the common benefit of the learned". Today the Vatican Library is still at the service of learned women and men, having evolved during its long history into a modern research library, preserving valuable collections and treasures, which constitute an essential part of the cultural heritage of mankind.

The year 2012 marks the beginning of the digitization project when a preliminary study of the requirements of photographic devices was completed. At the same time, all necessary improvements were made to the Vatican Library's server farm to hold the long-term preservation archive of high-resolution images that were acquired.

The aim of the project was to achieve the twofold objective of long-term digital preservation on the one hand and dissemination of digital contents through a web platform on the other. To achieve both objectives, an action-oriented approach was chosen which focused on the dynamic and innovative application of standards, best practices and shared principles relevant to the domain of international library science.

---

* Corresponding author: mantovani@unisal.it

The article outlines the most important characteristics of this experience and focuses on the evaluation of digital assets that allow digital humanities to study ancient roots of knowledge in a new way. The availability of such a vast collection of digital objects consequently leads to effective data mining operations for the process of extracting, learning and predicting. The library, in its initiatives for the study of the digital contents produced so far, is promoting case studies relating to computer vision for the automatic detection of iconographic contents in selections of digitized manuscripts of the BAV described below.

## 2. Long-Term Data Preservation (LTDP): Pipeline and FITS Format

In the digital age, the problem of long-term preservation is still a challenge. For the BAV, ensuring data legibility and accessibility to digital resources is deemed essential and forms an important part of its preservation strategy. Long-term digital preservation requires the use of technologies and processes which can guarantee the survival of digital information objects while maintaining their integrity and identity [16]. This is the challenge the BAV faces today.

First and foremost, LTDP requires quality control of the digital objects to be ingested within the archiving system. As a first step of the pipeline, all images acquired at the Vatican Library are checked by a validator, in addition to the human post-production control, which is carried out by the staff of the Library's Photographic Lab. The validator adopted for the analysis of TIFF files produced by acquisition devices is the open-source JHOVE[1]. This tool, typically used during routine operations of digital repositories and for digital preservation activities, is able to identify and validate image files and their related technical metadata.

After quality control of the master copies, the LTDP pipeline for the ingestion of digital objects is structured to undertake different tasks which are automatically tracked thanks to a software controller[2] called Inside [11]. Upon completion of validation by JHOVE, Inside enables the automatic conversion of TIFF master copies into the format that the BAV deems suitable for long-term preservation, which is described below.

In considering the suitability of particular digital formats for the purpose of preserving digital information as a relevant resource for future generations, it is useful to analyze some important factors that could be crucial for making the correct choices.

### 2.1. Seven sustainability factors

The Library of Congress in Washington which, in addition to being the U.S. national library, is a center of excellence for the development of metadata for bibliographic resources and for the best practices of digital preservation, designates seven sustainability factors to be considered when choosing the format to entrust digital assets to[3].

These factors are relevant because they influence feasibility and the cost of preserving information content; furthermore, they ensure the effectiveness and durability of digital objects over time, beyond any eventual technological changes.

They are significant whatever strategy is adopted as a basis for future preservation actions. It is appropriate to briefly synthesize these sustainability factors in order to understand the format chosen by the Vatican Library for its digital preservation.

1. **Disclosure:** Disclosure refers to the availability of information and instructions for uses related to long-term format management. Preservation is not feasible

without understanding how the information contained is encoded as bits and bytes in digital files.

2. **Adoption:** Adoption refers to the degree of diffusion of the format. A format that has been employed or used by archival institutions provides evidence of adoption.

3. **Transparency:** Transparency refers to the format's ability to be directly parsed with basic tools, including human readability using the simple mediation of a text-only editor.

4. **Self-documentation:** Formats that are self-documenting are a guarantee of long-term readability. If there are explanatory elements within them the risk of obsolescence is reduced. This is a very important feature that is also applicable to the embedded metadata that describes the digital object, in relation to its content, its technical characteristics and the information that specifies its life cycle (e.g. conversions, updates, etc.).

5. **External dependencies:** External dependencies refer to the degree to which a particular format depends on a particular hardware, operating system, or software to render or use, and the expected complexity in managing those dependencies in future technical environments.

6. **Impact of patents:** Patents and licensing costs related to a digital format can affect, if not prevent, the sustainability of that format in the long run. Although the costs for licenses to decode formats are currently low, the existence of patents can slow down the development of open-source encoders and decoders. Furthermore, prices of commercial software for transcoding content into outdated formats can result in high licensing costs.

7. **Technical protection mechanisms:** The techniques used to protect data and documents in an LTDP archive should allow access to file formats in the present as well as in the future. Therefore, if files are protected in an LTDP archive with mechanisms such as encryption, a security system should be implemented to document their possible decryption.

8. It is these principles which inspired the Vatican Library in the choice of digital format to be adopted for the LTDP archive and the reasons for which the FITS (Flexible Image Transport System)[4] format was selected.

### 2.2. FITS Format

FITS [12] is an image and data format widely used in science, mainly for archiving astronomical images, and has been in use for more than 50 years; thus, we may argue that its longevity is undisputed because it has been in use since the 1970s. FITS was conceived as a format to be used for the transfer of data between astronomical observatories and research centers (hence the "T" in the acronym). It was later used as a tool for archiving and analyzing scientific data by institutions such as NASA (National Aeronautics and Space Administration), ESA (European Space Agency), ESO (European Southern Observatory). Its implementation started in 1976 and three years later the format was officially presented in the "International Image Processing Workshop" in Trieste (Italy)[5]. Since then, the format has evolved and the first official version was standardized in 1993 by the International Astronomical Union FITS Working Group, while the last release of the format dates back to 2016. Its development has always been guided by the principle "once FITS, always FITS": no change in the specifications should ever conflict with retrospective FITS files so as to avoid past files becoming potentially unreadable by future reading software.

Outside of these research areas, the FITS format is still not that popular, and is mostly unknown to the general public. In the context of cultural heritage, BAV is one of the early adopters of FITS: interest regarding this format is emerging, following the use case of the Library as demonstrated by the publication of UNI 11845:2022 standard which is described below.

In the scientific world there are huge archives of FITS images: as an example, we can mention the archive of the PAN-STARRS system containing 1.6 Petabytes of data, the ESO archive containing 44.8 million files for a total of 1.01 Petabytes of data[6] and finally the FITS archive of the Vatican Apostolic Library, currently containing more than 8 million[7] FITS images for a total of about 500 Terabytes of data, to be increased due to the ongoing digitization project.

FITS was adopted because it has a series of distinctive characteristics that rely on the sustainability factors and make it especially suitable for the purpose of long-term preservation, particularly for its simple structure: it is made up of two fundamental parts, an initial part with ASCII encoding, called Header Data Unit (HDU), and a binary part containing the data (Figure 1).



```
SIMPLE   =                      T / Java FITS: Tue Feb 20 12:36:46 CET 2024
BITPIX   =                      8 / bits per data value
NAXIS    =                      3 / number of axes
NAXIS1   =                   7760 / size of the n'th axis
NAXIS2   =                  10328 / size of the n'th axis
NAXIS3   =                      3 / size of the n'th axis
UNIKEY   =                      T / Compliant with UNI 11845:2022
EXTEND   =                      T
HDUNAME  = 'Capp.Giulia.VIII.39_0303_fa_0148r' / Original image filename
GRPID1   =                     -1 / The image is part of a group
GRPLC1   = '        '              / File group URI
LONGSTRN= 'OGIP 1.0'              / The OGIP long string convention may be used
CREATOR = 'FITSBatchConv v2.0.0' / Software that created this FITS file
INSTRUME= 'Phase One IQ180'
PROGRAM = 'Capture One 7 Macintosh'
DATE-OBS= '2013-01-21T13:10:21'
DATE     = '2024-02-20T13:49:20' / Date and time of FITS file creation
EXPTIME = 0.008000023756566464 / Exposure Time
REFERENC= 'https://digi.vatlib.it/view/MSS_Capp.Giulia.VIII.39' / Publication UR
ORIGIN  = 'Biblioteca Apostolica Vaticana'
OBJECT  = 'Capp.Giulia.VIII.39'
COLORMAP= 'RGB     '              / Colors mapping
CTYPE1   = '        '              / Linear transformation on axis 1
CTYPE2   = '        '              / Linear transformation on axis 2
CTYPE3   = 'RGB     '              / Name of the coordinate represented by axis 3
CRPIX1   =                    0.0 / Location of reference point along axis 1
CRPIX2   =                    0.0 / Location of reference point along axis 2
CRPIX3   =                    0.0 / Location of reference point along axis 3
CRVAL1   =                    0.0 / Coordinate CTYPE1 at the reference point CRPIX1
CRVAL2   =                    0.0 / Coordinate CTYPE2 at the reference point CRPIX2
CRVAL3   =                    0.0 / Coordinate CTYPE3 at the reference point CRPIX3
CUNIT1   = 'mm      '              / Units of CRVAL1 and CDELT1
CUNIT2   = 'mm      '              / Units of CRVAL2 and CDELT2
CUNIT3   = '        '              / Units of CRVAL3 and CDELT3
CDELT1   = 0.054978354978354974 / Coordinate increment at reference point
CDELT2   = 0.054978354978354974 / Coordinate increment at reference point
CDELT3   =                    1.0 / Coordinate increment at reference point
CHECKSUM= 'ZN1EcNjCZNjCbNjC'      / checksum for the current HDU
DATASUM = '3564769436'            / checksum of the data records
IMGURESL= 'INCH    '              / Resolution Unit
IMGXRESL=                  462.0 / Horizontal resolution
IMGYRESL=                  462.0 / Vertical resolution
BAV01    = 'Alamire '              / Contributor
BAV02    = '        '              / Lens model
BAV03    =      11.000140161779888 / Aperture Value
BAV04    =                     50 / ISO speed
```

*Figure 1. Example of FITS Header.*

In relation to sustainability, it is easy to demonstrate its peculiar characteristics of self-documentation, transparency, freedom of use and autonomy.

- **Self-documentation:** The first blocks of information composing the FITS file, which are encoded in ASCII-and therefore easy to access, contain a set of metadata, essential for using the data contained in the file itself. In particular, the first bytes contain the number of bits per pixel, the number of axes and the number of elements per axis in the image contained in the file. This information is sufficient for understanding and accessing the content of the individual HDUs contained in the FITS files.
- **Transparency:** The ASCII encoding allows immediate reading access to all metadata contained in the FITS file header; in addition, the structure of the FITS file makes it possible to point directly to the bytes corresponding to the various pixels of the image by exploiting the information contained in the first bytes of the file. Consequently, the metadata are easily accessible by means of data analysis tools and can be read instantly even through any available hex viewer.
- **Freedom of use:** The FITS format is an open format, so its technical specifications are publicly available and there is no legal restriction on the use of the format. Otherwise, it would not be possible to guarantee the readability of files in the long term.
- **Autonomy**: No special software or hardware is required to use the format. In order to gain access to a file's content, it is sufficient to access the file's binary data.

Therefore, the evolution, or conversely the obsolescence, of the hardware and software, will never jeopardize the use of a FITS file archive.

### 2.3. Requirements for LTDP archives and UNI 11845:2022

The implementation of the LTDP archive of FITS files revolves around a series of requirements that have been identified for managing the BAV's digital objects. These requirements are closely related to the possibilities that the FITS format offers. They pertain to:

- **Homogeneity**: Homogeneity provides the LTDP archive with efficient management of data validation processes and actions for the recovery of corrupt or truncated data, as well as helping in information retrieval processes. In addition, it facilitates the use of the archive contents even if the documentation describing the archive is itself lost.
- **Validation**: The validation process includes a series of procedures designed to verify the consistency and structure of the data archive and to pinpoint the eventual issue of data corruption or data loss.
- **Conversion**: The long-term versatility and usability of a digital image archive is guaranteed by the ability to convert the FITS files into other exchange formats if requested by end users (e.g.: PDF format, JPEG format, TIFF format). The conversion process is therefore crucial in building a long-term preservation archive.
- **Information retrieval**: The extraction of metadata from the files contained in the archive is an essential process for the maintenance and usability of the archive itself. The LTDP archive allows an efficient metadata extraction by

accessing only those portions of the FITS file that contain the metadata itself.

- **Origin and history**: Whatever the nature of the digitized object, it is essential to keep knowledge of the file acquisition and creation methods and all modifications undergone by the file itself for the purpose of long-term preservation. This information is an essential part of the image archive, and it is invoked or directly stored in the FITS files.

- **Semantic content**: Whatever the nature of the digitized object, a basic description has to be included either as referenced or wrapped data, in relation to the semantic content that can therefore enable the archived data to be interpreted in the long term.

- **Structural models**: The LTDP archive "FITS format based" can preserve the logical relations between the digital data and the related metadata, as well as between the digital data itself when there is a need to keep memory of the configuration pertaining to physical and/or logical arrangement of the data represented and managed in one or more FITS files. This includes, for example, cases such as:
    - Preservation metadata (history of the data) and descriptive metadata (semantic content): typically implemented as XML files in specific namespaces, these metadata can be embedded as binary files in specific HDUs of the FITS file.
    - Embedding of multiple objects into a single FITS file.
    - Relationships between multiple FITS files.

These requirements and the related technical documentation have been standardized and published by the body for Italian National Unification: the official Italian Standardization Body entrusted with the development, publication and promotion of national and international standards[8].

The creation of UNI 11845:2022 [13] within the UNI Technical Commission "Documentation and information"[9] occurred according to the principles and rules of international standards. The technical committee that carried out this work was attended by experts, in addition to the BAV, from the Italian Ministry of Culture (Directorates of Library and Archival Heritage), academic institutions and scientific bodies (University of Rome Tor Vergata, National Institute for Astrophysics, European Space Agency), Italian national associations of archivists and librarians, the archive of the Bank of Italy and Italian companies specialized in LTDP and the preservation of digital assets.

LTDP requirements stated in this standard are paving a new way in the technical evaluation of LTDP archives, especially for cultural assets and are, moreover, demonstrating a change of paradigm: from the conception of a long-term archiving that provided for the ingestion of data considered as inert and closed in the preservation system, to an archive structure where access to the data implies constant data curation (Figure 2).

This means that access is high-performance, enabling data analysis for the preservation of documents and data that can be updated and documented with the semantic content of the original objects.

## 3. The digital library and its discovery services

The second objective of the Vatican Library's digitization project concerns implementation of the digital library platform [14].

*Figure 2. Example of conversion: from FITS format to a derivative format by using the BAV application.*

The conceptual design of the functional architecture is mostly focused on enhancing the discovery services for the digital collection. The implementation of a digital library platform is not a mere display of digitized images nor a way of linking bibliographic records in a catalog. Firstly, a dissemination pipeline must solve the problem of the image aggregation of complex digital objects, such as books. A folder stored in a file system physically contains all the files of a distinct semantic unit (e.g., images of a manuscript, per each folio) however its unitary structure is obviously inadequate for a viewer to represent the digital object in its aggregation. Folders contain files whose ordering is simply by file name and not according to the physical or logical sequences of contents.So even just viewing one image after another, to digitally compose the volume requires the application of standardized procedures. The physical sequence of a manuscript, for example, may involve different numbering or irregularities that should be documented and managed for viewers. Physical as well as logical sequences (about the structure and contents of the digital object) are typically managed by metadata syntax. Without going into technical details of the standards that manage the aggregation of data in complex digital objects, behind the scenes of digital library functions there has been a METS profile[10], designed by the BAV, since the first generation of the online digital library[11] and is now used for the DVL (DigiVatLib)[12] (Figure 3). The DVL, therefore, offers navigation of the digitized contents per each of its distinct collections and each of these refers to descriptive syntaxes, which are peculiar to each type of resource. For example, searching for a digitized manuscript entails indexes corresponding to the specific elements describing codicological information along with the individual works within a codex, but a table of contents is also included as a pointer to browse the digitized resource. The DVL therefore provides dual access to information, which includes both the consultation of images and the consultation of all the information associated with a digitized resource, in its native description, including any bibliographic references.
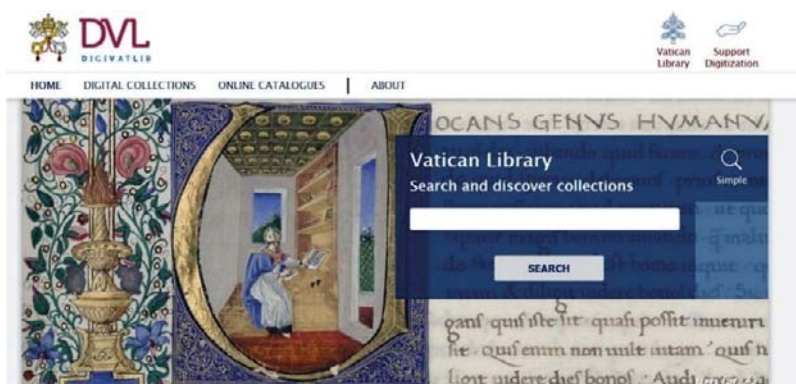
*Figure 3. DVL – DigiVatLib home page.*

Another widely used service is the online availability of inventories and finding aids. These references constitute an individual collection with its specific search menu and their permalinks are included in the descriptive metadata of manuscripts, in order to connect the description to the related references in inventories. The ingestion into the DVL of each digital object implies its automatic connection to the bibliographic records available in the online catalogs. The recognition of the structure of the persistent identifier, associated to the digital object, allows the DVL to direct the objects in each collection: manuscripts, incunabula, graphics, coins and medals, printed books[13].

This is one of the main tasks performed by the back-end system, called DWF (DigiWorkFlow), which is responsible for orchestrating the complete workflow, i.e. from the ingestion of digital images to the structuring of the METS, along with the association of descriptive and structural metadata and the viewer. The latter conforms to the specifications of the interoperability protocol IIIF (International Image Interoperability Framework)[14].

The DVL was implemented in 2016 and at that time, its most innovative feature was the adoption of the IIIF [15], which radically changed the paradigm of use of digital objects. The Vatican Library's scientific interest in the IIIF data model[15] dates back to 2015. At that time, many digitization projects that libraries were implementing, provided access to valuable collections but considered them as individual repositories, like *silos.* Access to these rich collections held in cultural heritage institutions was crucially dependent on each database configuration and service functions. A means to break down these barriers was to employ the techniques of the semantic web able to interconnect repositories of digital objects (Figure 4). In essence, the IIIF data model is based on the semantic web approach to allow scholars to compare and study digital objects by using tools able to offer a simultaneous display of resources despite their different provenances.

The term "interoperability" in the IIIF acronym stands for the technological effort that seeks to meet this need, to express the opportunity of exchanging information or services between IT systems, by facilitating their mutual interaction. The interoperability of digital libraries began in 2012 at Stanford University Libraries. A group of experts from the Digital Library Systems and Services Division studied how to access digital images of manuscripts on the web, along with all the data and documentation pertaining to them (the description of a catalog, detailed notes commenting on particular elements, transcription of folios, etc).

*Figure 4. DWF – DigiWorkFlow: example of the Archival Information Package (general information).*

This study led to the creation of a standardized "interoperability protocol" for the free circulation of digital data and images in the web, directly traceable through search engines and independent from the repositories of individual libraries.

In a nutshell, by using an IIIF-compliant viewer, it is possible to display a digitized manuscript on a laptop, simply by calling up its web address (what is technically called URI: Uniform Resource Identifier of the digital object) and then placing it alongside another manuscript (or even several others) for all relevant comparisons.The application of the interoperability protocol to the world of manuscripts immediately appeared of great interest for studies of philology, bibliology, paleography, with particular regard to critical editions and the possibility of carrying out virtual reconstructions of collections dispersed in various libraries or simply of fragmentary materials preserved in different places.
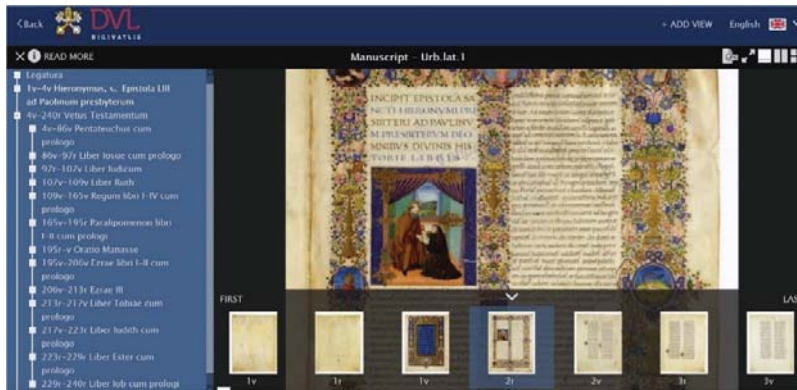


*Figure 5. DVL: example of a digital resource.*

The Vatican Apostolic Library (Figure 5) is one of the most prominent early adopters of IIIF technology [10]. In 2015, the experimentation of a small but significant first case study by the BAV was completed which involved applying the protocol on a small number of digitized manuscripts and managing the annotations: transcriptions, notes, and marginalia on folios. After this first pilot project, only a year later, the DVL was

implemented, fully compatible with the IIIF. In the years 2016-2019, the Vatican Apostolic Library promoted a special research activity on some specific groups of manuscripts, the results of which were available in the said "IIIF-mode" by comparing and annotating on the folios of the manuscripts transcriptions, comments on the texts, glosses and illuminations, identification of copyists, illuminators, and owners. The Andrew W. Mellon Foundation, the New York foundation that supported the development of the IIIF in its most significant applications, sponsored this three-year research. The initiative made use of the technical consultancy of the same Stanford experts who first studied this revolutionary conception of cultural heritage, which became, in its digital expression, a sort of unified heritage, a "mare magnum" on the open horizon of knowledge. An innovative service in which the Vatican Library can express its mission, confirming its centuries-old tradition. The outcome of this research was the implementation of a platform capable of managing selections of resources with the main purpose of producing annotations (transcriptions, comments, comparative analysis of texts and images) performed by curators. The peculiar characteristics of this implementation have demonstrated how IIIF capabilities can be used to study manuscript collections. Furthermore, the features required in this project were consistent with the design patterns and development practices embraced by IIIF communities to be integrated in open-source tools.



*Figure 6. Example of an IIIF manifest of a digitized manuscript.*

## 4. Thematic Pathways on the Web

The BAV use case on IIIF annotations is entitled *Thematic Pathways on the Web*[16] (which produced over 26,000 annotations[17]). As previously mentioned, this is an implementation carried out in collaboration with Stanford University Libraries for manuscripts selected under a specific theme. The project platform has implemented software such as Mirador[18] and Spotlight[19]. The latest is an open-source app based on the Ruby on Rails framework,[20] for managing platforms with search functions that exhibit collections and objects taken from digital libraries.

The content of all annotations (created in Mirador) is indexed along with the

metadata (exposed as URI-addressable resources through the OPAC (Online Public Access Catalog) and read by Spotlight), thus constituting a semantically enriched system that allows scholars to carry out integrated searches for all the available information related to a resource, as well as managing the narrative part of the thematic pathway. Each of these is a sort of virtual exhibit through texts, images and through the discovery of content in the innovative ways offered by the IIIF. Essentially, a thematic pathway in the BAV Spotlight platform is composed of three different types of information:

- A general description (introduction, historical information, etc.) of the chosen theme that represents the narrative in which the thread of the discourse is intertwined with the contents enjoyed via IIIF.
- A set of descriptive and content metadata for each manuscript chosen.
- A large pool of annotations, comments, insights on detailed parts of a manuscript (e.g. texts, commentaries, miniatures, etc.) and transcriptions of information units.

Above all, the objective is to deal with themes identified in the digital collections and showcase the selected digital resource, so as to conceive an exhibition itinerary. These digital exhibitions promote a new perspective to the study of manuscripts by means of web communication and IIIF (Figure 7). The first four exhibits carried out in the three-year research period focused on the following topics:

A.   Course in Paleography (Greek and Latin, from Antiquity to the Renaissance)

The rich collection of manuscripts preserved in the Library makes it possible to follow the evolution of the Greek and Latin scripts all the way from antiquity to the Renaissance. A careful selection of images of the manuscripts, accompanied by transcriptions and comments (such as IIIF annotations) is offered as teaching material for a course in Greek or Latin. The availability of online manuscript images, together with the possibilities offered by the IIIF API (Application Programming Interface), enables a complete transformation of teaching practice in this field.
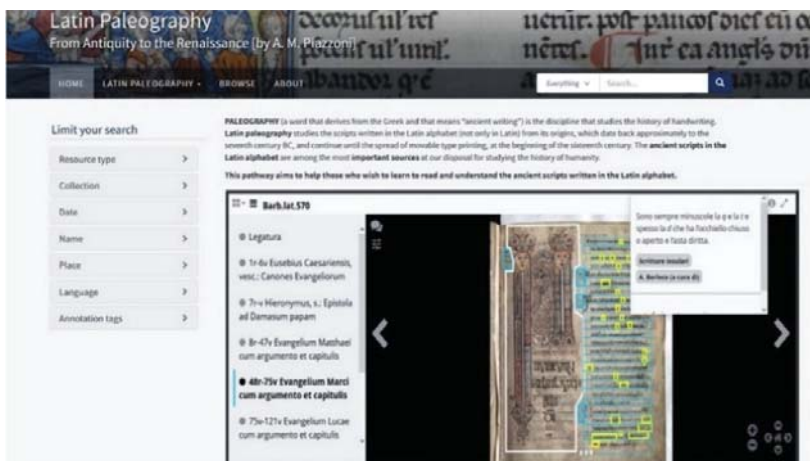


*Figure 7. Thematic Pathways on the Web: example of page.*

B.   Latin Classics

The Vatican Library owns one of the most important collections of manuscripts with texts by Classical Latin authors, many of them richly illustrated. The aim of this pathway is to describe 81 manuscripts directly from the original codices: metadata and annotations pertaining to the study of texts and illuminations have been provided. The work throws light not only on the illustrations of the texts but especially on the relationship between text, illuminations, comments and glosses regarding the main text in a manuscript (marginal or interlinear notations).

C.   Vatican Palimpsests

The Vatican Library has identified more than 380 manuscripts in its own collections, which include palimpsests, that is, erased then recycled parchment folios. This pathway intends to present this rich and scarcely explored material to the public by making in-depth archaeological research on the palimpsests of twenty-four select manuscripts and recover their lost identities with the help of IIIF technology.

D.   The Library of a "Humanist Prince"

The library of Federico da Montefeltro (1422-1482), Duke of Urbino (from 1474), is known as a typical humanist collection. In the first years, Federico bought or ordered manuscripts in Florence (both in writing and in illumination), but later preferred Ferrara or Paduan artists and scribes active in Urbino. This pathway points out the characteristics of the two schools, very different in style, and the most important artists (half of the chosen manuscripts is representative of the Florentine school while the other half of the Ferrara and Paduan schools). The *Thematic Pathways on the Web* platform can be used for hosting many digital contents. The most recent virtual exhibitions in the web showcase are: *Super Hanc Petram* (of the Vatican Medal Collection) and more recently *Traveling with Dante*[21].

## 5.  Beyond AI, again IIIF

The adoption of artificial intelligence, computing vision and machine learning is rapidly growing. The impact of these new IT frontiers in the world of digital humanities means there is an increasing number of experiments and services: from the improvement of optical character recognition of texts to new services for data analysis of digital collections.

The Vatican Library is promoting experiments both in the field of character recognition[22] and in research related to the machine learning potential of a convolutional neural network [17]. A convolutional neural network (ConvNet/CNN) [18] is a deep learning algorithm that is particularly useful for acquiring input images and analyzing them against an associated dataset, in order to recognize known *patterns* by identifying objects, classes and categories (while in the training phase it builds the *reference patterns*). The artificial neural network experiment focused on the automatic detection of iconographic content of a selection of digitized manuscripts. The most advanced content of this case study is strongly related to interoperability, giving it an innovative feature since the machine learning carried out is proposed as IIIF annotations. This result is not only remarkable from a mere technical point of view but above all for the opportunity that it offers to integrate current discovery tools with data mining aimed at detecting the contents of digital objects.

### 5.1. YOLO and its biases: a case study

The convolutional neural network that was implemented made use of YOLO (You Only Look Once)[23]: an algorithm capable of detecting objects present in images. YOLO "looks once" at the image because it can predict its meaning after just one interpretation. The detection is carried out in accordance with the instances of objects belonging to a specific class recognized within an image. The algorithm basically recognizes the existence of objects in an image using bounding boxes and assigns types or classes to the identified objects. For example, it takes an image as input and generates one or more bounding boxes, each labeled as a reference class. In principle, a convolutional neural network can handle classification, multi-class localization and detection of objects with multiple occurrences. Image classification algorithms predict the type or class of an object in an image (input) from a predefined set of classes for which the algorithm has been trained. The output is a class or label that represents a particular object, accompanied by a confidence index of that prediction[24].

Object localization algorithms detect the presence of an object in the image and represent its position in a bounding box. The Vatican Library implemented the YOLOV5 version (using PyTorch[25]). YoloV5 is one of the most relevant pre-trained networks that uses COCO (Common Objets in COntext)[26], i.e. a large dataset of images: over one and a half million instances of objects segmented within images, corresponding to 80 categories of objects. The artificial neural network was able to identify and frame the iconographic contents with the limit of the identification of COCO classes.

In the application of COCO in the domain of illuminations, this case study highlighted the issue related to the needs of a dataset consistent with the corpus of data to be analyzed: the current datasets available in pre-trained networks do not have classes corresponding to the categorization of the objects identified in illuminations or iconographic details within manuscripts. With regard to the choice of manuscripts (the images to be segmented): instead of establishing selection criteria, for the purpose of this case study, it was preferred to opt for a random selection by considering only two aspects:

- A quantity of images adequate for the training of the artificial neural network.
- The availability of descriptive metadata of the manuscripts to be processed.

The data sample consists of 1,874,999 images, corresponding to 5,186 manuscripts. The successful segmentation process led to the detection of 54,629 objects. For each, the network added a specific probability index.

The manuscripts were divided into the following groups, in order to proceed with the analysis of the results or the relevance of the classification of the iconographic content identified in the segmented images:

- Manuscripts (115 in total) with at least 100 records,
- Manuscripts (232 in total) with at least 50 entries,
- Manuscripts (483 in total) with at least 20 records,
- Manuscripts (747 in total) with at least 10 records,
- Manuscripts (1054 in total) with at least 5 records.

The sample analyzed, based on the presence of illuminations, is 75 manuscripts belonging to the first group[27], in which the artificial neural network identified and classified 21,933 objects.

The interpretation of the network, depending on the COCO dataset, obviously

classified iconographic contents corresponding to the classes to which it belongs. It is important to dwell on the attribution of classes corresponding to shapes, in analogy to the patterns of the images of which COCO is composed.

As for an example of biases, we may refer to the selection of circular or ovoid details found in the digitized folios, such as halos or stamps, but also holes in parchment subsumed by the artificial neural network under inconsistent categories with respect to their real meaning (Figure 8).

The systematic human validation of the results obtained from the automated neural network was limited to some specific categories. From a quantitative point of view, at the end of these checks, the data obtained were the following:

- o 2575 overall total results of which:
  - 1260 results were actually consistent with the class assigned by the network (approximately 49% of the total);
  - 1315 results with class reassignment or with class added from scratch and, if applicable, related box from scratch (approximately 51% of the total);
  - 1034 results (irrelevant and later deleted);
  - 309 new classes added.



Figure 8. Example of segments interpreted in YOLO using COCO as 'Frisbee'.

From a qualitative point of view, the errors (1315+1034) that required human intervention are obviously interesting. In the group of irrelevant and, therefore, deleted results, these were mostly cases of outlines or shapes interpreted as iconographic content. Classes that were reassigned or added, are either due to the failure of the network, a lack of the corresponding class in COCO or a mismatch of concepts (Figure 9).

### 5.2. From results of the artificial neural network to automated IIIF annotations

The results of the artificial network relating to the 75 selected manuscripts (with at least 100 objects detected, as previously reported) - which for each segment (object) denotes the class and the probability index (*confidence*) of this prediction - were

transformed into IIIF annotations and indexed within the Spotlight platform in which the aforementioned BAV's *Thematic Pathways on the Web* are disseminated.
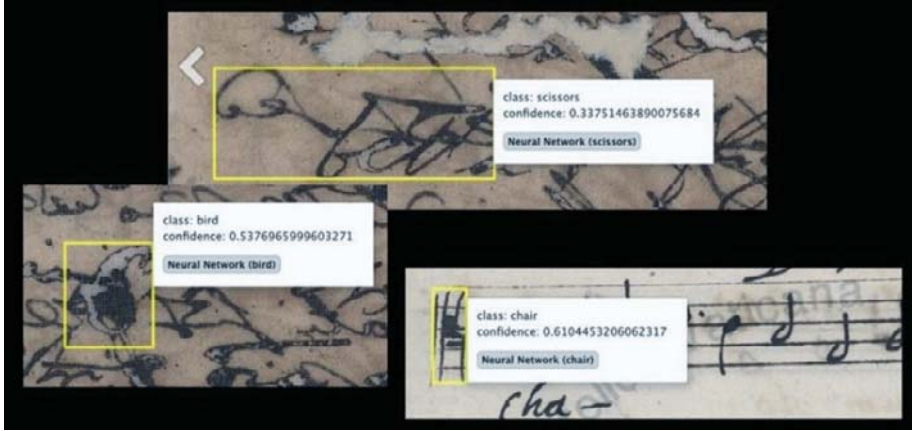


*Figure 9. Example of biases.*

While in the latter the annotations for each thematic pathway are part of a curatorial project that sees the commentary, transcriptions, classifications as manually filled annotations, in this case study, the script (called AI4I[28]) converted the output of the artificial neural network to an annotation that is integrated into the set of all data collected, using IIIF, on that specific digital object.

The assigned class is transformed into the IIIF annotation tag (which is preceded by a source-specific prefix: 'AI4I (neural network)' while the body of the annotation reports the result in its original form.

The coordinates of the segmented images become the IIIF annotation region that the script transforms into canonical annotations[29] (Figure 10).

```
tag="person"
testo=" class: person<br> confidence: 0.5682786703109741"
canvas="https://digi.vatlib.it/iiif/MSS_Urb.lat.1/canvas/p0003"
x=400
y=500
w=300
h=500
coordinates="xywh="+str(x)+","+str(y)+","+str(w)+","+str(h)
coordinateshtml="M"+str(x)+","+str(y)+"h"+str(w/2)+"v0h"+str(w/2)+"v"+str(h/2)+
"v"+str(h/2)+"h-"+str(w/2)+"h-"+str(w/2)+"v-"+str(h/2)+"z"
manifest="https://digi.vatlib.it/iiif/MSS_Urb.lat.1/manifest.json"
id="Urblat1ann14"
color="#ff8c00"
strokewidth="3"
```

*Figure 10. Example of transformation: data from YoloV5 and the relation to the IIIF 'json manifest' of the manuscript to which it belongs.*

The *AI4I script* transforms these values into the information shown in Figure 11.

```
{
  "@context":"http://iiif.io/api/presentation/2/context.json",
  "@type":"oa:Annotation",
  "motivation":[
    "oa:tagging",
    "oa:commenting"
  ],
  "resource":[
    {
      "@type":"oa:Tag",
      "chars":"'+tag+'"
    },
    {
      "@type":"dctypes:Text",
      "format":"text/html",
      "chars":"'+testo+'"
    }
  ],
  "on":[
    {
      "@type":"oa:SpecificResource",
      "full":"'+canvas+'",
      "selector":{
```

*Figure 11. Example of an IIIF annotation.*

The *AI4I script* also includes other elements that are taken from the descriptive metadata of the manuscripts, such as context-specific information: shelfmark, dates, provenance, language of the manuscript.

The opportunity of integrating an environment of indexed data with the results produced by an algorithm could lead to improving a corpus of images and classes consistent with the domain of the original resources: to date, a corpus focused on illuminations is not yet available.

The success of this experiment can be found in the ability to use artificial intelligence within an integrated IIIF system, where automatic processing is only a tool to facilitate the retrieval of data and scholars' research within a hybrid environment, in which the data produced by human intelligence determine the use of artificial intelligence data (and not vice versa). In itself, the recognition of an illuminated content by a 'machine' does not add anything if a human being is already able to recognize it.

However, the value of these data can be found in their association with other elements such as descriptive metadata in order to enable data mining analysis to identify, with the help of algorithms, the evolution of styles, or schools or to recognize different manuscripts/illuminations.

## 6. Conclusion

The aim of this paper was to illustrate the process of production, conservation, development and increasingly wide sharing of the collections and digital resources of the

Vatican Apostolic Library, an institution of ancient origin that also today "is an outstanding means for the Church to contribute to the development and dissemination of culture, in support of the work of the Apostolic See. Through its various sections, it is responsible for collecting and preserving a vast patrimony of learning and art and of making it available" [19].

The evolution of the Vatican Library's digital services in the span of a decade fosters knowledge of the collections made possible by a complex technological world that exploits data analysis. If, previously, access to the collections was only through an OPAC, the integration of bibliographic information with digital objects opens up a varied scenario where the interoperable vision of the IIIF for the sharing of research and resources breaks in. With growing awareness of the value of digital assets, the long-term preservation of heritage data and documents is considered as indispensable for the usability of the digital library.

The world of artificial intelligence has so far seen in the BAV experimental cases aimed at achieving objectives related to an integrated use of algorithms, in the coexistence of metadata, content production (such as IIIF annotations) and human-generated research tools. In particular, the experimentation described in its workflow, demonstrated a method for setting up an 'AI-to-IIIF' pipeline capable of ensuring: an efficient image partitioning method (the construction of a segmentation algorithm independent of YOLOV5), the ability to implement a CNN, the transformation of the results into *canvas* within an 'annotation list' in the IIIF manifests of the selected manuscripts, the incorporation of annotations in the Spotlight platform. The annotations obtained, led to two important results. The first, a precious 49% of success of the artificial neural network without any human intervention, and the second, an equally relevant 51% of new attributions: corrections, reassignments of categories, new identifications that constitute an important element for a new training of the network for the detection of segmented images, using more consistent classes based on the aforementioned results. As is known, a heuristic method of machine learning initially proceeds through trial and error: what is defined as 'supervised learning'. In this type of learning, the artificial intelligence is trained to provide a known answer: the values are entered and compared with the correct answer. In this way the algorithms can gain experience in order to avoid previous errors and adapt their response for subsequent values, so that results are gradually improved until they give an accurate answer. The results of this case study are available at: https://ai4mss-poc.vatlib.it. On this web site (which is conceived as one of the thematic pathways of the exhibits included in the BAV Spotlight) scholars can evaluate all the annotations, tags, metadata of the manuscripts, analyses of the results, as well as groupings by century and provenance of the identified classes.

The current challenges and opportunities that the development of advanced technologies and cybernetics pose to our culture and society today, starting from humanistic institutions such as libraries, make us responsible, for the "cultivation of the human spirit", and the commitment, especially towards the new generations, of promoting and "protecting wisdom, in other words, knowledge that is human and humanizing" [20-22].

**Notes**

[1] The original JHOVE project was a collaboration between JSTOR and Harvard University Library with funding from the Andrew W. Mellon Foundation for the Electronic-Archiving Initiative. JHOVE is currently maintained by the Open Preservation Foundation; Cf. https://github.com/openpreserve/jhove [Accessed: Jan. 28, 2024].

[2] Inside, on the one hand, pertains to check long-term preservation procedures and,

on the other, tracks the process of dissemination of digital collections. It was designed and implemented by the Coordination of IT Services of the BAV.

[3] Cf. Library of Congress:
https://www.loc.gov/preservation/digital/formats/sustain/sustain.shtml    [Accessed: Jan. 28, 2024].

[4] "FITS Standard Web site: "FITS Standard Document: The official reference document that defines the requirements for FITS format data files". Available at: https://fits.gsfc.nasa.gov/fits_standard.html [Accessed: Jan. 28, 2024].

[5] Cf. https://archive.stsci.edu/fits/users_guide/node8.html [Accessed: Jan. 28, 2024].

[6] Cf. https://www.eso.org/public/science/archive/ [Accessed: Jan. 28, 2024].

[7] Total amount of FITS files at the time of writing of this article (January 2024).

[8] The UNI is recognized by the Italian State and by the European Union and is one of the national bodies of the International Standards Organization (ISO) and European Committee for Standardization (CEN).

[9] Cf. https://www.uni.com/normazione/organi-tecnici/dettaglio-commissione/?id=37 [Accessed: Jan. 28 ,2024].

[10] Metadata Encoding and Transmission Standard. It is the well-known administrative metadata schema, developed by the Library of Congress and internationally adopted in best practices for digital libraries.

[11] The first implementation, thanks to the support of Heidelberg University, was the DWork: an application able to manage the process flow of digitization and web presentation of the digitized works.

[12] Cf. https://digi.vatlib.it [Accessed: Jan. 28, 2024]. The DVL (DigiVatLib – Digital Vatican Library), as well as its back-end application, the DWF (DigiWorkFlow), were designed by the Coordination of IT Services of the BAV and engineered by NTT Data Company (Japan).

[13] At the time of writing of this article (January 2024), the collection of printed books, currently composed only of cropped images related to details of pages in rare books, is evolving to manage a complete flow of digitization of printed books.

[14] International Image Interoperability Framework. Cf.:
https://iiif.io [Accessed: Jan. 28, 2024].

[15] Cf. https://iiif.io/api/model/shared-canvas/1.0/ [Accessed: Jan. 28, 2024].

[16] Cf. https://spotlight.vatlib.it [Accessed: Jan. 28, 2024].

[17] It is the largest experimentation with annotations produced so far.

[18] Cf. https://projectmirador.org/ [Accessed: Jan. 28, 2024]. Mirador is an open-source image (and recently video) viewer with the ability to zoom, view, compare and annotate digital assets compatible with the IIIF protocol.

[19] Cf. https://github.com/projectblacklight/spotlight [Accessed: Jan. 28, 2024]

[20] Cf. https://rubyonrails.org/ [Accessed: Jan. 28, 2024]. Ruby on Rails is an open-source framework for web applications written in the Ruby programming language.

[21] Exhibit on the occasion of the 7th centenary of the death of Dante Alighieri (1265-1321) curated by the Vatican Library, promoted by the Dante Scientific-Organizational Committee (Pontifical Council for Culture) and supported by the Cortile dei Gentili.

[22] An optical recognition project of Greek writing, for the achievement of a dataset of one million characters. The project, led by the Japanese company Toppan and based on images of BAV manuscripts, is currently underway.

[23] This is the name of an open-source object detection algorithm that processes images very quickly, in virtually real time. Cf. https://github.com/ultralytics/yolov5 [Accessed: Jan. 28, 2024].

[24] The probability index is between 0 and 1: the closer it gets to 1, the more certain

the prediction is considered by the network.

25 PyTorch is an open-source software. It is a machine learning framework based on the Python programming language and the Torch library used to create artificial neural networks.

26 Cf. https://cocodataset.org [Accessed: Jan. 28, 2024].

27 Manuscripts were considered according to the illuminations detected.

28 AI4I: Artificial Intelligence for Illuminations.

29 This includes the connotative properties of IIIF annotation (with the values, in the motivation, related to 'commenting' and 'tagging'. Cf. Simplest Annotations https://iiif.io/api/cookbook/recipe/0266-full-canvas-annotation/ [Accessed: Jan. 28, 2024].

### References

[1]    Cardinal Stickler, A. M. (1989) *The Vatican Library. Its History and Treasures*, New York, Belser.

[2]    Piazzoni, A. M., Jatta, B. (2010). *Conoscere la Biblioteca Vaticana. Una storia aperta al futuro*, Vatican City, Biblioteca Apostolica Vaticana.

[3]    Ceresa, M. (2012) La Biblioteca Vaticana tra Riforma Cattolica, crescita delle collezioni e nuovo edificio (1535-1590), in: Storia della Biblioteca Apostolica Vaticana vol. II, Vatican City, Biblioteca Apostolica Vaticana.

[4]    Montuschi, C. (2014) La Vaticana nel Seicento (1590-1700). Una biblioteca di biblioteche, in: Storia della Biblioteca Apostolica Vaticana vol. III, Vatican City, Biblioteca Apostolica Vaticana.

[5]    Jatta, B. (2016) La Biblioteca Vaticana e le arti nel secolo dei Lumi (1700-1797), in: Storia della Biblioteca Apostolica Vaticana vol. IV, Vatican City, Biblioteca Apostolica Vaticana.

[6]    Rita, A. (2020) La Biblioteca Vaticana dall'occupazione francese all'ultimo Papa Re (1797-1878), in: Storia della Biblioteca Apostolica Vaticana vol. V, Vatican City, Biblioteca Apostolica Vaticana.

[7]    Pope Francis (2022) Praedicate evangelium. Apostolic Constitution on the Roman Curia and Its Service to the Church in the World, art. 243. Vatican City.

[8]    Ammenti, L. (2020) *Per litteras ad astra. Storia dell'automazione della Biblioteca Apostolica Vaticana dalla carta al digitale*, Rome, Aracne.

[9]    Piazzoni, A. M. (2018) *The Process for the Digitization of Manuscripts in the Vatican Library.* Vatican City, Biblioteca Apostolica Vaticana.

[10]   Manoni, P., Núñez Gaitán, Á., Schuler, I. (2018) *The Vatican Apostolic Library's Digital Preservation Project*, available at: https://library.ifla.org/id/eprint/2113/1/160-manoni-en.pdf [Accessed: Jan. 28, 2024].

[11]   LISCI, Daniele (2024). *The software controller Inside*, in: Piazzoni, Ambrogio M. (ed.) *The Process for the Digitization of Manuscripts in the Vatican Library.* Vatican City, Biblioteca Apostolica Vaticana.

[12]   Giuffrida, G. (2024) *The FITS format: analysis and use* in: Piazzoni, Ambrogio M. (ed.) *The Process for the Digitization of Manuscripts in the Vatican Library*. Vatican City, Biblioteca Apostolica Vaticana.

[13]   UNI/TC014/SC01 (2022). Norma tecnica *UNI 11845:2022: Processi di gestione della conservazione a lungo termine di immagini digitali con l'uso del formato FITS*, available at: https://store.uni.com/uni-11845-2022 [Accessed: Jan. 28, 2024].

[14]   Manoni, P. (2024) The DigiVatLib in: The Process for the Digitization of

Manuscripts in the Vatican Library, Vatican City, Biblioteca Apostolica Vaticana.

[15]     Crame, T. (2017) *An Introduction to IIIF*, available at: https://re-sources.digirati.com/iiif/an-introduction-to-iiif/ [Accessed: Jan. 28, 2024].

[16]     Manoni, P. (2020) L'adozione del IIIF nell'ecosistema digitale della Biblioteca Apostolica Vaticana*, DigItalia*, 2, pp. 96-105. DOI: 10.36181/digitalia-00017 [Accessed: Jan. 28, 2024].

[17]     Manoni, P. (2023) AI4MSS: un esperimento di intelligenza artificiale alla Biblioteca Apostolica Vaticana, in: Guardando oltre i confini: partire dalla tradizione per costruire il futuro delle biblioteche, Rome, AIB, 2023 - DOI: 10.1400/294065 [Accessed: Jan. 28, 2024].

[18]     Venkatesan, R. (2020) *Convolutional Neural Networks in Visual Computing: A Concise Guide*, Cleveland: CRC Press.

[19]     Pope Francis (2022). *Praedicate evangelium*. *Apostolic Constitution on the Roman Curia and Its Service to the Church in the World*, art. 243. Vatican City.

[20]     Pope Francis (2019) *Christus vivit. Post-Synodal Apostolic Exhortation to young people and the entire people of God*, n. 223. Vatican City.

[21]     Pope Francis (2024) *Artificial Intelligence and Peace. Message for the 57th World Day of Peace*, Vatican City, available at: https://www.vatican.va/content/francesco/en/messages/peace/documents/20231208-messaggio-57giornatamondiale-pace2024.html [accessed Jan. 28, 2024].

[22]     Pope Francis (2024). *Artificial Intelligence and the Wisdom of the Heart: Towards a Fully Human Communication. Message for the 58th World Day of Social Communications,* 12th January 2024, Vatican City, available at: https://press.vatican.va/content/salastampa/en/bollettino/pubblico/2024/01/24/240124b.html [accessed Jan. 28, 2024].

**Biographical Notes**

**Paola Manoni** is the Head of the Coordination of IT Services at the Vatican Library. Her research and work areas are focused on metadata, digital libraries, interoperability, machine learning and long-term preservation for cultural heritage. She is currently chairing the UNI Technical Committee for "Documentation and Information" as well as the Sub Committee for "Technical interoperability" in the same organization. She represents the Vatican Library in many international initiatives such as the IIIF Consortium and, in this context, she is also co-chairing the "IIIF Design" Working Group.

**Mauro Mantovani** is the current Prefect of the Vatican Library. A Salesian Catholic priest, he is Professor of Theoretical Philosophy at the Salesian Pontifical University, where he was Dean of the Faculties of Philosophy and of Social Communication Sciences, and Rector Magnificus (2015-2021). His research focuses on borderline issues between philosophy, theology and science. He is a Councillor of the Pontifical Academy of St. Thomas and Member of many other Pontifical Commissions and Scientific and Academic Committees.

**Summary**

This article outlines the two main purposes of the Vatican Library digitization project

(long-term digital preservation and dissemination of contents), focusing on the exploitation of digital assets and providing the humanities with a new way to study the ancient roots of knowledge. The article describes the LTDP (Long Term Data Preservation) pipeline, and the choice of preservation format based on sustainability factors and focuses on the LTDP archives requirements in light of the UNI 11845:2022 standard. The dissemination pipeline describes the digital library platform and demonstrates the Vatican Library's adoption of the IIIF protocol (International Image Interoperability Framework). Finally, the article describes a pilot project relating to the use of AI for the recognition of iconographic details in illuminations, with the peculiarity of managing the results of the artificial neural network within IIIF annotations.

## Riassunto

Questo articolo delinea i due principali scopi del progetto di digitalizzazione della Biblioteca Vaticana (conservazione digitale a lungo termine e diffusione dei contenuti), concentrandosi sulla valorizzazione del patrimonio culturale digitale. Il processo di digitalizzazione consente alle discipline umanistiche una nuova via per lo studio delle antiche radici della conoscenza. L'articolo descrive quindi la pipeline LTDP e le scelte del formato di conservazione effettuate in base a fattori di sostenibilità. Vengono inoltre discussi i requisiti degli archivi LTDP alla luce della norma UNI 11845:2022. La pipeline di diffusione descrive la piattaforma della biblioteca digitale e dimostra l'adozione da parte della Biblioteca Vaticana del protocollo IIIF (International Image Interoperability Framework). Infine, l'articolo descrive un progetto pilota relativo all'utilizzo dell'AI per il riconoscimento di dettagli iconografici nelle illuminazioni, con la peculiarità di gestire i risultati della rete neurale artificiale all'interno delle annotazioni IIIF.