

Towards a Comprehensive Assessment of the Quality and Richness of the Europeana Metadata of Food-related Images

Yalemisew Abgaz, Amelie Dorn, Gerda Koch, Jose Luis Preza Diaz

Adapt Centre Dublin City University, Austrian Academy of Sciences, European-Local Austria
Dublin Ireland, Vienna Austria, Graz Austria
Yalemisew.abgaz@adaptcentre.ie, kochg@europeana-local.at
{Amelie.Dorn,JoseLuis.PrezaDiaz}@oeaw.ac.at,

Abstract

Semantic enrichment of historical images to build interactive AI systems for the Digital Humanities domain has recently gained significant attention. However, before implementing any semantic enrichment tool for building AI systems, it is also crucial to analyse the quality and richness of the existing datasets and understand the areas where semantic enrichment is most required. Here, we propose an approach to conducting a preliminary analysis of selected historical images from the Europeana platform using existing linked data quality assessment tools. The analysis targets food images by collecting metadata provided from curators such as Galleries, Libraries, Archives and Museums (GLAMs) and cultural aggregators such as Europeana. We identified metrics to evaluate the quality of the metadata associated with food-related images which are harvested from the Europeana platform. In this paper, we present the food-image dataset, the associated metadata and our proposed method for the assessment. The results of our assessment will be used to guide the current effort to semantically enrich the images and build high-quality metadata using Computer Vision.

Cultural image analysis, Semantic enrichment, Computer Vision, Ontology, Knowledge design

1. Introduction

As a result of open access policy (European Commission, 2011) adopted by Galleries, Libraries, Archives and Museums (GLAMs), a huge collection of cultural and historical resources is now available on the internet to promote access. Many GLAMs started to publish digital resources and the associated metadata to support ease of search and retrieval by both humans and computer agents (Abgaz et al., 2018; Stork et al., 2019). However, a significant portion of the resources still lacks quality and rich metadata. In some cases, the available metadata only describes basic bibliographic information such as title and the publication year of the resource. Often, the interpretation and utilisation of the data by users other than subject experts are hampered by the lack of domain knowledge and machine-readable rich semantics to understand the dataset.

By rich semantics, we mean that the availability of multiple descriptors of a resource including bibliographic information, domain-specific annotation, links to interconnected resources, etc. By quality, we refer to a multitude of metrics including the correctness, reuse of existing terms, use of multiple languages, etc. defined in (Zaveri et al., 2015; Debattista et al., 2016; Debattista et al., 2018). The availability of rich semantics enables the exploitation of the metadata in several creative ways by both humans and machines, whereas ensuring the quality enables to build dependable systems which produce high-quality results.

There have been efforts made to provide joint platforms and standard tools to aggregate and publish data from GLAMs. Europeana.eu¹ is one of such platforms established by the European Union as a virtual aggregator of digitised collections from more than 3,500 institutions across Europe. This platform brings together contributing institutions and Europeana local platforms to aggregate content, facilitate knowledge transfer and innovation, distribute cultural her-

itage content and engage users to participate in the use and contribution of the resources via a centralised platform supporting multilingual and multi-faceted search and retrieval of the available resources (Haslhofer and Isaac, 2011; Isaac and Haslhofer, 2013). Cultural and historical collections including images, pictures, paintings, photographs, specimen, etc. are at the primary focus of Europeana. The Europeana effort started in 2008 and the current collection still suffers from a lack of rich metadata for many of its objects. As these metadata emerge from many different contributors, there are still many discrepancies (both in coverage and semantics) in the richness and quality of the metadata despite the effort made to standardise using the European Data Model (EDM)² (Haslhofer and Isaac, 2011).

In this position paper, we present our proposed approach for analysing the quality and semantic richness of selected images related to food by taking the Europeana collection as a cases study. Even if there is a consensus on the importance of analysing the quality and richness of the whole Europeana collection, in this paper, we will focus on analysing the coverage and the quality of the semantic annotation of food-related images using food-related domain-specific ontologies and thesauri. By historical images, we refer to the collection of images, pictures, paintings and photographs that represent some historical or cultural importance. It is observed that even if the collection is enriched with metadata of some kind, the historical, cultural and domain-specific aspects of the data are underrepresented by the available metadata. The metadata is not semantically enriched to reflect the detailed content of the images. This problem is partially demonstrated during the evaluation of the quality of search results obtained from the platform when users search the collection using historical and cultural aspects. It further requires a meticulous investiga-

¹<https://www.europeana.eu/portal/en>

²<https://pro.europeana.eu/resources/standardization-tools/edm-documentation>

tion to identify the strength and weaknesses of the metadata in representing the detailed aspects of the cultural images. In this research paper, we present our research questions followed by our proposed approach. The questions are:

- How much semantic annotation is available for food-related images and what is the quality of the available metadata?
- How rich is the domain-specific annotation in using multiple vocabularies?
- What aspects (technical, social, cultural, political, etc.) of the images are semantically well annotated?
- What are the gaps that are observed in the metadata and how can we address it using semantic enrichment?

Our focus is on historical images related to food. So far we collected 65 buckets of food images representing particular food topics. These images are used to analyse the semantic richness and quality of the metadata in depth. The Europeana images contain associated metadata which can be downloaded through a special functionality provided to us by the Europeana Local-Austria. We have collected all the metadata (semantic annotations) of the images in JSON and RDF formats. We will use this metadata throughout to evaluate the quality and richness of the metadata.

This paper is organised in five sections. Section 2. introduces Europeana and the coverage of the collection followed by some discussion of relevant research in Section 3. Section 4. presents the data collection process, the target food image collection and the metadata. Section 5. presents the proposed approach and metrics to be used and, finally, we present the conclusion and future work in Section 6.

2. Background

Europeana is an aggregator platform which provides central access to resources from GLAMS. The platform allows users to search all the collections that are distributed across several institutions in Europe from a single search interface. However, Europeana does not host the original digital objects on its servers but provides metadata about the items and dereferenceable links to the institutions that hold the collections. This approach allows Europeana to maintain the level of aggregation required to support search and retrieval of information, and it enables the institution to keep and continuously improve the collection and the associated metadata while the original data stays in the content providers' websites. Europeana uses metadata from the providers and maps this metadata using EDM (Isaac, 2013; Innocenti, 2014) to provide a single common interface for efficient and searchable information.

Currently, Europeana offers access to about 60 million items including books, music, artworks and more³. But Europeana's aim is not only to aggregate the metadata but also to involve content providers in the very challenging task of improving the quality of the metadata by achieving good quality metadata for 70% of their collections. This is achieved through the use of enrichment tools to improve

the existing metadata and by assisting content providers to follow new quality frameworks.

3. Related Research

Previous research has been conducted to determine the quality of the Europeana metadata. Peter et al. Király et al. (2019; Kirly and Bchler (2018) evaluated the data quality in Europeana focusing on its multilinguality. The authors defined metrics for evaluating the multilinguality using metrics such as completeness, consistency, conformity and accessibility. Even if this paper provides good coverage of the metrics used in determining multilinguality, its focus is only on a language-related quality measure. In our proposed method, we would like to widen the scope and include other quality metrics available elsewhere and also measure the diversity of the metadata concerning the coverage of the subject matter presented in the image collection. Other metrics proposed in (Gavrilis et al., 2015) present a quality measure in metadata repositories. The authors proposed five metrics together with some contextual parameters concerning metadata generation and use. The quality measures the authors use include completeness, accuracy, appropriateness, consistency and auditability. These metrics also overlap with the metrics used to evaluate the multilinguality of the metadata. However, they incorporate contextual parameters such as a requirement for higher accuracy using weightings of the metrics. They evaluated their metrics using Europeana data of the archaeology aggregator CARARE (Connecting ARchaeology and ARchitecture for Europeana).

Generic and comprehensive data quality measures are also proposed by (Debattista et al., 2016) incorporating 24 metrics distributed across four major categories. These metrics also measure the quality of metadata and present the results using percentages. The approach also provides a customisable implementation of the metrics which can be used based on the specific requirements of the evaluation. We will initially consider all the metrics that are covered in the paper and later filter those that are not applicable. A followup paper (Debattista et al., 2018) has also used a Europeana dataset to demonstrate the applicability of the proposed metrics, however, the detail of the analysis reported in the paper is not sufficient to make any concrete decision regarding the quality of the metadata. Thus, it is important to use the proposed metrics to drill down and investigate the selected quality issues.

4. Data Collection

For this study, we use digital images and their associated metadata collected from the Europeana online image collection. In the whole repository, there are more than 58 million digital objects (images, texts, audios, videos and 3D objects) available with the associated metadata describing mainly the bibliographic information of the objects.

4.1. Food Images

For this study, we restrict our focus on digital images including paintings, photos, drawings and sketches. Since conducting a deep analysis on the full collection is beyond the scope of our project, we narrow down the focus only

³<https://www.europeana.eu/portal/en/about.html>

Search Topic	Items	Topic	Items	Search Topic	Items
Alimentation sistemas culinarios	838	Food and Nederland s	122	Lebensmittel+	1773
Breakfast	100	Food and Norway	280	Lunch	363
Cafe	123	Food and party	29	Painting and Food	182
Comedor	36	Food and people	465	Painting and Fruit	484
Dessert	532	Food and Portugal	60	Panaderia	307
Drawings and Illustrations	98	Food and shop	1968	Photograph and Breakfast	45
Eating	880	Food and shopping	152	Photograph and Dinner	28
Food and autumn	11	Food and society	366	Photograph and Eating	20
Food and Belgium	127	Food and Spain	48	Photograph and Food	63
Food and celebration	64	Food and sports	12	Photograph and Fruit	207
Food and cuisine	27	Food and spring	55	Photograph and Lunch	41
Food and culture	9445	Food and summer	25	Print and Food	5
Food and customs	16	Food and Sweden	758	Produccion y alimentos	4060
Food and dancing	27	Food and Switzerland	52	Reposteria	394
Food and Denmark	32	Food and traditions	47	Soup	300
Food and family	216	Food and winter	37	Speisesaal	244
Food and Finland	100	Food and woman	64	Still life	354
Food and France	227	Food and work	130	Still life and Food	8765
Food and Germany	180	Food+Austria	4986	Lebensmittel	627
Food and Luxembourg	33	Food+machine	396	Godigital	6
Food and man	280	Frhstck	38	Gastronomy	1100
Food and market	119				
Total					42969

Table 1: The distribution of the images across different buckets

to food-related images. This is due to the following reasons. First, food is associated with our daily life and it is one of the most familiar topics for humans to deal with. Second, food represents the culture and the history of both traditional and modern society. We also have food-related images that cover a long period from the early centuries to the modern-day. Third, food is highly interconnected to several other disciplines including health, fitness, nutrition, economics, business, culture, society, agriculture, technology, politics, etc. This allows us to analyse the richness of the metadata associated with food and to evaluate the coverage of these aspects of food in the available metadata. Finally, since this analysis is being conducted in the context of the ChIA⁴ project (accessing and analysing cultural images with new technologies), the focus is on testing the quality of the existing semantic enrichment of cultural food images to improve access and enhance analysis using artificial intelligence applications such as chatbots to support interactive search. Results from this project will not only enable wider access possibilities for Europeana images but also provide increased semantic capabilities for Digital Humanity researchers to work with image-related data.

So far we have collected images from the Europeana platform including photos, paintings, drawings using 64 non-exclusive buckets. These images are collected by using several food-related keywords prepared by experts from sociolinguistic, computer science, and digital humanity domains. A total of 42,969 images are collected and included in the analysis. Table 1 summarises the distribution of food images across the search topics.

4.2. Metadata

We use a platform provided by the Europeana Local-Austria team to download both the images and the metadata. For all the selected images, the metadata is available in a JSON and RDF format which is provided in the EDM standard. Depending on the provider, additional metadata is also available for most of the images. This indicates that



Figure 1: Sample image with its metadata.

there is some uniformity in the usage of bibliographic data across all the images, however, the use of additional metadata fields and ontologies largely depends on the provider of the image. A sample image is shown in Figure 1 and a snippet of the associated metadata is given in the text below.

```
{
  "object": {
    "about": "/2059513/data_foodanddrink_efd_LGMA_0933",
    "aggregations": [
      {
        "about": "/aggregation/provider/2059513/
          data_foodanddrink_efd_LGMA_0933",
        "edmDataProvider": {
          "def": [
            "Local Government Management Agency"
          ],
          "edmIsShownBy": "http://griffiths.askaboutireland.ie/gv4/dev/
            fandd_images/selection_of_breads_and_butter.jpg",
          "edmObject": "http://griffiths.askaboutireland.ie/gv4/dev/
            fandd_images_thumbs/selection_of_breads_and_butter.jpg",
          "edmProvider": {
            "def": [
              "Europeana Food and Drink"
            ]
          },
          "edmRights": {
            "def": [ "http://creativecommons.org/licenses/by-sa/3.0/" ],
            ...
          },
          "concepts": [
            {
              "about": "http://data.europeana.eu/concept/base/48",
              "prefLabel": {
                "de": ["Bild (Fotografie)"],
                "fi": ["Valokuva"],
                "ko": [" "],
                ...
              }
            }
          ]
        }
      }
    ]
  }
}
```

Since all the metadata related to an image is downloaded into a single file, the number of metadata files in the collection is equal to the size of the images. The metadata in RDF

⁴<https://chia.acdh.oew.ac.at/>

can be directly used by the selected quality assessment tool. This metadata will be analysed for its quality using specific metrics and following a sampling approach, we also consider a manual evaluation of the descriptive nature of the associated metadata compared to the actual image. Even if this task consumes a significant amount of time, it is worth to check the quality going a little beyond what the automated analysis tools provide. This metadata is further used to analyze the richness of the metadata in describing the concepts/aspects depicted in the image. This looks into potential ontologies, vocabularies and thesauri in the food domain and checks how many of them are used across the images to semantically annotate the images.

5. Proposed Assessment Approach

We considered two types of quality measures applicable to the assessment of the quality and richness of the metadata. The first is using quantitative measures where objective metrics are used to analyze quality based on some mathematical formula, and the second one is a qualitative approach where an expert judgement is required to determine the quality. In this work, we will use both methods in such a way that existing widely used objective metrics are selected and used to evaluate the quality and the richness of all the metadata of the selected images. The qualitative evaluation focuses on a deep analysis of the metadata by comparing it with the corresponding image and evaluate how much of the explicit and implicit information contained in the target image is represented by the metadata. In this particular approach, we will use experts from the food domain to qualitatively evaluate the selected images and the corresponding metadata to evaluate both the quality and the richness of the metadata. This approach complements the quantitative approach with expert judgement on the accuracy and correctness of the metadata and identifies the gap between the potentially useful information contained in the image and what is represented in the metadata.

5.1. Quality Analysis Tools and Metrics

Several researchers have identified and proposed metadata quality metrics including the 67 metrics and 18 quality dimensions (Zaveri et al., 2015) and 27 metrics implemented (Debattista et al., 2016). The later metrics are also implemented in a linked data quality assessment framework (Luzzu). Due to its comprehensive and deployable tool, we conduct an initial experiment with the Luzzu framework to quantitatively analyze the quality of our dataset. The metrics included in the Luzzu framework are categorised into four major categories (Debattista et al., 2016): representational, where the focus is on the design of the data in terms of common best practices and guidelines; contextual category, which focuses on the relevance, correctness, understandability and timeliness; intrinsic category, which focuses on correctness and coherence of the data including syntactic validity, semantic accuracy, consistency, conciseness and completeness; and accessibility category, which focuses on the (re)usability of linked data resources by both machines and humans. All these categories contain relevant metrics for our dataset. However, not all the metrics are directly useful for the work we are conducting, such as the

length of the characters in a URI. Thus, we carefully select the metrics we use to assess the quality of the metadata

5.2. Semantic Richness Analysis

Zavier et al (Zaveri et al. (2015) further identified metrics that are used to determine the richness of the metadata: detection of good quality interlinks, the existence of links to external data providers and dereferenced back-links. However, in (Debattista et al., 2016) interlinking is included in the accessibility metrics. In analysing the richness of the metadata, even if these metrics measure how richly the metadata is connected with other sources, our main interest is to check whether these external links are connected to domain-specific ontologies, vocabularies, thesauri which give detailed context and meaning to the contents of the images. This requires a further analysis of the external links included in the metadata and evaluating whether these links point to domain-specific or bibliographic metadata.

To achieve this objective, we identify major domain-specific ontologies (Dooley et al., 2018), vocabularies⁵ (Harpring, 2018; Caracciolo et al., 2013; Leatherdale et al., 1982) and thesauri in the areas of the topics of the selected datasets. Mainly, we narrowed down our focus to food-related metadata to evaluate the semantic richness in providing useful information for supporting educators, scientists and even content providers to focus more on the semantic enrichment using domain-specific metadata which makes the collection more relevant to the users.

6. Conclusion

In this paper, we present the current work we are conducting to evaluate the quality and the richness of the metadata of a selected set of food image collections from Europeana to identify gaps of the current semantic enrichment. To this end, we selected 42,969 images and the associated metadata for the evaluation. We proposed both qualitative and quantitative evaluation methods with existing scientifically proven methods and metrics. So far, we have identified most of the relevant metrics, selected the framework and acquired the relevant data. Our next step will be to apply the method and evaluate the quality and richness of the dataset using the proposed methods. One of the challenging tasks is the qualitative evaluation of the richness and the contextual accuracy of the metadata compared to the contents of the images. To address this issue, we will incorporate evaluators from the three categories of Europeana users: the educators, scientists, and content providers to evaluate the richness and the correctness of the metadata.

Acknowledgements:

This research is funded by the Austrian Academy of Sciences under the funding scheme: go!digital Next Generation (GDNG 2018-051). The ChIA project is carried out in collaboration with the ADAPT SFI Research Centre at Dublin City University. The ADAPT SFI Centre for Digital Media Technology is funded by Science Foundation Ireland through the SFI Research Centres Programme and is co-funded under the European Regional Development Fund (ERDF) through Grant # 13/RC/2106.

⁵<http://www.getty.edu/research/tools/vocabularies/aat/help.html>

7. Bibliographical References

- Abgaz, Y., Dorn, A., Piringer, B., Wandl-Vogt, E., and Way, A. (2018). Semantic modelling and publishing of traditional data collection questionnaires and answers. *Information*, 9(12).
- Caracciolo, C., Stellato, A., Morshed, A., Johannsen, G., Rajbhandari, S., Jaques, Y., and Keizer, J. (2013). The agrovoc linked dataset. *Semant. Web*, 4(3):341348, July.
- Debattista, J., Auer, S., and Lange, C. (2016). Luzzua methodology and framework for linked data quality assessment. *J. Data and Information Quality*, 8(1), October.
- Debattista, J., Lange, C., Auer, S., and Cortis, D. (2018). Evaluating the quality of the lod cloud: An empirical investigation. *Semantic Web*, 9(1):131–150.
- Dooley, D. M., Griffiths, E. J., Gosal, G. S., Buttigieg, P. L., Hoehndorf, R., Lange, M. C., Schriml, L. M., Brinkman, F. S. L., and Hsiao, W. W. L. (2018). Foodon: a harmonized food ontology to increase global food traceability, quality control and data integration. *npj Science of Food*, 2(23).
- European Commission. (2011). Commission recommendation of 27 october 2011 on the digitisation and online accessibility of cultural material and digital preservation. Technical report, European Commission.
- Gavrilis, D., Makri, D.-N., Papachristopoulos, L., Angelis, S., Kravvaritis, K., Papatheodorou, C., and Constantopoulos, P. (2015). Measuring quality in metadata repositories. In Sarantos Kapidakis, et al., editors, *Research and Advanced Technology for Digital Libraries*, pages 56–67, Cham. Springer International Publishing.
- Harpring, P. (2018). Getty vocabularies: Linked open data version 3.4. semantic representation.
- Haslhofer, B. and Isaac, A. (2011). data.europeana.eu: The europeana linked open data pilot. *International Conference on Dublin Core and Metadata Applications*, 0:94–104.
- Innocenti, P. (2014). *Migrating Heritage: Experiences of Cultural Networks and Cultural Dialogue in Europe*. Ashgate.
- Isaac, A. and Haslhofer, B. (2013). Europeana linked open data - data.europeana.eu. *Semantic Web*, 4:291–297, 01.
- Isaac, A. (2013). Europeana data model primer. Technical report, European Commission.
- Király, P., Stiller, J., Charles, V., Bailer, W., and Freire, N. (2019). Evaluating data quality in europeana: Metrics for multilinguality. In Emmanouel Garoufallou, et al., editors, *Metadata and Semantic Research*, pages 199–211, Cham. Springer International Publishing.
- Kirly, P. and Bchler, M. (2018). Measuring completeness as metadata quality metric in europeana. In *2018 IEEE International Conference on Big Data (Big Data)*, pages 2711–2720.
- Leatherdale, D., Tidbury, G. E., Mack, R., Food, of the United Nations., A. O., and of the European Communities., C. (1982). *AGROVOC : a multilingual thesaurus of agricultural terminology / Donald Leatherdale ; with the collaboration of G. Eric Tidbury and Roy Mack*. Api-
mondia, by arrangement with the Commission of the European Communities [S.I.], english version. edition.
- Stork, L., Weber, A., Miracle, E. G., Verbeek, F., Plaat, A., [van den Herik], J., and Wolstencroft, K. (2019). Semantic annotation of natural history collections. *Journal of Web Semantics*, 59:100462.
- Zaveri, A., Rula, A., Maurino, A., Pietrobon, R., Lehmann, J., and Auer, S. (2015). Quality assessment for Linked Data: A survey. *Semantic Web Journal*.