

# On assessing weaker logical status claims in Wikidata Cultural Heritage records

Alessio Di Pasquale<sup>a</sup>, Valentina Pasqual<sup>b,\*</sup>, Francesca Tomasi<sup>b</sup> and Fabio Vitali<sup>a</sup>

<sup>a</sup> *Department of Computer Science, University of Bologna, Italy*

*E-mails: alessio.dipasquale@studio.unibo.it, fabioitali@unibo.it*

<sup>b</sup> *Digital Humanities Advanced Research Centre (DH.arc), Department of Italian Studies and Classical Philology, University of Bologna, Italy*

*E-mails: valentina.pasqual2@unibo.it, francescatomasi@unibo.it*

**Abstract.** This work analyses the usage of different approaches adopted in Wikidata to represent information with weaker logical status (WLS, e.g., uncertain information, competing hypotheses, temporally evolving information). The study examines four main approaches: non-asserted statements, ranked statements, non-existing valued objects, and statements qualified with properties `P5102:nature of statement`, `P1480:sourcing circumstances`, and `P2241:reason for deprecated rank`. We analyse their prevalence, success, and clarity in Wikidata. The analysis is performed over Cultural Heritage artefacts stored in Wikidata, divided into three subsets (i.e., visual heritage, textual heritage, and audio-visual heritage), and compared with astronomical data (stars and galaxies entities). Our findings indicate that (1) the presence of weaker logical status information is limited, with only a small proportion of items reporting such information, (2) the usage of WLS claims varies significantly between the two datasets in terms of prevalence and success of such approaches, and (3) precise assessment of WLS statements is made complicated by the ambiguities and overlappings between WLS and non-WLS claims allowed by the chosen representations. Finally, we list a few proposals to simplify and standardise this information representation in Wikidata, hoping to increase its clarity, accuracy and richness.

**Keywords:** Wikidata, ranked statements, weaker logical status, uncertainty, Cultural Heritage

## 1. Introduction

Since 2012, Wikidata [1] has been one of the most outstanding platforms for collecting and sharing Linked Open Data through the web.

Wikidata encompasses a multitude of facts, including some that may be contrasting since they come from different and disagreeing sources. Expecting global consensus on the “true” data would be unrealistic since many facts are disputed or uncertain. Wikidata allows conflicting data to coexist and provides mechanisms to organize this plurality, going beyond the triple-based representation of factual information, for instance, including contextual metadata and constraints over those statements [2, 3]. For instance, Wikidata contributors can add time-sensitive information through qualifiers and ranks to represent temporally evolving information (e.g., the number of followers of a YouTube channel that is updated year after year) or multiple coexisting (and possibly competing) claims over the same subject (e.g., maintaining both the old as well as a new theory over some topic). In many such cases, multiple

---

\*Corresponding author. E-mail: valentina.pasqual2@unibo.it.

1 information items are present. Yet, newer or better information does not replace older or less true assertions. How- 1  
2 ever, they coexist next to each other, and one or more mechanisms are used to signal their simultaneous presence 2  
3 and, when appropriate, the currently adopted stance. 3

4 We understand these statements as enjoying a somehow *weaker logical status* than asserted statements: they are 4  
5 neither true nor false, but they are, e.g., true from a specific moment onward but not earlier, or true up to a given 5  
6 moment but not afterwards, or accepted as true by most people but not everybody, etc. 6

7 It is a cultural necessity in many (if not all) fields of knowledge to access available data about a complex topic 7  
8 entirely and objectively as they evolve, as different scholars or models interpret them and represent available hy- 8  
9 potheses rather than a positive certainty. For instance, Cultural Heritage scholars study attributions, the temporal 9  
10 context of events, the temporal evolution of content, and the contradictions of opinions and assertions so that ex- 10  
11 pressing weak statements, i.e., claims we are not sure about, becomes a necessary tool to increase precise awareness 11  
12 of the currently available data for those who consult or reuse it. Interpretation thus plays a central role in humanities 12  
13 disciplines. Yet, Cultural Heritage knowledge graphs and domain ontologies frequently limit the formalisation of 13  
14 these phenomena or only partially represent them ([4, 5], cf. Section 2). Recently, a rekindled interest has been 14  
15 shown in the formalisation of uncertain statements [6–8], claiming that interpretation constitutes a focal point in 15  
16 humanities data and metadata. Interestingly, these works prove how different motivations for nuanced statements 16  
17 with varying degrees of truth create a small and consistent number of approaches to express them. We conclude that 17  
18 studying the very idea of weak logical status claims *per se*, independently of their different justifications, can help 18  
19 shed light on these commonalities and their relative merits and issues. WLS claims are used not only for missing 19  
20 or incomplete information but also for the correct representation of personal opinions or beliefs, for temporally 20  
21 constrained information, for geographically constrained information, etc. 21

22 Wikidata supports several patterns to represent situations best expressed with weaker logical status claims. In this 22  
23 paper, we analyse some of these patterns as they are employed in actual collections, both in the humanities and, as a 23  
24 comparison, in hard sciences. A factor that increases complexity is that many of these uses have partially overlapping 24  
25 semantics, i.e., Wikidata contributors can use them for other purposes beyond weaker logical status claims, and this 25  
26 muddles the correct identification and interpretation of the situations we are interested in. We, therefore, want to 26  
27 discuss both the designed use of each approach, its actual usage and success in Wikidata. Finally, we discuss the 27  
28 impact of their ambiguous applications due to the coexistence of multiple uses for the same techniques. 28

29 In particular, we analysed four main families of approaches to the weaker logical status of statements, *asserted* 29  
30 *vs. non-asserted statements*, *ranked statements*, *unknown objects*, and *qualified statements*. In this paper, we try to 30  
31 answer the following research questions: 31

- 32 – RQ1 - How widespread are these approaches in the current state of Wikidata? 32
- 33 – RQ2 - How does the cultural domain of the Wikidata topics (and, presumably, of the individuals contributing 33  
34 to the data regarding the Wikidata topics) affect and reflect on the relative success of some WLS types over 34  
35 others? 35
- 36 – RQ3 - Does the actual usage of the surveyed approaches match their designed use declared by Wikidata? 36  
37 37

38 In addition, we wonder about how we could improve the clarity and cleanliness of such differentiation. 38

39 To perform such analysis, we accessed and downloaded two large sets of topics from Wikidata, one belonging to 39  
40 the Cultural Heritage (visual works of art such as paintings and statues, text documents, and audio-visual entities) 40  
41 and another from astronomy (celestial bodies such as stars and galaxies). Both use multiple fuzzy assertions and 41  
42 hypotheses and, therefore, need assertions with weaker status (e.g., attributions uncertainties or physical locations 42  
43 moving over time for paintings vs. spectral class or radial velocity for stars). 43

44 The decision to use a comparative dataset in this study is motivated by the wish to explore the similarities and 44  
45 differences between astronomical and humanities academic practices. Both fields involve studying unique objects, 45  
46 such as stars or books. Yet, the way data is treated differs, with astronomical observations becoming scientific data as 46  
47 soon as they are used as evidence of phenomena [9], while humanities rarely can go beyond learned interpretations. 47

48 The data sources also vary, with humanities researchers using historical documents, literature, art, and oral tra- 48  
49 ditions, each having varying levels of reliability and introducing systemic and insurmountable uncertainty. In as- 49  
50 tronomy, uncertainty is often related to instrumental limitations and observational conditions. Methodologically, 50  
51 astronomy relies on empirical observation, mathematical modelling, and experimental validation. At the same time, 51

1 humanities research is frequently interpretative and qualitative, and the necessary proof to obtain historical certainty is often unattainable [10]. This difference leads to distinct epistemological foundations, with the humanities acknowledging the subjectivity and cultural bias in interpretations [11], and astronomy seeking to minimise uncertainty through rigorous data collection and adherence to physical principles [12].

2 Our study's hypotheses and assumptions include the idea that annotators in Cultural Heritage and astronomy may approach data incompleteness and uncertainty differently, with Cultural Heritage favouring qualitative, context-rich representations of competing hypotheses and astronomy leaning towards more quantitative, data-centric representations. This difference may reflect broader epistemological stances in their respective communities. Additionally, our study assumes that these distinct approaches to handling data incompleteness and uncertainty may impact the ease of integrating data from these fields in interdisciplinary research, with Cultural Heritage data potentially requiring more effort for reconciliation due to its contextual and subjective nature.

3 Overall our findings show that the amount of weaker logical status statements in Wikidata seems suspiciously low, as only 0,4% of visual artworks report attribution debates, a fairly low figure compared to, e.g., a more reasonable 8,5% coming from the RKD images collection<sup>1</sup>, a difference that could be attributed to the difficulty and ambiguities in the procedures to report such complex information. We propose a way to simplify, streamline, and homogenise such complexity, hoping to increase the abundance, richness, and correctness of the representation of such phenomena in Wikidata.

4 The paper is structured as follows: in Section 2, the state of the art is presented focusing on the representation of WLS claims in RDF (particularly in the context of Cultural Heritage) and how the representation of complex data scenarios in Knowledge Graphs (KGs) is evaluated. In Section 3, we present the approaches provided in Wikidata to encode WLS claims. Section 4 outlines the research objectives, the data acquisition process is briefly described, and the analysis of our Wikidata sample dataset is presented. In Section 5, we present our proposal for improving annotating weaker-logical status knowledge quality. Finally, in Section 6, we summarise our findings and outline our conclusions about the work.

## 27 2. State of the art

28 Public Knowledge Graphs such as Wikidata [1], DBpedia [13], Yago [14], and Google Knowledge Graph constitute publicly available collections that can be used for research, either expressing specialist knowledge or general knowledge. In particular, Wikidata is a *collaborative* public platform built and maintained by a community of contributors.

29 Weaker logical status statements are natural in many contexts covered by these KGs, but the support for their representations varies considerably. Guidelines, data modelling, and harmonisation (a particularly relevant need for open platforms) can help express them, i.e., concurrent opinions or uncertain claims. In the field of Cultural Heritage studies, the knowledge competition is intriguing. However, some online databases or data models only partially address this issue.

30 Although domain ontologies represent the domain of cultural heritage hardly ever integrate support for representing interpretations (i.e., hermeneutics) into their models [6], there are some exceptions [5, 15].

31 For instance, CIDOC CRM [5] is a conceptual model developed and maintained by the International Council of Museums (ICOM), widely adopted by many knowledge graphs in the Cultural Heritage domain [16, 17]. It offers a formal approach to express weaker logical status claims through instances of classes representing n-ary relations.

32 CIDOC-CRM [5] adopts n-ary relationships for WLS claims, e.g., via the `crm:E13_AttributeAssignment` class<sup>2</sup>. For instance, the painting "Girl reading a letter at an open window"<sup>3</sup>, has been attributed over time to Rembrandt, Hooch, and finally to Vermeer (the currently accepted attribution). Listing 1 shows each attribution (`:aa1`, `:aa2` and `:aa3`) as a `crm:E13_AttributeAssignment` which requires the use of three predicates: `crm:P140_assigned_attribute_to` to indicate the item to which an attribute or

33 <sup>1</sup><https://rkd.nl/en/explore/images>

34 <sup>2</sup><https://cidoc-crm.org/Entity/e13-attribute-assignment/version-6.2.1>

35 <sup>3</sup>[https://en.wikipedia.org/wiki/Girl\\_Reading\\_a\\_Letter\\_at\\_an\\_Open\\_Window](https://en.wikipedia.org/wiki/Girl_Reading_a_Letter_at_an_Open_Window)

relation is assigned, and `crm:P141_assigned` to indicate the attribute that was assigned or the item, `crm:P177_assigned_property_of_type` to indicate the type of property or relation that this assignment maintains to hold between the item to which it assigns an attribute and the attribute itself.

```

:aa1 a crm:E13_Attribute_Assignment ;
    crm:P177_assigned_property_of_type crm:P14_carried_out_by ;
    crm:P141_assigned ulan:500011051 ; # Rembrandt
    crm:P140_assigned_attribute_to :painting-pr ;
    crm:P4_has_time-span :XVIII_cent.

:aa2 a crm:E13_Attribute_Assignment ;
    crm:P177_assigned_property_of_type crm:P14_carried_out_by ;
    crm:P141_assigned ulan:500020229 ; # Hooch
    crm:P140_assigned_attribute_to :painting-pr ;
    crm:P4_has_time-span :1821.

:aa3 a crm:E13_Attribute_Assignment ;
    crm:P177_assigned_property_of_type crm:P14_carried_out_by ;
    crm:P141_assigned ulan:500032927 ; # Vermeer
    crm:P140_assigned_attribute_to :painting-pr ;
    crm:P4_has_time-span :1860;
    crm:P14_carried_out_by ulan:500326948. # Thore

```

Listing 1: CIDOC CRM use of n-ary relation for encoding concurring attributions

Europeana [15] stores approximately 50 million heterogeneous digitised items from museums, libraries, and archives across Europe. Data is collected by content providers (i.e., cultural institutions) using the EDM data model [4] and the use of proxies [18] allows to express conflicting information and track data provenance. However, concurrent statements are not visible on the online pages, and no mechanism is in place to determine which proxy will be made visible when multiple exist.

An interesting instance of an EDM collection is the RKD catalogue, a comprehensive collection of data about Dutch works of art throughout history. By design, RKD allows and gathers contested and discarded attributions of paintings and portraits. Although, at the moment, there is no SPARQL endpoint available for querying the collections, users can browse RKD data through an online catalogue. Interestingly, about 83,600 artwork descriptions in Wikidata have been linked to the RKD dataset via the predicate `P350:RKDimages ID`<sup>4</sup>, representing ~7,5% of the total of visual artworks in Wikidata.

Despite the support of representational definitions of weaker logical status claims in EDM, CIDOC-CRM and RDK data models, these weaker forms of information are often poorly reported (*reticence*) or are expressed in textual annotations rather than being modelled in the data structure (*dumping*) [19].

The widespread adoption of Wikidata within the Cultural Heritage community has been well-documented [20]<sup>5</sup>. Wikidata is seen not only as a valuable tool for data publishing, alignment and enrichment but also as a means of gaining valuable insights into Cultural Heritage data and the community itself [21]. Given the significance of comprehensive data in knowledge bases, there has been a focus on improving and evaluating their schema and data quality [22]. In this context, weaker logical status claims may make good use of reification methods and several studies have been performed to improve their usage, e.g., by [23], who compared the efficiency of several reification methods (e.g., singleton properties, n-ary relations, named graphs and standard reification) on Wikidata data.

Handling WLS in Semantic Web data can be placed within the broader topic of representing and reasoning over data enriched with metadata or contextualised data. The matter has been discussed at length from many different angles. A primary objective is that of reconciliation or integration of multiple data sources. Indeed, effective representation and reasoning about knowledge with heterogeneous viewpoints is one of the objectives for applications concerned with distributed knowledge sources. Yet, semantic web ontologies force a unique, global view of the represented world, in which the axioms are meant to be interpreted as universally true. The same domains are often

<sup>4</sup><https://w.wiki/7wfW>

<sup>5</sup>The list of cultural institutions involved in Wikidata can be found at <https://www.wikidata.org/wiki/Wikidata:GLAM>

1 modelled differently depending on the intended use of an ontology. The problem of reconciliation, therefore, is to  
2 bring different world views together to create a single, unified model for representation and reasoning. This may  
3 be obtained through formal Interoperability Systems [24] extending the expressive reach of Description Logic, or  
4 bridge rules mapping separate contexts determining how the local concepts in the two ontologies map onto each  
5 other [25], or extended representation models such as RDFS with Annotations [26]. Different approaches, such as  
6 colouring [27] or NDFluents [28], or RDF+ [29], on the other hand, are less interested in obtaining reconciliation  
7 and more in representing adequately the semantics of inferences about heterogeneous claims.

8 To the best of our knowledge, current research has not extensively tackled WLS representation in RDF. However,  
9 the representation of complex data scenarios in knowledge bases (and in particular, in Wikidata) has been evaluated  
10 according to multiple metrics. For instance, Piscopo and Simperl [30] survey quality metrics from 28 scientific pub-  
11 lications on the topic and categorise quality assessments into three dimensions: intrinsic (accuracy, trustworthiness,  
12 consistency), context (relevance, completeness and timeliness) and representation (ease of understanding and inter-  
13 operability). Among quality measures, evaluation of completeness, defined by Faerber et al. [31] as the “presence  
14 of all required information in a given dataset”, has been approached through various methods and assessments as  
15 comparing data for similar entities [32], measuring entity relatedness [33], evaluating thoroughness of information  
16 by determining the completeness of specific attributes of objects [34], assessing low-quality statements through the  
17 analysis of items’ discussion pages, deprecated statements and constraint violations [35], and assessing and compar-  
18 ing data quality across large knowledge bases [31, 36]. Additionally, Arnaut et al. [37] surveyed negative knowledge  
19 in Wikidata, analysing deleted statements, count predicates, deprecated statements, negated predicates and noValues  
20 to measure Wikidata completeness from this point of view.

21 Overall, among current Cultural Heritage KGs, WLS representation seems to be slightly tackled, showing Wiki-  
22 data as one of the few platforms providing designed approaches to represent such knowledge. However, little or no  
23 evaluation has been conducted specifically on the representation of weaker logical status claims in Wikidata, nor  
24 has a comprehensive analysis been carried out to assess the amount of knowledge related to WLS status in Cultural  
25 Heritage. In the next section, we detail our proposal to address these shortcomings.

### 26 27 28 **3. Representing weaker logical statuses in Wikidata**

29  
30 Wikidata represents weaker logical status statements (e.g., for uncertain or debated assertions) using at least  
31 three different approaches: ranked statements (Section 3.1), statements with specific qualifiers (Section 3.2) and  
32 statements with a non-existing valued object (Section 3.3).

#### 33 34 *3.1. Ranked statements*

35  
36 Ranking of assertions is modelled by the Wikibase data model<sup>6</sup> to express different degrees of the preferability  
37 of individual claims.

38 Claims in Wikidata are expressed through *statements*, a custom reification approach<sup>7</sup> [23] to express contextual  
39 information (e.g., qualifiers, rankings, references) about it. Statements connect the claim’s subject and predicate to  
40 a Statement entity, which refers to the claim’s object and can be further used as the subject of other triples.

41 Statements do not assert the corresponding claim, but an additional triple must be added to assert the claim’s  
42 content. The additional triple (which uses a different prefix) flatly relates the statement’s subject to the statement’s  
43 intended object through the statement’s predicate, thus enabling simple query support for asserted facts. The separa-  
44 tion between statements and their assertion is selectively provided, allowing easy support for both claims presented  
45 as facts (where both the statement and the assertion triple exist) and claims not meant to be considered facts (the  
46 statement exists, but no assertion triple is added).

47 The ranking mechanism is enriched with the representation of asserted and non-asserted statements. Rankings  
48 [38] communicate the scientific community’s or Wikidata annotators’ consensus. Disputes are separately hosted

50 <sup>6</sup><https://www.mediawiki.org/wiki/Wikibase/DataModel#Statements>

51 <sup>7</sup><http://www.wikidata.org/wiki/Help:Statements>

in plain text on the corresponding discussion page. Many possible combinations of variously ranked competing statements can be found in the Wikidata collection, with various and debatable interpretations. Ranking is assigned to individual statements using values such as *preferred*, *normal* and *deprecated*).

Note that whether or not a statement is asserted is determined solely by its rank and the absence of higher-ranked statements using the same predicate. The Wikidata engine automatically asserts the statement and it is not the editors' conscious choice.

### 3.1.1. Normal statements

The normal ranking is the default ranking for statements. A statement ranked normal can be either asserted or not depending on the existence and intended meaning of competing statements against it. For instance, in Listing 2, "The Scream" by Edvard Munch belongs to the Expressionist period<sup>8</sup>, and this is expressed as an asserted normal statement, to signify that the annotator does not give a WLS status to the statement. In Listing 4, on the other hand, the first statement (lines 1-5) is ranked normal but not asserted since the preferred statement is present and asserted instead.

```

1 # "The Scream" belongs to the Expressionist movement
2 wd:Q471379 wdt:P135 wd:Q80113 .
3 wd:Q471379 p:P135 s:Q471379-c3e5c17d-4730-a5dc-85cb-efc9766b7c80 .
4 s:Q471379-c3e5c17d-4730-a5dc-85cb-efc9766b7c80 a wikibase:Statement,
5   wikibase:rank wikibase:NormalRank ;
6   ps:P135 wd:Q80113 .

```

Listing 2: Normal rank

### 3.1.2. Deprecated statements

Deprecated statements are meant for claims with a weak logical status and do not represent a correct value in the editors' view. Deprecated statements are always automatically non-asserted independently of the ranking of the other concurring statements. Wikidata designed use for deprecated ranking is stated to be "used for statements that are known to include errors (i.e. data produced by flawed measurement processes, inaccurate statements) or represent outdated knowledge (i.e. information that was never correct, but was at some point thought to be)". Additionally, Wikidata negates the use of deprecated ranks for claims which describe "correct historical information, such as previous values of a statement [...]"<sup>9</sup>.

For instance, Listing 3 expresses the concept that "The Lamentation"<sup>10</sup>, a print by Albrecht Dürer, was reported to be created in 1504. The deprecated rank and the lack of an asserted triple indicate that this date is invalid.

```

1 # creation date thought to be 1504
2 wd:Q18338462 p:P571 s:Q18338462-FDDCD91B-3919-450A-B00D-FE3ADA773A11 .
3 s:Q18338462-FDDCD91B-3919-450A-B00D-FE3ADA773A11 a wikibase:Statement ;
4   wikibase:rank wikibase:DeprecatedRank ;
5   ps:P571 wdt:P571 "1504-01-01T00:00:00Z"^^xsd:dateTime .# creation date: 1504

```

Listing 3: Deprecated rank

### 3.1.3. Preferred statements

Preferred statements are meant for claims with a stronger status and representing the currently presumed correct value of a predicate. They are always also asserted. For instance, as shown in Listing 4, a retracted attribution of the painting "Madonna with the Blue Diadem"<sup>11</sup> to Raphael is represented only by a statement ranked as normal and no assertion triple, while the attribution to Gianfrancesco Penni enjoys a preferred rank, and the assertion triple.

<sup>8</sup><http://www.wikidata.org/entity/Q471379>

<sup>9</sup>[https://www.wikidata.org/wiki/Help:Ranking#Deprecated\\_rank](https://www.wikidata.org/wiki/Help:Ranking#Deprecated_rank)

<sup>10</sup><http://www.wikidata.org/entity/Q18338462>

<sup>11</sup><http://www.wikidata.org/entity/Q738038>

Even though the first attribution is ranked normal rather than deprecated, we must consider it a superseded claim. This example shows that the nature of normal statements varies depending on whether they coexist or not with competing preferred and/or deprecated claims, and similarly, the presence or absence of assertion triples may vary. The preferred rank designed use is “most current statement”, implying that other concurring statements should represent outdated statements, and “statement that best represents consensus (be it scientific consensus or the Wikidata community consensus)”, implying that other concurring statements should represent concurring discarded statements<sup>12</sup>.

```

1      # attribution to Raphael
2      wd:Q738038 p:P170 s:q738038-121B92D0-E6E1-4514-960C-AE34F50054E5 .
3      s:q738038-121B92D0-E6E1-4514-960C-AE34F50054E5 a wikibase:Statement ;
4      wikibase:rank wikibase:NormalRank ;
5      ps:P170 wd:Q5597 .          # creator: Raphael
6
7      # attribution to Gianfrancesco Penni
8      wd:Q738038 wdt:P170 wd:Q2327761 .      # creator: Gianfrancesco Penni (assertion)
9      wd:Q738038 p:P170 s:Q738038-7729b786-4d4f-a0ca-2ded-4ea2c6307e1c .
10     s:Q738038-7729b786-4d4f-a0ca-2ded-4ea2c6307e1c a wikibase:Statement;
11     wikibase:rank wikibase:PreferredRank ;
12     ps:P170 wd:Q2327761.          # creator: Gianfrancesco Penni

```

Listing 4: Preferred and Normal ranks

### 3.2. Qualifiers

Statements, independently of rank, can be decorated with additional triples annotating contextual information or specifications about the claim itself<sup>13,14</sup>. Those annotations may be *additive* when they provide additional information about the fact (e.g., to specify the character played by an actor when listing them as a cast member of a movie) or *contextual* when they limit the contexts in which the underlying fact is true (e.g., the claim is a hypothesis) [39]. Wikidata states the designed use of qualifiers as “to represent a plurality of perspectives on Wikidata, including data which may provide contradicting information. In disputes, community consensus ultimately determines the value of a property. However, other points of view can be added as additional values using qualifiers (as well as sources). Ranks can also be used; if a consensus exists, it should be indicated by a preferred rank”<sup>15</sup>. Additionally, among the designed uses of qualifiers with a single value, Wikidata allows the usage “to constrain the validity of the value(s)”<sup>16</sup>. No specific designed use is provided for uncertainty-related qualifiers (e.g. possibly).

Following the example from Aljalbout et al. [40], we examined the 150 most frequently used qualifiers in Wikidata and their most commonly used values. The most used qualifiers to use WLS values are `P1480:sourcing circumstances`<sup>17</sup> (47th most used one) and `P5102:nature of statement`<sup>18</sup> (134th most used one). Additionally, the Wikidata model provides the properties `P2241:reason for deprecated rank`<sup>19</sup> (42nd most used qualifier) and `P7451:reason for preferred rank`<sup>20</sup> (114th most used qualifier) to annotate contextual information about superseded and preferred claims, respectively.

For instance, in Listing 5, we see that the painting “Abstract Speed + Sound”<sup>21</sup> by Giacomo Balla is described as *possibly* part of a triptych. Using a qualifier with a normal ranking seems to imply that the statement is considered true and, therefore, asserted.

<sup>12</sup>[https://www.wikidata.org/wiki/Help:Ranking#Preferred\\_Rank](https://www.wikidata.org/wiki/Help:Ranking#Preferred_Rank)

<sup>13</sup><https://www.wikidata.org/wiki/Help:Qualifiers>

<sup>14</sup>The complete list of available qualifiers in Wikidata is available at <https://w.wiki/6TrP>

<sup>15</sup>[https://www.wikidata.org/wiki/Help:Qualifiers#For\\_disputed\\_items\\_&\\_community\\_consensus](https://www.wikidata.org/wiki/Help:Qualifiers#For_disputed_items_&_community_consensus)

<sup>16</sup>[https://www.wikidata.org/wiki/Help:Qualifiers#For\\_single\\_values](https://www.wikidata.org/wiki/Help:Qualifiers#For_single_values)

<sup>17</sup>The most frequently used values are: *circa, presumably, allegedly, inference, uncertainty, possibly, near, probably, conventional date, disputed*

<sup>18</sup>The most frequently used values are: *originally, attribution, hypothesis, often, allegedly, expected, possibly, disputed, rarely, mainly*

<sup>19</sup>[http://www.wikidata.org/wiki/Property\\_talk:P2241](http://www.wikidata.org/wiki/Property_talk:P2241)

<sup>20</sup>[http://www.wikidata.org/wiki/Property\\_talk:P7451](http://www.wikidata.org/wiki/Property_talk:P7451)

<sup>21</sup><http://www.wikidata.org/entity/Q19882431>

```

1      1      wd:Q19882431 wdt:P361 wd:Q79218 .      # part of: triptych (assertion)
2      2      wd:Q19882431 p:P361 s:Q19882431-1ac26ff2-4981-ff79-4fae-9d411ae34296 .
3      3      s:Q19882431-1ac26ff2-4981-ff79-4fae-9d411ae34296 a wikibase:Statement;
4      4      wikibase:rank wikibase:NormalRank ;
5      5      ps:P361 wd:Q79218 ;      # part of: triptych
6      6      pq:P5102 wd:Q30230067 .      # circumstance: possibly

```

Listing 5: A qualified statement in Wikidata

Wikidata provides a list of 96 recommended values for *nature of statement* and 83 recommended values for *sourcing circumstances* in their respective *Property Talk* pages. In contrast, no recommended terms are provided for *reason for deprecated rank* nor *reason for preferred rank*. However, terms that were used with these properties can be retrieved via a simple SPARQL query<sup>22</sup>, showing respectively 384 and 83 distinct terms. Even at first glance, it is possible to notice an extensive range of types and specificities (e.g., qualifiers such as *possibly*, *presumably*, and *probably* versus, say, *prosopographical phantom*, *project management estimation* or *archive footage*), and many are not connected to weaker logical status assessments. In addition, semantic overlaps can be noticed on many of these terms, e.g., between *allegation* and *allegedly*, or between *hypothesis*, *hypothetical entity*, *hypothetically* and *scientific hypothesis*. These overlaps support arbitrariness of choice for contributors, increasing the ambiguity of the resulting annotation.

### 3.3. Missing values

There are three types of basic information structures used to describe entities in Wikibase (called SNAKs, or *Some Notation about Knowledge*<sup>23</sup>) in Wikidata: actual values (URIs or literals), `someValue` placeholders and `noValue` placeholders. They are used to represent that the statement is associated with an unknown value (mapped as `someValue`) or with a non-existing value (mapped as `noValue`), which is a more precise assessment than simply not recording the statement at all. The same syntactic tool is known to generate precision and correctness issues (e.g., see Hernandez [41]) since the RDF standard specifically defines blank nodes with existential semantics. At the same time, SPARQL does not follow such semantics. Wikidata declares `someValue` and `noValue` claims intended use as: “There are times when for a given property an item has either no value (the absence of that property) or an unknown value. These data values may still provide important information about the item and, if so, should be recorded in Wikidata. For instance, we should say that Elizabeth I of England (Q7207) has no value for the child (P40) property, which is a positive statement that she had no children. We should also say that William Shakespeare (Q692) has an unknown value for the date of birth (P569) property, which is a positive statement that that information has not been preserved.”. Additionally, Wikidata defines the `noValue` SNAK as “in some cases, we want to emphasize that a property value has not just been left out (or not entered yet) but that it really does not exist”<sup>24</sup> and `someValue` SNAK as “the information that a property has some value can be important and useful, even if the value is not known”<sup>25</sup>.

#### 3.3.1. Unknown values

Unknown valued statements are claims whose object exists but is not known<sup>26</sup>. For instance, in “The Book of Lismore”<sup>27</sup> there is an unknown value for the `P195:collection` property, which is a positive statement that the information existed but it has not been preserved. As mentioned, unknown values are represented in RDF via blank nodes as shown in Listing 6.

<sup>22</sup>List of terms used in Wikidata with *reason for deprecated rank* <https://w.wiki/6Tpt> and with *reason for preferred rank* <https://w.wiki/7VGf>

<sup>23</sup><https://www.wikidata.org/entity/Q86719099>

<sup>24</sup><https://m.mediawiki.org/wiki/Wikibase/DataModel#PropertyNoValueSnak>

<sup>25</sup><https://m.mediawiki.org/wiki/Wikibase/DataModel#PropertySomeValueSnak>

<sup>26</sup>[https://www.wikidata.org/wiki/Help:Statements#Unknown\\_or\\_no\\_values](https://www.wikidata.org/wiki/Help:Statements#Unknown_or_no_values)

<sup>27</sup><https://www.wikidata.org/wiki/Q1371647>



```

1 wd:Q1371647 wdt:P195 _:15518d67963a082b352304a1ab8e016e. # unkown collection
2
3 wd:Q1371647 p:P195 s:Q1371647-B07F6386-A7D0-4C9D-8E77-CC2BD523354E .
4 s:Q1371647-B07F6386-A7D0-4C9D-8E77-CC2BD523354E ps:P195 _:0088bc50e53b3902bea74cc2380cbd09 ;
5 pq:P3831 wd:Q768717 . # the role of this collection is to be a private collection

```

Listing 6: Unknown-valued statement in Wikidata

### 3.3.2. Non-existing values

Non-existing valued statements<sup>28</sup> are claims whose object is not existent (or not available in Wikidata). For instance, the pilot episode of X-files<sup>29</sup> has a non-existing value for the *follows* (P155) property, considering that the pilot starts the series. Non-existing values do not create additional values but are represented by making the statement node (and, for asserted claims, the entity itself) an instance of a class “wdno:P???”<sup>30</sup>, where the “???” corresponds to the relevant property id<sup>30</sup>. Non-existing values are not conceptually a WLS claim, but we list them in this survey because there exists in practice some overlap between unknown valued and non-existing valued claims. For instance, the “Les amours de Cartouche”<sup>31</sup>, a literary work from the 18th century, has been recorded with an unknown author (supposedly to mark its anonymous author), as shown in Listing 7. The example is incorrect as it should use an unknown value. This leads to confusion about the usage of missing values, further contributing to complications.

```

1 wd:Q123914909 a wdno:P50. # author: unknown (noValue)
2 wd:Q123914909 p:P50 s:Q123914909-592511E6-FF3D-454B-A2C2-9D7A9207C6A0 .
3
4 s:Q123914909-592511E6-FF3D-454B-A2C2-9D7A9207C6A0 a wdno:P50 ;
5 pq:P1932 "no value" .

```

Listing 7: non-existing valued statement in Wikidata

### 3.4. Discussion

Even before checking on the actual usage patterns of these approaches, we can immediately notice the richness of annotations made possible by them, the subtle nuances they afford, and the variety of (potential) sources of ambiguities, overlapping connotations and representation vagueness. In particular, we can summarise three specific problems that are worth further discussion:

1. Although the separate uses of normal, preferred and deprecated rankings are clear and practical, there are uncertainties when they coexist on the same predicate, especially for the different representations of normal statements when preferred ones are also present or when all three rankings are present.
2. The sheer number of qualifiers, the differing levels of their respective specificities, and the manifest semantic overlapping of many of them make it hard to guarantee homogeneity and precision in their use. Contextualising qualifiers, be they temporal, provenance or otherwise, does not add to the base information but changes the context within which such information is true. As Patel-Schneider [39] suggests, contextual qualifiers should not be shown to consumers. Still, basic tools (visualisers, contextualisers, reasoners) should be written to take such context into account correctly, and low-level tools should remove facts that are not valid in the selected contexts.
3. The subtlety in the semantic differences between providing no statement, specifying a `noValue` and providing a `someValue` for a property of a Wikidata item, as well as their other types of applications makes the use of missing values potentially ambiguous.

<sup>28</sup><https://www.wikidata.org/wiki/Help:Statements>

<sup>29</sup><http://www.wikidata.org/entity/Q7194381>

<sup>30</sup>[https://www.mediawiki.org/wiki/Wikibase/Indexing/RDF\\_Dump\\_Format#Novalue](https://www.mediawiki.org/wiki/Wikibase/Indexing/RDF_Dump_Format#Novalue)

<sup>31</sup><http://www.wikidata.org/entity/Q123914909>

In a way, WLS claims can be seen simply as logical disjunctions of competing claims each of which is separately annotated with context, provenance, confidence, temporal boundaries, etc.: “according to  $\alpha$ ,  $s p o_1$ ” and “according to  $\beta$ ,  $s p o_2$ ” can be seen as “[ $s p o_1$ ] $_{\alpha}$   $\vee$  [ $s p o_2$ ] $_{\beta}$ ” with some added annotations connecting the first branch to  $\alpha$  and the second to  $\beta$  (e.g., through reification, named graph, or blank nodes). This approach has limitations from the practical and the conceptual points of view. Practically, RDF has no real way to express disjunctions without some additional baggage to encode predicate calculus employing the systematic use of reification [42]. Conceptually, focusing on the inner statements to the exclusion of the contextualising information may miss the point that in many scholarly domains, it is not the full list of competing claims to be of interest but the very existence of the diatribe in the first place. Disjunctions would not help here.

Another way to formally understand WLS claims is to link them to modal statements in modal logic [43], which can be used to understand the coexistence of strong logical status claims, expressed as atomic formulas  $p(s, o)$ , and weak logical status ones, expressed as modal formulas  $K_{\alpha}p(s, o)$  or  $B_{\beta}p(s, o)$ , where  $K_{\alpha}$  and  $B_{\beta}$  are modal operators guided by specific modal axioms<sup>32</sup>. Various types of modal logics exist and have been used to introduce different operators and represent different semantics, such as possibility and necessity (the *strictu sensu* modal logic), or obligation and permission (*deontic logic*), or temporally bounded predicates (*temporal logic*), or belief (*doxastic logic* or knowledge (*epistemic logic*). Overall, they form a complete formal mechanism to study the characteristics and principles of WLS claims that does not imply the need to proceed to a reconciliation of different world views.

Yet, all these reflections are empty and pointless unless we examine how contributors use these approaches to express real WLS claims in their Wikidata contributions. The following section covers this topic.

#### 4. Usage patterns of WLS in Wikidata datasets

To generate some analysis about the actual usage of WLS claims and to provide an initial answer to our research questions, we collected three datasets of Wikidata items: one about Cultural Heritage items (visual arts, text documents and audio-visual entities), another about Astronomical objects (galaxies and stars) and one with a selection of random entities reflecting the actual distribution of entities in classes in the whole Wikidata as discussed in Section 1.

The datasets were selected to be approximately comparable in size, and the number of individual statements and under evidence that many types of entities rely on weaker logical status claims when entities undergo re-evaluations due to new pieces of evidence or the recording of different opinions.

##### 4.1. Data Acquisition

```
SELECT DISTINCT ?artwork ?type
WHERE {
  ?artwork wdt:P31 ?type.
  ?type (wdt:P279*) wd:Q838948.
  hint:Prior hint:rangeSafe true
}
```

Listing 8: SPARQL query retrieving Wikidata entities to subclasses of work of art (Q838948)

The first dataset contains Cultural Heritage items (CH), a complete snapshot of the Wikidata records of these cultural assets. All Wikidata entities belonging to the class *work of art*<sup>33</sup> or any of its sub-classes were collected using a SPARQL query (Listing 8). The statements for all selected entities were downloaded in JSON format<sup>34</sup>. Data is stored in numerous JSON files, and each contains a complete representation of at most 50 Wikidata entities with their labels, descriptions and statements. This Cultural Heritage dataset has been semi-automatically divided into three sub-datasets due to the wide diversity of cultural properties and their associated claims:

<sup>32</sup>e.g.,  $\mathbf{T}$  ( $K_{\alpha}\phi \rightarrow \phi$ ) for epistemic logic or  $\mathbf{N}$  ( $\vdash \phi \implies \vdash B_{\alpha}\phi$ ) for doxastic logic.

<sup>33</sup><http://www.wikidata.org/entity/Q838948>

<sup>34</sup>via [http://www.wikidata.org/wiki/Wikidata:Data\\_access](http://www.wikidata.org/wiki/Wikidata:Data_access)

- 1 – *Audio-Visual heritage* (CHav): This collection holds information about audio-visual materials that have cultural, historical, or artistic value. They include movies, videos, recordings of music or spoken words, and other audio-visual materials that record a particular event in a specific time or place. The dataset contains 1,251,626 entities and 17,141,394 statements organised in 25,033 JSON files.
- 2 – *Visual heritage* (CHv): This collection holds information about visual artefacts with cultural, historical, or artistic value. They include paintings, drawings, sculptures, photographs, decorative arts, etc. The dataset contains 1,078,855 entities and 12,850,825 statements organised in 21,579 JSON files.
- 3 – *Textual heritage* (CHt): This collection holds information about written and printed materials with historical or cultural significance. They include books, manuscripts, letters, and other written documents. The dataset contains 625,110 entities and 4,584,444 statements organised in 12,503 JSON files

We also downloaded Wikidata entities of architecture-related classes; they were later discarded due to their fairly lower number as well as for the presence of many statistical ambiguities that could make their evaluation useless (e.g., many entities belonging to these classes should not be considered relevant to Cultural Heritage collections).

The second dataset, chosen to verify our assumptions using a different collection with a similar size, is a collection of astronomical entities organised into two datasets:

- *Stars* (ANs): This collection holds a random selection of 1,199,950 Wikidata entities (of the ~3.3 million existing) belonging to the class *Star*<sup>35</sup>, The dataset contains 27,470,140 statements in 23,999 JSON files<sup>36</sup>.
- *Galaxies* (ANg): This collection holds a random selection of 1,200,000 Wikidata entities (of the ~2 million existing) belonging to the class *Galaxy*<sup>37</sup>, The dataset contains 14,439,421 statements in 24,000 JSON files.

We decided to limit the number of astronomical entities to 1,200,000 to approximately balance them to each other (although the CHt is about half in size with 625,110 entities), as well as the average number of statements for each entity (CHav: 13.7, CHv: 11.9, CHt: 7.3, ANs: 22.9, ANg: 12).

The third dataset is a selection of randomly chosen entities from Wikidata. This dataset was acquired to compare WLS claims in the other datasets with a randomised subset designed to mimic the overall distribution of WLS claims in Wikidata.

- *Random* (R): This dataset comprises 1,159,800 Wikidata entities (starting from a selection of 1.2 million entities from which we removed duplicates) chosen randomly from the most numerous 100 classes to reflect the proportional distribution of entities found in Wikidata<sup>38</sup>. This dataset encompasses 61,798,072 statements distributed across 23,196 JSON files.

In Table 1, we summarise basic information about these collections. All these datasets can be accessed and downloaded from Zenodo<sup>39</sup> [44] and all Python scripts are accessible in GitHub<sup>40</sup>.

## 4.2. Analysis

In the following, we will describe as WLS statements all Wikidata statements showing the use of each approach described in Section 3, regardless of whether they have been used to make weaker logical status claims. Table 1 shows a tabular presentation of our analysis.

<sup>35</sup><http://www.wikidata.org/entity/Q523>

<sup>36</sup>the ANs dataset was meant to be composed of 24,000 files with 50 entities each, but after running our tests we noticed that a file was corrupt and we chose to discard that contribution.

<sup>37</sup><http://www.wikidata.org/entity/Q318>

<sup>38</sup><https://w.wiki/7iCR>

<sup>39</sup><https://doi.org/10.5281/zenodo.7624783>

<sup>40</sup>[https://github.com/alessiodipasquale/Wikidata\\_WLS](https://github.com/alessiodipasquale/Wikidata_WLS)

	Cultural Heritage			Astronomy		Random (R)
	Audio-visual (CHav)	Visual (CHv)	Textual (CHt)	Stars (ANs)	Galaxies (ANg)	
Entities	1,251,626	1,078,855	625,110	1,199,950	1,200,000	1,159,800
Statements	17,141,394	12,850,825	4,584,444	27,470,140	14,439,421	61,798,072
Weaker Logical Status (WLS)	<b>50,193</b>	<b>227,218</b>	<b>17,216</b>	<b>7,532,169</b>	<b>721,504</b>	<b>1,101,014</b>
% WLS / Statements	<b>0.29%</b>	<b>1.77%</b>	<b>0.37%</b>	<b>27.42%</b>	<b>5.00%</b>	<b>1.78%</b>
Non-asserted statements	43,211	9,056	14,055	7,532,107	721,503	1,089,469
Ranked as Deprecated	7,622	3,057	1,568	2,768,829	189,691	721,870
Deprecated with a reason	4,949	769	715	2	0	8,993
Non-existing values	50,611	1,969	1,356	4	0	3,857
Unknown value	4,896	106,521	1,843	0	0	5,139
Qualified statements	2,406	114,674	1,556	532	1	7,716
WLS qualified statements	2,086	111,641	1,318	62	1	6,406
WLS qualifiers w/o <i>circa</i>	719	3,988	330	35	0	1,724

Table 1

Entities, statements and types of WLS statements

Even though critical analysis is a pivotal part of humanities discourses, plainly stated statements with no competing claims are largely the most represented information in the CH dataset: the vast majority of statements here (>99%, in particular 99.74% in CHav, 99.92% in CHv and 99.69% in CHt) are plainly asserted statements with no WLS additions. In contrast, the Astronomical datasets show a reasonably different situation, 83% overall of plainly asserted statements, specifically ANs at 72.58% and ANg at 95%. The overall distribution of the Random (R) dataset showcases a low percentage of WLS claims (1.78%), closer to the CH and the AN datasets. Yet, interestingly, almost the whole percentage is made of non-asserted statements (98.95%) matching a similar distribution in the AN dataset.

When analysing the Random (R) dataset, we notice that the ranking system's simplicity leads to a clear predominance of deprecated items and, consequently, of non-asserted claims. The other approaches appear to be under-utilised in a proportion closer to the AN dataset. Possibly, this is a reflection that, in the CH community, historical uncertainty and the representation of interpretation are more frequent and typical than in other disciplines.

To further explore these data, we can notice that:

**Non-asserted statements:** Of the approaches previously listed (cf. Section 3), non-asserted statements (i.e., variously ranked statements with no corresponding asserted triples) are largely the most frequent approach for representing competing information in both AN and R. The situation is fairly different in the CH collections, non-asserted statements being the most frequently used approach in CHt (81.64%) and CHav (only 86.09%) and almost unused in CHv (3.99%).

**Deprecated statements:** Deprecated claims are visibly a small portion of the overall non-asserted statements, occurring only in 20% of the non-asserted statements of the Cultural Heritage entities, in 30% of the non-asserted statements of Astronomical entities and in the 66% of the non-asserted statements of Random entities. At the same time, about half of the deprecated statements were annotated with the corresponding *reason for deprecated rank* qualifier (in particular, 45.59% CHt, 25.15% CHv, 64.93% CHav – compare this with basically 0% in both AN datasets and 1.24% in R dataset), proving that scholars in the humanities have a solid interest in annotating provenance of WLS claims on CH data. Yet, only less than 1% of preferred statements have been annotated with the corresponding qualifier *reason for preferred rank*.

**Unknown values:** Unknown valued statements are not used at all in Astronomical data (absolute 0 in both ANg and ANs out), poorly adopted in the R dataset (0.47%), and sparsely used in the Humanities as well (9.75% in CHav and 10.71% in CHt). Higher is the result for the CHv dataset, with 46.88% of the overall WLS claims using this approach.

**Non-existing values:** Even if they do not represent WLS claims, we examined them in our datasets for contiguity to unknown values. Non-existing values are almost unused in Astronomical data (exactly 4 occurrences in ANs and an absolute 0 in ANg out of more than 7 million WLS claims) and very sparsely used in the Humanities and Random datasets as well: 1.969 statements in CHv, 1.356 statements in CHt and 3.857 statements in R dataset. Fairly higher

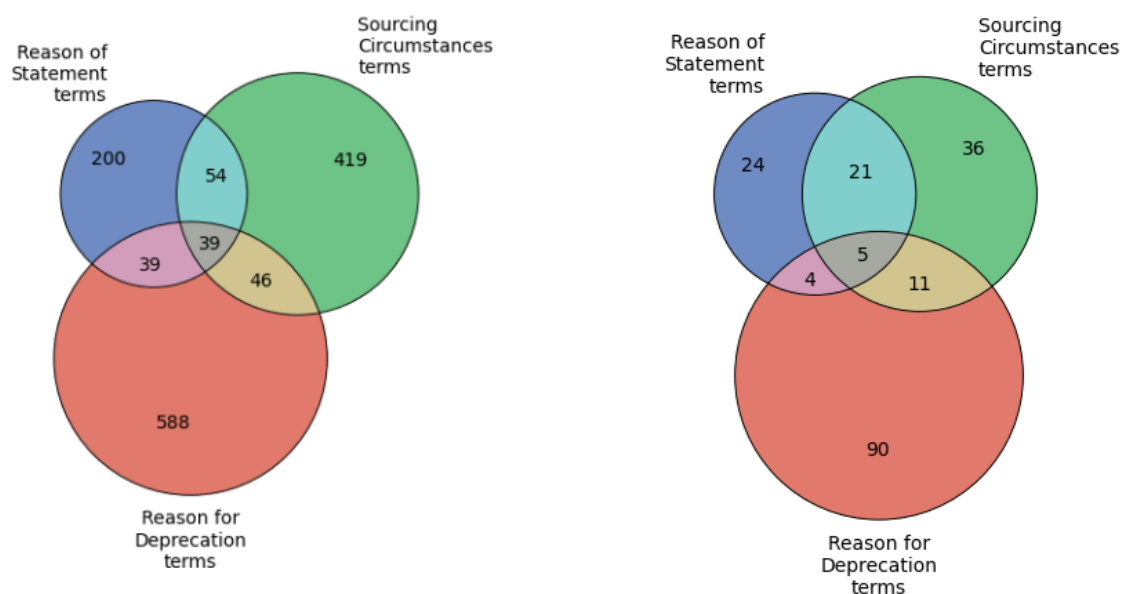


Fig. 1. Terms used in qualifiers *nature of statement*, *sourcing circumstances* and *reason for deprecated rank* throughout Wikidata (left) and in the CH datasets (right)

is the result for the CHav dataset, with 50,611 statements using this approach. This outlier value is probably justified and will be commented on later in this section.

**Qualifiers:** Statements qualified with *nature of statement* and *sourcing circumstances* predicates are the least employed out of the surveyed ones, being used in 7.66% of the WLS statements in CHt, in 0.58% of the WLS statements in R and in 4.16% of the CHav statements, present in 0.0008% of the ANs statements and only in one ANg statement. Yet, they are used in 49.13% of the WLS statements of the CHv dataset. This value will be commented later on in this section.

We further surveyed the terms actually used as values for the qualifiers.

We witnessed the use of respectively 200 different values for qualifier *nature of statement*, 419 for *sourcing circumstances* and 588 for *reason for deprecated rank*. These values largely exceed the proposed values specified in the corresponding Wikidata property talk pages (respectively, 194 values for *nature of statement* and 175 for *sourcing circumstances*) or property constraints as for the 384 values for *reason for deprecated rank*). Furthermore, the three sets of actual terms show a considerable overlap of values between them (in our datasets, but also over all of Wikidata), as shown in Figure 1. This seems to imply that the semantics associated with these values, and indeed the properties themselves, may have been unclear to contributors, who then, in some cases, selected the qualifier in non-predictable ways. Therefore, we decided to group all three sets into a single category (shown as *WLS qualified statements* in Table 1).

Since the R dataset is not disciplinary, we deemed that the variety of situations occurring across disciplinary boundaries would inevitably pollute any analysis deeper than mere counting, and therefore, in the following sections, we will focus only on the disciplinary datasets.

We further surveyed the terms actually used as values for the qualifiers.

Overall, the three sets contain a variety of terms such as generic contextual information items, e.g., provenance details, as well as domain-specific terms not relevant to our purposes (e.g., *show election*, *declared deserted*, or *text exceeds character limit*), as well as qualifiers we can truly consider suggesting weaker logical statuses (e.g., *possibly*, *disputed*, *expected*, etc.).

Therefore, we ignored the suggested values provided by the *Property Talk* pages and focused on the actual values found in our datasets. We surveyed the list of terms and selected a subset of 101 terms that seem to concretely refer to

WLS claims. This subset of WLS terms appears to be widespread in CH and Random datasets (2,086 occurrences in CHav, 111.641 occurrences in CHv, 1,318 occurrences in CHt and 6,406 occurrences in R), while almost not employed in Astronomical datasets (62 occurrences in ANs and only 1 in ANg).

The distribution of approaches to represent WLS claims in the CH dataset is not homogeneous, as unknown values and WLS-qualified statements are both highly used in the CHv dataset, while non-asserted statements for CHav and CHt. An obvious outlier is the use of one specific qualifier. Indeed, the value *circa*<sup>41</sup> is by far the most employed value in CHv, appearing 107,653 times in *sourcing circumstances*. This brings the overall count of this value completely out of scale concerning other values (e.g., the second most frequent WLS term in CHv is *probably*, occurring only 1.676 times). By removing specifically the value “circa” from the others in the last line of Table 1, we see a much more homogeneous distribution of values across the three CH datasets. On the contrary, many other terms in the list are present only once in the whole dataset and contribute very little to the overall impact of the qualified statements.

Another outlier seems to be the number of non-existing valued statements, which are present in the CHav dataset with a much higher proportion than elsewhere. In this dataset, non-existing valued statements seem to be heavily employed correctly in specific properties that appear frequently here and not elsewhere, such as *P364:original language of film or TV show*, *P155:follows*, and *P156:followed by*. This is the correct use of the non-existing valued predicates. At the same time, in the other datasets, these properties do not appear with the same frequency, and we observe a more heterogeneous distribution of approaches (cf. Figure 2).

In theory, the approaches to represent WLS claims are **not** meant as alternatives to each other and to be used exclusively. It would be perfectly acceptable and reasonable to use them on the same statement for the same entity, e.g., to describe a deprecated non-existing valued statement that results as non-asserted. Yet, approaches co-occurrence in the surveyed datasets is poorly represented, and datasets demonstrate very few cases of use of multiple WLS approaches for the same statements. In particular, no co-occurrence can be found in the AN dataset because almost all WLS claims are expressed via ranked statements except for a little co-occurrence of deprecated statements marked with a WLS qualifier in the ANs dataset (0.1%). Co-occurrences between approaches representing WLS information seem to be poorly implied in CH datasets. Almost no co-occurrence could be found between unknown and deprecated statements (0.1% in CHav, 0.04% in CHv and none (0%) in CHt), as well as the co-occurrence of deprecated and WLS qualified statements (0.04% in CHav, 0.01% in CHv, 0.07% in CHt), as well as the co-occurrence of unknown and WLS qualified statements (in 0% in CHav, 1.14% in CHv and 0.13% in CHt).

To summarise, it becomes manifest that the prevalence of each approach is quite diverse, even between the datasets of the same domain. Specifically, in CHav the most commonly used approach representing WLS information is non-asserted (86.09%), in CHv it is the WLS Qualified statement (49.13%) followed by unknown value (46.88%), and in CHt it is non-asserted (81.64%). In the Astronomy datasets, non-asserted statements overwhelmingly represent WLS claims, but deprecated statements have a much larger impact on them than in the Cultural Heritage domain.

The property analysis provides valuable insights, too, as shown in Figure 2. We divided the actual usage of WLS approaches by the property where they appear. The x-axis contains, for each dataset, the ten most frequent properties in which WLS statements appear. The y-axis shows in logarithmic scale the number of occurrences of such statements, organised by colour: non-asserted statements (with rank normal), non-asserted statements (with rank deprecated), statements with qualifiers (only WLS-related qualifiers), and non-existing valued statements.

The datasets were analysed by systematically evaluating the properties associated with the surveyed approaches. Each dataset was analysed to identify (1) the most prominent properties of each dataset and (2) the most prominent properties of each dataset with each approach.

Normal ranked, yet non-asserted statements appear in large numbers in CHav for *P8687:Social media followers*, *P348:software version identifier*, *P175:performer* and *P1476:title*. They represent peculiar uses of the non-asserted normal ranks for statements that represent multiple, independent values for the same property, none of which is “more important” than the others. Similar reflections can be made for *P18:image* on dataset CHv, and

<sup>41</sup><http://www.wikidata.org/entity/Q5727902>

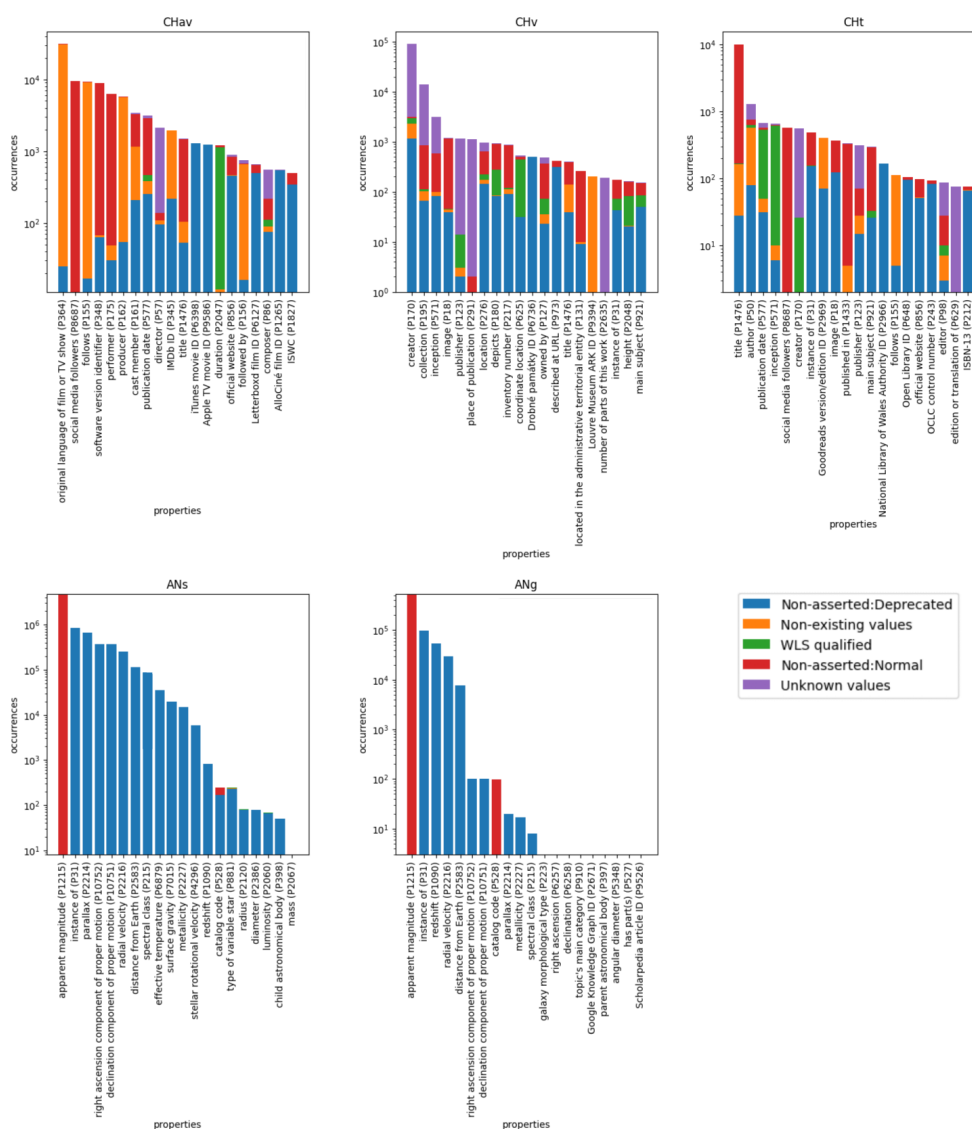


Fig. 2. Top 10 most recurrent properties implied in WLS claims in each disciplinary dataset

properties P1433:*published in* and P921:*main subject* in dataset CHt. The property P1215:*apparent magnitude* dominates this category for astronomical data. Most of the remaining properties employ a deprecated rank for evolving or uncertain information. Despite the different designed uses for deprecated and preferred rankings, Figure 2 shows that Non-asserted:Deprecated and Non-asserted:Normal claims highly co-occur with the same properties.

Qualified statements are largely present in CHv and CHt on properties P571:*inception*, P577:*publication date*, and P625:*coordinate location*, where, as mentioned, the *circa*<sup>42</sup> value for qualifier dominates the occurrences.

Unknown valued statements are primarily used in CHv and CHt datasets and only sparsely in CHav dataset. Their usage seems to be mainly implied in the description of agents in roles in all CH datasets (e.g., P170:*creator*, P98:*editor*, P123:*publisher*, P50:*author*, P86:*composer*, P57:*director*). In CHv and CHt datasets, their usage includes also locations (e.g., P195:*collection*, P291:*place of publication*), time (e.g., P571:*inception*) and the

<sup>42</sup><https://www.wikidata.org/entity/Q5727902>

artworks' description (e.g., P2635:*number of parts of this work*, P629:*edition or translation of*). The significant prevalence of unknown values when annotating agents in roles related to artworks is evident in the CHv dataset, reflecting the paramount relevance of authorship attributions given by scholars in art history.

We can also notice the predominance of non-existing valued statements in CHav (P364:*original language of film or TV show*, P155:*follows*, P156:*followed by*, P162:*producer* and P345:*IMDb ID*), which goes to prove the peculiarity of the use of non-existing valued statements in the CHav dataset previously described. The dataset CHt has a considerable number of non-existing valued statements, too, but only on properties P1476:*title* and P50:*author*, for untitled and/or anonymous documents.

Besides this, we registered some co-occurrences of the use of unknown and non-existing valued statements with the same properties (e.g., P57:*director* in CHav, P170:*creator*, P291:*publisher*, P180:*depicts*, P571:*inception*, P127:*owned by* in CHv and P50:*author*, P98:*editor*, P123:*publisher*, P577:*publication date*, as shown in Figure 2).

To summarise, we list some of the complexities and ambiguities we identified in both the CH and the AN datasets besides their designed use described by Wikidata (cf. Section 3). The list comprehends a more fine-grained distinction of WLS situations.

#### – Ranked statements

- \* *Evolving situation*: The claim is not true at the moment but was correct at some point in the past, and keeping this information is deemed interesting to maintain. For instance, the number of P8687:*social media followers* of artists and politicians, the change of P276:*location* of a movable cultural object such as a painting or a statue, or the change of its P6216:*copyright status*, may change over time. This change is recorded via differently ranked statements. For instance, the print “At the Races: Anteriei”<sup>43</sup> *star* recently shifted from copyright to the public domain. In this case, the deprecated statement was correct up to a given moment in time but is not correct anymore.
- \* *Evolving knowledge*: Because of a new observation or theory, a previous value is considered superseded. This situation is mainly connected to new observations, theories, measurements, guesses and interpretations. For instance, the introduction of a new accepted attribution of a work of art means that the previous one is now deemed as false or at least deprecated, or, in astronomy, the object “15 Orionis”<sup>44</sup> was previously considered an P31:*instance of an infrared source*<sup>45</sup>, but it is now fully considered as a *star*<sup>46</sup>; in this case, the deprecated statement has always been incorrect, but it has been decreed as such only after a specific moment in time.
- \* *Less favoured versions*: Similar claims are ranked not because they are either false or true but because one of them is preferred over the others so that they are marked as preferred and asserted while the others are non-asserted. For instance, the P1476:*titles* of textual works are often provided in different languages, and the title in the original language is marked as the preferred version, while the translated titles in other languages are not asserted. In this case, the deprecated statement is not incorrect, but it has been demoted to prioritise another one. This is not strictly a WLS situation but uses the same ranking approach as truly WLS ones.

#### – Qualified statements

- \* *Uncertainty*: For instance, the painting “Madame Antoine Arnault”<sup>47</sup> has P170:*creator* set to Jean-Baptiste Regnault<sup>48</sup> with a P5102:*nature of statement* qualifier *disputed*; Here, the statement is not certain, and competing (and incompatible) statements may be present or at least expected.

<sup>43</sup><http://www.wikidata.org/entity/Q79471408>

<sup>44</sup><http://www.wikidata.org/entity/Q6675>

<sup>45</sup><http://www.wikidata.org/entity/Q67206691>

<sup>46</sup><http://www.wikidata.org/entity/Q523>

<sup>47</sup><https://www.wikidata.org/entity/Q109252498>

<sup>48</sup><https://www.wikidata.org/entity/Q453485>



- 1 \* *Caution*: For instance, the “Frontispiece to Christopher Saxton’s Atlas of the Counties of England and Wales State I”<sup>49</sup> has the P170:creator property set to Remigius Hogenberg<sup>50</sup>, with the contributor cautioning through a P5102:nature of statement qualifier that this is only an attribution<sup>51</sup>. Here, the statement is not certain, but it implies that the proposed value may be wrong rather than positively asserting disagreements on it.
- 2  
3  
4  
5
- 6 \* *Imprecision*: For instance, the hypothetical entity “IRAS 17163-3907”<sup>52</sup> has an observed P2060:luminosity property set to “500.000 solar luminosity” with a P1480:sourcing circumstances qualifier circa; similarly, the painting “Girl Reading a Letter at an Open Window”<sup>53</sup> by Johannes Vermeer is dated (P571:inception) 14th century with a sourcing circumstances qualifier circa. For instance, the star “Altair” (Q12975) has a P1102:flattening property set to 0.2 with a nature of statement qualifier greater than; Here, the statement is certain but the value is inherently loose. One may wonder if this is truly a WLS statement or a positive statement of an imprecise value.
- 7  
8  
9  
10  
11  
12

### 13 – Missing value statements.

14

- 15 \* *Data entry errors*: Data include errors probably introduced during the annotation. For instance, the novel “Invisible Monsters”<sup>54</sup> is both attributed to Chuck Palahniuk (the actual author) and an unknown and probably erroneous entity. Here, there is a clear error in the dataset. Whether a someValue or a noValue is used is not important as they would both be errors.
- 16  
17  
18
- 19 \* *Dumping from pre-existing databases*: Some non-existing values may result from an error in the conversion or an empty field of a record after importing an existing database into Wikidata. For instance, the painting “Marshy Landscape”<sup>55</sup> has a non-existing valued statement for the P528:catalogue code property. Again, this represents an error in the dataset, so the corresponding statement should be omitted.
- 20  
21  
22
- 23 \* *The value does not exist*: For instance, the first and last entities of a sequence use properties P155:follows and P156:followed by with a non-existing value. For instance, the first episode of a TV series or the last song of a recording should have non-existing values for the corresponding properties. This is a correct use of a noValue, not a WLS claim.
- 24  
25  
26
- 27 \* *Model fitting*: When the model does not fully support the situation to be described, some arrangements were taken, such as the use of a non-existing value for the property original language of film or TV shows P364 when the entity is a silent movie. For instance see “Silent Tests”<sup>56</sup>, whose P364:original language of film or TV show predicate is non-existing valued and additionally qualified with P518:applied to part dialogue<sup>57</sup>). Here, a non-existing value is correctly used for a value not felt necessary in the model (e.g., a specific property language of dialogue to be used in sound fields and omitted for silent movies). This is, again, a correct use of the noValue claim, yet not a WLS claim.
- 28  
29  
30  
31  
32  
33
- 34 \* *The value exists but is not known*: For instance, the painting “The Welcome Home”<sup>58</sup> is marked to have an unknown P170:creator as a someValue blank node. This is probably the only true WLS use of missing value statements.
- 35  
36  
37

38 The previous list shows a series of situations where the same approaches are used for different purposes. All such purposes (except data entry errors) are legitimate. Yet, we fear that users may have trouble differentiating the purpose of each use because the approaches chosen are not sufficiently precise enough to distinguish the specific situation clearly and unambiguously. Rather than suggest forcing all different situations into a single over-encompassing

39  
40  
41  
42

---

43 <sup>49</sup><https://www.wikidata.org/entity/Q105949375>

44 <sup>50</sup><https://www.wikidata.org/entity/Q18576859>

45 <sup>51</sup><https://www.wikidata.org/entity/Q230768>

46 <sup>52</sup><https://www.wikidata.org/entity/Q540167>

47 <sup>53</sup><https://www.wikidata.org/entity/Q700251>

48 <sup>54</sup><https://www.wikidata.org/entity/Q2600527>

49 <sup>55</sup><https://www.wikidata.org/entity/Q6773948>

50 <sup>56</sup><https://www.wikidata.org/entity/Q390207>

51 <sup>57</sup><https://www.wikidata.org/entity/Q131395>

<sup>58</sup><https://www.wikidata.org/entity/Q110041706>

approach, Section 5 lists some increasingly impactful solutions to solve these ambiguities without overly revolutionising the data model.

### 4.3. Discussion

The datasets presented in the previous section and our analysis of their content allow us to reach some conclusions on the research questions specified in the introduction.

**RQ1 – How widespread are these approaches in the current state of Wikidata?** – The current state of WLS claims in Wikidata is poor. Even though Wikidata focuses on collecting and referencing the facts claimed elsewhere<sup>59</sup> [3], rather than conjectural or controversial information<sup>60</sup>, in many cases it is objective and scientifically precise to represent the complexity of uncertainty and evolving knowledge, rather than omitting information because they are not completely established. In these cases, Wikidata seems to be doing poorly, as <1% of the claims we analysed in CH datasets show weaker logical status characteristics, 5% in the ANg dataset and 27.41% in the ANs data. Of course, finding a reference that backs the uncertainty of a claim (e.g. it is disputed) can be rarer than a reference to facts that are unequivocal for their annotators. Thus, it is natural that WLS claims in Wikidata generally appear with a much lower percentage than certain facts. Nonetheless, CH datasets show a much lower figure than, e.g., the ANg and ANs datasets. Does this show an intrinsic difference in the two cultural domains, or is there something else underneath? To provide an answer to this further question, we turned to the RKD database.

RKD<sup>61</sup> holds detailed data about Dutch and Flemish paintings, drawings and prints throughout the ages, from XVI Centuries artworks to modern ones. Overall, more than 260,000 items belonging to the image collections are described in the database, and through an EDM-inspired data model, particular attention is given to multiple competing assertions, e.g., incompatible authorship attributions. Namely, RKD contains more than 317,000 recorded attributions, i.e., an average of 1.2 attributions per artwork. Thus, deprecated authorship attributions are present in about 8.5% of the works in the RKD image collection (e.g., about 290,000 current attributions vs. 27,000 discarded ones in the RKD images collection), a conspicuously higher figure than the meagre 1.77% WLS statements of the CHv dataset.

One may wonder that Dutch and Flemish collections are not representative of the full scale of worldwide types of artworks represented in the CHv dataset. Yet, they provide an interesting starting point for a further comparison. We created a sub-dataset of CHv and further analysed it to improve our understanding of this issue. First of all, it should be noted that, as mentioned, about 83,600 artwork descriptions out of the 267,238 available in RKD have been linked to Wikidata<sup>62</sup>, representing ~7.5% of the total of visual artworks in Wikidata. Thus, a dataset in Wikidata with the same artworks of RKD inevitably risks being polluted by RKD data itself. Since RKD is highly specialized in Dutch paintings from the 17th to the 20th century, the Wikidata sub-dataset we created contains European artworks painted in the same temporal period and *explicitly excluding RKD artworks*<sup>63</sup>. Wikidata stores 501,049 paintings in the interval 17-20th century not present in RKD, for a total of 340,661 attributions<sup>64</sup>. The results of such comparison are shown in Table 2. Out of the total number of Wikidata statements, only 0.13% of the items are discarded attributions (448)<sup>65</sup>. This fact may indicate a radical under-representation of complex attributions within Wikidata entities. We are bound to conclude that WLS statements are not particularly widespread nor successful in Wikidata collections within the Cultural Heritage domain, and they arguably misrepresent the complexity and variety of situations in this domain.

**RQ2 – How does the cultural domain of the Wikidata topics (and, presumably, of the individuals contributing to the data regarding the Wikidata topics) affect and reflect on the relative success of some WLS types over others?**

<sup>59</sup><https://www.wikidata.org/wiki/Wikidata:Verifiability>

<sup>60</sup>[http://www.wikidata.org/wiki/Help:Ranking#What\\_ranks\\_are\\_not](http://www.wikidata.org/wiki/Help:Ranking#What_ranks_are_not)

<sup>61</sup>see <https://rkd.nl/en/>

<sup>62</sup><https://w.wiki/7wfw>

<sup>63</sup><https://w.wiki/7VRg>

<sup>64</sup>The count of attributions is calculated over the number of claims having the predicate `P170:creator`

<sup>65</sup>The number of discarded attributions is calculated over the number of claims having `P170:creator` as predicate and not being asserted

	RKD images 17-20th c. Dutch paintings	Wikidata 17-20th c. paintings
Paintings	267,238	501,049
Attributions	317,165	340,661
Current attributions	289,918	340,213
Discarded attributions	27,247	448
% Discarded	8.5%	0.13%

Table 2

Comparison between attributions in RKD images collection, CHv dataset and CHv selection of paintings from 17th to 20th century

– Our data analysis highlighted several peculiarities between the Cultural Heritage and Astronomical datasets. The two families of datasets present many different representational artefacts: while the CH datasets seem to employ, with variable proportions, all the listed approaches, the astronomical datasets employ almost exclusively ranked statements. Additionally, while WLS statements in AN datasets affect a fairly small number of properties, they cover a much wider range of properties in CH, as shown in Figure 2. These aspects highlight key differences in what the two communities consider weaker logical status: we may hypothesise that deprecations in astronomical data mostly reflect the result of newer and better data. In contrast, the humanities community uses WLS statements for a much larger set of uncertainties due to ignorance, scholarly interpretations and disagreements as hypothesised in Section 1. Thus, it may occur that the specification of the `P5102:nature of statement` and the `P2241:reason for deprecated rank` qualifiers may seem overkill in astronomical data, and a real necessity for some annotations in the humanities.

**RQ3 – Does the actual usage of the surveyed approaches match their designed use declared by Wikidata?** – Wikidata provides a set of designed uses for WLS claims annotation as described in Section 3. In addition to them, Wikidata contributors have, over time, adopted frequent annotation patterns that are only sometimes aligned with designed uses. Thus, there is much noise and ambiguity in how Wikidata contributors have used approaches provided by Wikidata to represent WLS information in the datasets we studied. This makes it difficult to differentiate and search WLS data. The variety of cases listed at the end of Section 4.2 summarises an incomplete yet vast collection of WLS and non-WLS situations modelled through the same WLS representation approaches. Therefore, it is difficult to search for specific data patterns over the entire dataset and even to interpret individual entities correctly. In particular, such ambiguities can be specifically listed for the surveyed approaches: (1) Ranked statements are used for both representing WLS information as the evolution of opinions in critical debate (evolving knowledge), historical information (evolving situation) and non-related WLS information such as, e.g., less preferred variant. Additionally, despite the different designed uses for preferred and deprecated statements, in practice, they frequently co-occur in the CH dataset for the same properties, showing that annotators arbitrarily choose between these two approaches to represent such information (e.g. discarded attributions are sometimes represented with a non-asserted normal rank and sometimes with a deprecated rank). (2) The selection of terms provided with *nature of statement* and *sourcing circumstance*, despite being a very expressive pattern to represent WLS information but also its justification, is not exclusively related to WLS information, so that a subset of terms should be defined for this specific purpose (cf. 101 selected terms in Section 4.2). Additionally, no taxonomy is provided on types of WLS qualifiers. For this reason, automatic extraction of types of uncertainty (such as uncertainty, cautioning, and imprecision as discussed in Section 4.2) cannot be automatically performed. (3) Despite the designed use provided by Wikidata, the two types of missing values statements (noValue and someValue) present a significant co-occurrence within the same properties, indicating an unclear usage similar to the usage of ranked statements.

Furthermore, using the same approaches for WLS and non-WLS-related characterisations makes complex patterns hard to express and identify. For instance, if an artwork AW was *supposedly* moved from location X to location Y, but we are not certain, both locations X and Y must be represented as WLS, the first because of an evolving situation (AW is not at location X anymore) and the second because of uncertainty since the new location Y is only guessed. Therefore, none of these assertions can be asserted, and none can be ranked as preferred. We need a complete and thorough contextual annotation (e.g., why each claim is discarded), without which disambiguation and

full understanding of the state and truth of the relevant predicate is impossible. In Section 5, we suggest a possible pattern to represent such situation (cf. point 5, in particular, *normal rank + non-asserted*).

## 5. Towards a leaner and harmonic support for WLS in Wikidata

Getting down to detailing workable solutions to improve the situation for WLS statements in a project as large and as complex as Wikidata is always running the risk of becoming an exercise in futility. In this section, we respectfully suggest possible actions for WLS statements, starting from very conservative proposals with limited impact to more impacting changes.

We list possible remediation activities for the Wikidata data model and the collection to simplify and disambiguate WLS assertions from the rest. We approach such a complex endeavour with humility and caution, as it may be hard to assess the impact and difficulty of implementing each suggested step from our vantage point.

For this reason, we express our suggestions as an ordered list whose first items are meant as simple cleaning-up activities of little impact and then progress to bolder and more impacting actions that sometimes require not just a modification in the data model but possibly also the systematic update of small, but still numerically relevant, selections of the current datasets.

1. Require a `P7452:reason for preferred statement` qualifier in all preferred statements and a `P2241:reason for deprecated statement` qualifier in all deprecated statements. Provide simple-to-use interface widgets for their specification. Sure, such statements can only be saved with a qualifying proposition.
2. Require the specification of `P5102:nature of statement` and `P1480:sourcing circumstances` qualifiers for all WLS-related rankings: only asserted statements with normal rank are allowed to remain without qualifiers.
3. Create a new and separate *Certainty Degree* qualifier specifically for WLS statements, separating the reason for the chosen qualification from the certainty or confidence degree of the qualification. Such certainty degree should be scalar and use a limited number of values, avoiding any complexity in distinguishing between terms such as possibly, hypothetical, and dubious. A 5- or 7-item scale would suffice, e.g., *non accepted*, *highly unlikely*, *unlikely*, *possible*, *probable*, *almost surely*, and *accepted*. Different labels would be perfectly acceptable, even using numerical values instead of labels.
4. Reorganise the values of `P5102:nature of statement` and `P1480:sourcing circumstances` to remove values merely representing an uncertainty (replaced by the new *Certainty Degree* qualifier). To this end, an initial list of values is being created. The current list has been generated by following a Grounded Theory approach [45]: first, labels, definitions and usage data of suggested and used qualifiers have been collected and categorised to represent different macro-themes or concepts. These concepts allowed theories to emerge and be developed from the coded data with an iterative process that continued until the theory was “grounded” in the data. The resulting list in its current state, collecting the surveyed terms from the Wikidata *Property Talk* pages and the terms used in the CH datasets, contains 150 values referring to WLS claims and organised in 18 theories and can be accessed in the GitHub folder of the project<sup>66</sup>.
5. Restrict ranking for competing statements to just three (possibly four) different patterns and prevent any other variant:
  - *Preferred + Deprecated*: To be used whenever there are several competing statements, and some are chosen to be the best. Accepted statements are set to preferred (and asserted), while the rest are set to deprecated (and not asserted); there are no normal ranks. Both preferred and deprecated statements are fully qualified with `P5102:sourcing circumstances`, `P2241:reason for deprecated statement` and `P7452:reason for preferred statement` respectively, and the new *Certainty* qualifier. Preferred statements would be assigned an *accepted* or *almost surely* degree, while deprecated ones would be assigned a *not accepted* or *highly unlikely* certainty degree. Intermediate degrees would not be used.

<sup>66</sup>[https://github.com/alessiodipasquale/Wikidata\\_WLS](https://github.com/alessiodipasquale/Wikidata_WLS)

- 1 – *Normal rank + asserted*: This would be the default situation, to be used when no dispute or disagreement 1  
2 exists and the statement(s) are all equally accepted. All statements are also asserted. Since this is the default, 2  
3 no qualifier is necessary, but it is still possible to specify a `P5102:nature of statement` or a `P1480:sourcing` 3  
4 `circumstance` value. No certainty degree is necessary. 4
- 5 – *Normal rank + non-asserted*: To be used when there are several competing statements but none of them 5  
6 stands above the rest as being the most likely. For instance, this would be the case of a work of art not defi- 6  
7 nitely attributed to anyone but for which several competing hypotheses exist. However, none seem more con- 7  
8 vincing than the others. No statement is asserted, and `P5102:nature of statement` and/or a `P1480:sourcing` 8  
9 `circumstance` values are required. All statements would be assigned a value from the central ones, from 9  
10 *highly unlikely* to *probable*, excluding the extremes. 10  
11

12 A fourth pattern could be allowed for claims for which the only reported value is wrong, but no acceptable 12  
13 alternatives exist. In this case, we could use a deprecated statement for the reported wrong value and a non- 13  
14 existing valued statement with a normal rank to represent the non-existing correct value. 14  
15

## 16 6. Conclusions and future works 16

17  
18  
19  
20 Our work is the first systematic study about the representation of weaker logical status claims (WLS) over Cultural 20  
21 Heritage data in Wikidata. Through WLS claims, uncertain information, competing hypotheses, temporally evolving 21  
22 information, etc., for which a plain and direct assertion is inappropriate, can be expressed. We analysed four patterns 22  
23 used in Wikidata for WLS claims: asserted vs. non-asserted statements, ranked statements, missing values, and 23  
24 qualifiers. 24

25 In our analysis, we found several interesting facts. First, the number of statements expressed using a lower logical 25  
26 status is much lower than might have been expected by comparing similar sources. Secondly, the Wikidata data 26  
27 model is far from being too poor to express WLS claims; it offers users an overabundance of approaches, but 27  
28 their applications overlap and are also used for non-WLS applications. Finally, significant differences exist in how 28  
29 datasets from different domains employ these approaches for weaker logical status claims. Domain-specific non- 29  
30 WLS situations can be considered as a justification for much of this variety, and this contributed to the idea that 30  
31 WLS-specific features should be introduced in the Wikidata model to address specifically weaker logical status 31  
32 claims. We proposed a set of increasingly impacting modifications to the data model aiming towards a leaner and 32  
33 more accurate representation of these phenomena, expecting that they can improve data quality and information 33  
34 retrieval, specifically over uncertain, evolving and competing statements. 34  
35

36 We are still working toward a complete taxonomy of values for qualifying ranked predicates, as this seems to 36  
37 be, to our eyes, the most rapid and solid way to fully represent both the weaker logical status of a claim and its 37  
38 underlying nature and justification. We plan to publish this taxonomy with a proposal for mapping existing data 38  
39 points to this taxonomy so that no information is lost during conversion. 39  
40

## 41 42 Responsibility statement 42

43  
44  
45 Fabio Vitali and Valentina Pasqual jointly wrote the manuscript's introduction (Section 1). Valentina Pasqual 45  
46 authored Section 2, covering the state-of-the-art, and Section 3, focused on approaches to representing WLS in 46  
47 Wikidata. Valentina Pasqual collaborated with Alessio di Pasquale on the data analysis section (Section 4). Fabio 47  
48 Vitali is responsible for section 5, proposing new approaches to represent WLS in Wikidata. All authors contributed 48  
49 to the conclusions section (Section 6). Fabio Vitali and Francesca Tomasi provided critical revisions and feedback 49  
50 throughout the writing process, ensuring the coherence and accuracy of the manuscript. All authors actively partic- 50  
51 ipated in the manuscript review, providing intellectual contributions and final approval for the submission. 51

## References

- [1] F. Erxleben, M. Günther, M. Krötzsch, J. Mendez and D. Vrandečić, Introducing Wikidata to the Linked Data Web, in: *The Semantic Web - ISWC 2014*, P. Mika, T. Tudorache, A. Bernstein, C. Welty, C. Knoblock, D. Vrandečić, P. Groth, N. Noy, K. Janowicz and C. Goble, eds, Springer International Publishing, Cham, 2014, pp. 50–65. ISBN 978-3-319-11964-9. doi:doi.org/10.1007/978-3-319-11964-9\_4.
- [2] C. Möller, J. Lehmann and R. Usbeck, Survey on English Entity Linking on Wikidata: Datasets and Approaches, *Semantic Web* **13** (2022). doi:10.3233/SW-212865.
- [3] D. Vrandečić and M. Krötzsch, Wikidata: a free collaborative knowledgebase, *Communications of the ACM* **57**(10) (2014), 78–85. doi:10.1145/2629489.
- [4] M. Doerr, S. Gradmann, S. Henniecke, A. Isaac, C. Meghini and H. Van de Sompel, The europeana data model (EDM), in: *World Library and Information Congress: 76th IFLA general conference and assembly*, Vol. 10, 2010, p. 15.
- [5] M. Doerr, C.-E. Ore and S. Stead, The CIDOC conceptual reference model-A new standard for knowledge sharing, in: *26th international conference on conceptual modeling (ER 2007)*, 2007.
- [6] M. Piotrowski and M. Neuwirth, Prospects for computational hermeneutics, in: *Proceedings of the 9th AIUCD Annual Conference*, 2020. <http://amsacta.unibo.it/6316/>.
- [7] M. Fafinski and M. Piotrowski, Modelling Medieval Vagueness, in: *INFORMATIK 2020*, R.H. Reussner, A. Koziolok and R. Heinrich, eds, Gesellschaft für Informatik, Bonn, 2021, pp. 1317–1326. doi:10.18420/inf2020\_123.
- [8] M. Daquino, V. Pasqual and F. Tomasi, Knowledge Representation of digital Hermeneutics of archival and literary Sources, *JLIS.it* (2020), 59–76. doi:https://doi.org/10.4403/jlis.it-12642.
- [9] C.L. Borgman and M.F. Wofford, From Data Processes to Data Products: Knowledge Infrastructures in Astronomy, *Harvard Data Science Review* **3**(3) (2021), <https://hdsr.mitpress.mit.edu/pub/xfgywa6x>.
- [10] A. Blau, Uncertainty and the History of Ideas, *History and Theory* **50**(3) (2011), 358–372. doi:https://doi.org/10.1111/j.1468-2303.2011.00590.x. <https://onlinelibrary.wiley.com/doi/abs/10.1111/j.1468-2303.2011.00590.x>.
- [11] H.-G. Gadamer, *Truth and method*, A&C Black, 2013.
- [12] P.T. Darch and A.E. Sands, Uncertainty about the Long-Term: Digital Libraries, Astronomy Data, and Open Source Software, in: *2017 ACM/IEEE Joint Conference on Digital Libraries (JCDL)*, 2017, pp. 1–4. doi:10.1109/JCDL.2017.7991584.
- [13] S. Auer, C. Bizer, G. Kobilarov, J. Lehmann, R. Cyganiak and Z. Ives, Dbpedia: A nucleus for a web of open data, in: *The Semantic Web: 6th International Semantic Web Conference, 2nd Asian Semantic Web Conference, ISWC 2007+ ASWC 2007, Busan, Korea, November 11-15, 2007. Proceedings*, Springer, 2007, pp. 722–735. doi:doi.org/10.1007/978-3-540-76298-0\_52.
- [14] T. Rebele, F. Suchanek, J. Hoffart, J. Biega, E. Kuzey and G. Weikum, YAGO: A multilingual knowledge base from wikipedia, wordnet, and geonames, in: *The Semantic Web-ISWC 2016: 15th International Semantic Web Conference, Kobe, Japan, October 17–21, 2016, Proceedings, Part II 15*, Springer, 2016, pp. 177–185. doi:doi.org/10.1007/978-3-319-46547-0\_19.
- [15] V. Petras, T. Hill, J. Stiller and M. Gäde, Europeana—a Search Engine for Digitised Cultural Heritage Material, *Datenbank-Spektrum* **17** (2017), 41–46. doi:10.1007/s13222-016-0238-1.
- [16] E. Delmas-Glass and R. Sanderson, Fostering a community of PHAROS scholars through the adoption of open standards, *Art Libraries Journal* **45**(1) (2020), 19–23–. doi:10.1017/alj.2019.32.
- [17] E.E. Fink, American Art Collaborative (AAC) Linked Open Data (LOD) Initiative, Overview and Recommendations for Good Practices (2018). <https://repository.si.edu/bitstream/handle/10088/106410/OverviewandRecommendationsAccessible.pdf>.
- [18] A. Isaac et al., Europeana data model primer (2013). <https://pro.europeana.eu/page/edm-documentation>.
- [19] G. Barabucci, F. Tomasi and F. Vitali, Supporting complexity and conjectures in cultural heritage descriptions, in: *Proceedings of the International Conference Collect and Connect: Archives and Collections in a Digital Age*, CEUR Workshop, 2021, pp. 104–115. <http://ceur-ws.org/Vol-2810/paper9.pdf>.
- [20] A. Stinson, S. Fauconnier and L. Wyatt, Stepping Beyond Libraries: The Changing Orientation in Global GLAM-Wiki, *JLIS.it* **9**(3) (2018), 16–34–. doi:10.4403/jlis.it-12480. <https://www.jlis.it/index.php/jlis/article/view/95>.
- [21] M. Zhitomirsky-Geffet and S. Minster, Cultural information bubbles: A new approach for automatic ethical evaluation of digital artwork collections based on Wikidata, *Digital Scholarship in the Humanities* (2022). doi:10.1093/llc/fqac076.
- [22] M. Mora-Cantalops, S. Sánchez-Alonso and E. García-Barriocanal, A systematic literature review on Wikidata, *Data Technologies and Applications* **53**(3) (2019), 250–268. doi:doi.org/10.1108/DTA-12-2018-0110.
- [23] D. Hernández, A. Hogan and M. Krötzsch, Reifying RDF: What works well with wikidata?, *SSWS@ ISWC* **1457** (2015), 32–47.
- [24] S. Klarman and V. Gutiérrez-Basulto, Two-Dimensional Description Logics for Context-Based Semantic Interoperability, *Proceedings of the AAAI Conference on Artificial Intelligence* **25**(1) (2011), 215–220. doi:10.1609/aaai.v25i1.7854. <https://ojs.aaai.org/index.php/AAAI/article/view/7854>.
- [25] P. Bouquet, F. Giunchiglia, F. van Harmelen, L. Serafini and H. Stuckenschmidt, C-OWL: Contextualizing Ontologies, in: *The Semantic Web - ISWC 2003*, D. Fensel, K. Sycara and J. Mylopoulos, eds, Springer Berlin Heidelberg, Berlin, Heidelberg, 2003, pp. 164–179. ISBN 978-3-540-39718-2.
- [26] A. Zimmermann, N. Lopes, A. Polleres and U. Straccia, A general framework for representing, reasoning and querying with annotated semantic web data, *Journal of Web Semantics* **11** (2012), 72–95.
- [27] G. Flouris, I. Fundulaki, P. Pedititis, Y. Theoharis and V. Christophides, Coloring RDF Triples to Capture Provenance, in: *The Semantic Web - ISWC 2009, 8th International Semantic Web Conference, ISWC 2009, Chantilly, VA, USA, October 25-29, 2009. Proceedings*, A. Bernstein, D.R. Karger, T. Heath, L. Feigenbaum, D. Maynard, E. Motta and K. Thirunaryan, eds, Lecture Notes in Computer Science, Vol. 5823, Springer, 2009, pp. 196–212. doi:10.1007/978-3-642-04930-9\_13.

- [28] J.M. Giménez-García, A. Zimmermann and P. Maret, NdFluents: An Ontology for Annotated Statements with Inference Preservation, in: *The Semantic Web*, Springer International Publishing, 2017, pp. 638–654. doi:10.1007/978-3-319-58068-5\_39.
- [29] R. Dividino, S. Sizov, S. Staab and B. Schueler, Querying for provenance, trust, uncertainty and other meta knowledge in RDF, *Journal of Web Semantics* 7(3) (2009), 204–219.
- [30] A. Piscopo and E. Simperl, What We Talk about When We Talk about Wikidata Quality: A Literature Survey, in: *Proceedings of the 15th International Symposium on Open Collaboration*, OpenSym '19, Association for Computing Machinery, New York, NY, USA, 2019. ISBN 9781450363198. doi:10.1145/3306446.3340822.
- [31] M. Färber, F. Bartscherer, C. Menne, A. Rettinger, A. Zaveri, D. Kontokostas, S. Hellmann and J. Umbrich, Linked Data Quality of DBpedia, Freebase, OpenCyc, Wikidata, and YAGO, *Semantic Web* 9(1) (2018), 77–129–. doi:10.3233/SW-170275.
- [32] V. Balaraman, S. Razniewski and W. Nutt, ReCoin: Relative Completeness in Wikidata, in: *Companion Proceedings of the The Web Conference 2018*, WWW '18, International World Wide Web Conferences Steering Committee, Republic and Canton of Geneva, CHE, 2018, pp. 1787–1792–. ISBN 9781450356404. doi:10.1145/3184558.3191641.
- [33] M. Ponzá, P. Ferragina and S. Chakrabarti, A Two-Stage Framework for Computing Entity Relatedness in Wikipedia, in: *Proceedings of the 2017 ACM on Conference on Information and Knowledge Management*, CIKM '17, Association for Computing Machinery, New York, NY, USA, 2017, pp. 1867–1876–. ISBN 9781450349185. doi:10.1145/3132847.3132890.
- [34] L. Galárraga, S. Razniewski, A. Amarilli and F.M. Suchanek, Predicting Completeness in Knowledge Bases, in: *Proceedings of the Tenth ACM International Conference on Web Search and Data Mining*, WSDM '17, Association for Computing Machinery, New York, NY, USA, 2017, pp. 375–383–. ISBN 9781450346757. doi:10.1145/3018661.3018739.
- [35] K. Shenoy, F. Ilievski, D. Garijo, D. Schwabe and P. Szekely, A study of the quality of Wikidata, *Journal of Web Semantics* 72 (2022), 100679. doi:https://doi.org/10.1016/j.websem.2021.100679. https://www.sciencedirect.com/science/article/pii/S1570826821000536.
- [36] D. Abián, F. Guerra, J. Martínez-Romanos and R. Trillo-Lado, Wikidata and DBpedia: A Comparative Study, in: *Semantic Keyword-Based Search on Structured Data Sources*, J. Szymański and Y. Velegrakis, eds, Springer International Publishing, Cham, 2018, pp. 142–154. ISBN 978-3-319-74497-1.
- [37] H. Arnaout, S. Razniewski, G. Weikum and J.Z. Pan, Negative Knowledge for Open-world Wikidata, in: *Companion Proceedings of the Web Conference 2021*, WWW '21, Association for Computing Machinery, New York, NY, USA, 2021, pp. 544–551. ISBN 978-1-4503-8313-4. doi:10.1145/3442442.3452339.
- [38] Help:Ranking - Wikidata. <https://www.wikidata.org/wiki/Help:Ranking>.
- [39] P.F. Patel-Schneider, Contextualization via qualifiers., in: *CKG SemStats@ ISWC*, 2018.
- [40] S. Aljalbout, G. Falquet and D. Buchs, Handling Wikidata Qualifiers in Reasoning, *arXiv preprint arXiv:2304.03375* (2023).
- [41] D. Hernández, C. Gutierrez and A. Hogan, Certain Answers for SPARQL with Blank Nodes, in: *The Semantic Web – ISWC 2018: 17th International Semantic Web Conference, Monterey, CA, USA, October 8–12, 2018, Proceedings, Part I*, Springer-Verlag, Berlin, Heidelberg, 2018, pp. 337–353–. ISBN 978-3-030-00670-9. doi:10.1007/978-3-030-00671-6\_20.
- [42] D.V. McDermott and D. Dou, Representing Disjunction and Quantifiers in RDF, in: *Proceedings of the First International Semantic Web Conference on The Semantic Web*, ISWC '02, Springer-Verlag, Berlin, Heidelberg, 2002, pp. 250–263–. ISBN 3540437606.
- [43] J. Garson, Modal Logic, in: *The Stanford Encyclopedia of Philosophy*, Spring 2023 edn, E.N. Zalta and U. Nodelman, eds, Metaphysics Research Lab, Stanford University, 2023.
- [44] A.D. Pasquale, F. Vitali and V. Pasqual, Wikidata selection of Cultural Heritage, Stars, Galaxies and Random entities and claims, Zenodo, 2023. doi:10.5281/zenodo.7624783.
- [45] B.G. Glaser and A.L. Strauss, *The discovery of grounded theory: Strategies for qualitative research*, Routledge, 2017. ISBN 1351522167.
- [46] A. Blau, Uncertainty and the history of ideas, *History and Theory* 50(3) (2011), 358–372. <http://www.jstor.org/stable/41300100>.