





Wikidated 1.0: An Evolving Knowledge Graph Dataset of Wikidata’s Revision History

Lukas Schmelzeisen¹  , Corina Dima¹ , and Steffen Staab^{1,2} 

¹ University of Stuttgart, Germany

² University of Southampton, United Kingdom

{Lukas.Schmelzeisen,Corina.Dima,Steffen.Staab}@ipvs.uni-stuttgart.de

Abstract. Wikidata is the largest general-interest knowledge base that is openly available. It is collaboratively edited by thousands of volunteer editors and has thus evolved considerably since its inception in 2012. In this paper, we present Wikidated 1.0¹, a dataset of Wikidata’s full revision history, which encodes changes between Wikidata revisions as sets of deletions and additions of RDF triples. To the best of our knowledge, it constitutes the first large dataset of an evolving knowledge graph, a recently emerging research subject in the Semantic Web community. We introduce the methodology for generating Wikidated 1.0 from dumps of Wikidata, discuss its implementation and limitations, and present statistical characteristics of the dataset.

Keywords: Semantic Web · Wikidata · Edit History · Knowledge Graph Change · Knowledge Graph Evolution · Stream of RDF Triple Changes

1 Introduction

A *knowledge graph* is “a graph of data intended to accumulate and convey knowledge of the real world, whose nodes represent entities of interest and whose edges represent potentially different relations between these entities” [11]. Recently, knowledge graphs have received much attention in research and powered many diverse applications, such as web search [30], recommendations [9,24], question answering [12], and more [15,17].

Most research so far treats knowledge graphs as static in the sense that change over time is not modeled explicitly. However, in practice, knowledge graphs change over time: new knowledge may be added to the graph (in the form of new edges being added to existing entities, new entities emerging over time, etc.), and existing knowledge may be revised (in the form of imprecise labels being updated, duplicate entities merged, existing contents declared out-of-scope, etc.). This more general setting is only starting to be considered recently [14,21,25].

Copyright © 2021 for this paper by its authors. Use permitted under a Creative Commons Attribution 4.0 International License (CC BY 4.0).

¹ Dataset and code base available at <https://w3id.org/wikidated>.

For example, while there are a plethora of different approaches for knowledge graph embedding [36,7,31,2]—the task of representing entities and relations in a low-dimensional vector space—only a handful of these consider the situation where the representation needs to be updated due to a change in the underlying knowledge graph [4,37,6]. To highlight this distinction, we use the term *evolving knowledge graph* to refer to knowledge graphs that change over time, and the term *static knowledge graph* to refer to those that do not.

Currently, there are practically no published datasets recording a knowledge graph’s organic evolution over time that would enable such research and facilitate reproducible evaluation environments. Instead, existing research either simulates knowledge graph evolution on top of datasets of static knowledge graphs using simple heuristics such as statement ordering (e.g., Daruna et al. [6] split static knowledge graph datasets into chunks based on the order of triples in the dataset) or are based on the computed changes between major releases of knowledge graphs (e.g., Wu et al. [37] calculate the change sets between YAGO2.5, YAGO3, and YAGO3.1 [26]). While the former approaches can hardly be argued to constitute evolving knowledge graphs—in particular statement updates or deletion are not modeled—the latter ones fail to capture the inherent dynamics of how changes occur on the individual level, for example that popular entities receive frequent updates or that some updates might be (partially) reversed after a few days.

Wikidata [35] is “a collaboratively edited knowledge-base [...] whose aim is to curate and represent the factual information of Wikipedia (across all languages) in an interoperable, machine-readable format” [10]. With 90 million entities and 1.4 billion revision made by 20 thousand active users², Wikidata is the prime example of an evolving knowledge graph. In this paper, we present Wikidated 1.0, an evolving knowledge graph dataset covering the full revision history of Wikidata. To the best of our knowledge, Wikidated 1.0 is the first large dataset of an evolving knowledge graph. It records the fine-grained, organic evolution of Wikidata since its inception in 2012 until June 2021. It is suited for research into how knowledge graphs and their communities change over time—specifically for Wikidata, such as done in Sarasua et al. [27]—and enables reproducible evaluation environments for indexing and representation approaches of evolving knowledge graphs, such as incremental knowledge graph embedding [37,6].

In particular, our contributions are:

- We present our methodology for transforming a dump of Wikidata’s revision history into streams of RDF triple deletions and additions which form Wikidated 1.0, a dataset recording Wikidata’s evolution over time (Sect. 3).
- We present statistics over the dataset and visualize its characteristics (Sect. 4).
- We publicly release Wikidated 1.0 in two variants, the code base used to built it, and a Python API to access it¹.

In addition to the above, Sects. 2 and 5 discuss background and related work, respectively, and Sect. 6 concludes.

² Statistics from <https://www.wikidata.org/wiki/Wikidata:Statistics> (23 July 2021).

2 Background

In this section, we review the data model of Wikidata (Sect. 2.1) and its serialization as RDF (Sect. 2.2).

2.1 Data Model

Formally, the data model of Wikidata³ can be defined as a set of *entities* e_1, \dots, e_N , where N is the number of entities in Wikidata. Let $\text{id}(e_i)$ denote the *entity ID* of entity e_i . Each entity is either an *item* or a *property*⁴. Items are things or concepts in the real world about which facts should be stored; their IDs are numbers prefixed with “Q”—for example, the English writer **Douglas Adams** (Q42) or the concept of a **human** (Q5). Properties are abstract types of statements which are used to store facts about entities; their IDs are numbers prefixed with “P”—for example, the properties **instance-of** (P31) or **date-of-birth** (P569).

A *revision* defines an entity’s state at a specific point in time. Each revision is comprised of: (1) a *fingerprint*, which consists of multilingual sets of *labels*, *descriptions*, and *aliases* of the entity, (2) a set of *site links*, which are usually links to Wikipedia articles about the entity, and (3) a set of *statements*, which are records of facts about the entity. There is one exception: if an entity is found to be a duplicate of another one, a revision can also be a *redirect*. In the case of redirects, no fingerprint, site links, or statements are present, and the revision consists of just the entity ID that the redirect’s entity is deemed to be a duplicate of. Further, each revision (including redirects) carries metadata, such as the time it was created at, the contributor that authored it, and a comment string about the change.

Every time an entity is modified, a new revision is created. We therefore model entities as sequences of their revisions $e_i = (r_{i,1}, \dots, r_{i,n_i})$, where n_i is the number of revisions of entity e_i . Let $\text{id}(r_{i,j})$ denote the *revision ID* of revision $r_{i,j}$. Revision IDs are assigned by incrementing a global counter. They are thereby unique over all entities and induce a total ordering of all revisions in Wikidata⁵.

Finally, statements record facts about entities and fundamentally consist of a property and a *value*, which is either another entity or a literal. For example, statements about **Douglas Adams** include **instance-of human** and **date-of-birth “11 March 1952”**. As is the case in the latter example, literals can be of various data types, e. g., dates or geographical coordinates. An example of a statement about a property is that the **complementary-property** of **date-of-birth** is **date-of-death**. Note, that entities whose latest revision is a redirect and deleted entities can still be targets of the statements of other entities, but that Wikidata aims to replace instances of this with new revisions where this is not the case. There are

³ For more details, see <https://www.mediawiki.org/wiki/Wikibase/DataModel>.

⁴ As discussed in Sect. 3.3, technically, entities can also be *lexemes*, *forms*, or *senses*.

⁵ However, the special case, in which a revision has an earlier timestamp than one with a lower ID, can occur. We attribute this to slight miss-synchronizations of clocks on parallel servers. The difference is never larger than one second.

two special values of statements: **none**, which signifies that it is known that the property of that entity has no value, and **some**, which indicates that it is known that there is some value for the property of that entity, but it is unknown what it is. Each statement can be annotated by (1) a set of *qualifiers*, which refine a statement (e. g., to indicate that it has only been true for some period of time), (2) a set of *references*, which provide sources to support the statement, and (3) a *rank*, which can be used to assign preference to conflicting statements (e. g., to distinguish current from historical facts).

Dumps of Wikidata’s contents are available for download in various formats⁶. Most dump formats only provide the most recent revision of each entity; only the **pages-meta-history** XML dumps include the full revision history, but store revision contents as JSON blobs⁷. Notably, each revision stores its complete state at its creation time and there is no trivial way to identify what changed from one revision to the next.

In cases of vandalism, or when entities do not meet Wikidata’s notability policy⁸, administrators may delete the affected revisions or whole entities from Wikidata⁹. Since such deleted contents may contain copyrighted materials or sensible personal information, deleted entities and revisions are not accessible to the general public and are not contained in the official dumps of Wikidata. IDs of deleted entities and revisions are never reused for new ones, which leads to gaps in the incremental numbering. In case an entity is detected that is a duplicate of an existing one, Wikidata prefers not to delete the new entity, but to establish a redirect from the new entity to the existing one instead.

2.2 RDF

The *Resource Description Framework (RDF)* [28,5] is a metadata format and the standard way of exchanging information on the Semantic Web. For the Wikidated 1.0 dataset, we serialize Wikidata revisions as RDF graphs which allows for a straightforward definition of change between revisions.

Let I , B , and L be disjoint countably infinite sets of IRIs, blank nodes, and literals, respectively. A *RDF triple* is a triple $(s, p, o) \in (I \cup B) \times I \times (I \cup B \cup L)$, where s is called the subject, p the predicate, and o the object. A *RDF graph* is a set of RDF triples. Let $G_1 \setminus G_2$ denote the *set difference* between two RDF graphs G_1 and G_2 . Computing it is non-trivial, because in order to decide whether two RDF triples are equal, one needs to decide which triple components are equal to one another. This is straightforward for IRIs and literals, but hard for blank nodes, as they are only characterized through the RDF triples they participate in and do not have identifiers across RDF graphs. In general, finding a mapping between the blank nodes of two RDF graphs that minimizes the set difference is NP-hard [34,18,1,13]. For Wikidated 1.0, we circumvent this issue (see Sect. 3.2).

⁶ See https://www.wikidata.org/wiki/Wikidata:Database_download.

⁷ See https://doc.wikimedia.org/Wikibase/master/php/md_docs_topics.json.html.

⁸ Available at <https://www.wikidata.org/wiki/Wikidata:Notability>.

⁹ Requests for deletions and the decisions for each are recorded at https://www.wikidata.org/wiki/Wikidata:Requests_for_deletions.

```

1 wd:Q42 a wikibase:Item ;
2   rdfs:label "Douglas Adams"@en ;
3   schema:description "English writer and humorist"@en ;
4   skos:altLabel "Douglas Noel Adams"@en ;
5   wdt:P569 "1952-03-11T00:00:00Z"^^xsd:dateTime ;
6   p:P569 s:Q42-D8404CDA-25E4-4334-AF13-A3290BCD9C0F .
7
8 <https://en.wikipedia.org/wiki/Douglas_Adams> a schema:Article ;
9   schema:about wd:Q42 .
10
11 s:Q42-D8404CDA-25E4-4334-AF13-A3290BCD9C0F a wikibase:Statement ;
12   ps:P569 "1952-03-11T00:00:00Z"^^xsd:dateTime ;
13   prov:wasDerivedFrom ref:355b56329b78db22be549dec34f2570ca61ca056 .
14
15 ref:355b56329b78db22be549dec34f2570ca61ca056 a wikibase:Reference ;
16   pr:P248 wd:Q5375741 .

```

Listing 1. RDF serialization of Wikidata entity Q42 in Turtle syntax (abridged).

While Wikidata doesn’t store revisions in RDF internally, RDF serializations for them are available [8,10,15]. Listing 1 shows an example¹⁰. The shown revision of entity Q42 (Douglas Adams) is described by its fingerprint (Lines 2 to 4), a site link (Lines 8 to 9), a *simple statement* (Line 5), and a *full statement* (Line 6 and Lines 11 to 13). Both statements specify a *date-of-birth* (P569) of “11 March 1952”. The difference is that simple statements give a value “directly” while discarding statement annotations (i. e., qualifiers, references, and ranks), whereas full statements use *reification* (the insertion of a special statement node) to facilitate annotations. In this case, a reference to Encyclopædia Britannica Online (Q5375741, Lines 15 to 16) is used to annotate the statement (Line 13).

3 Constructing Wikidated 1.0 from Wikidata Dumps

In this section, we discuss our methodology for creating Wikidated 1.0 (Sect. 3.1), our implementation (Sect. 3.2), and limitations (Sect. 3.3).

3.1 Methodology

Fundamentally, Wikidated 1.0 is a transformation from a Wikidata dump’s stream of revisions into a stream of *incremental revisions*. We define an incremental revision as a tuple of the (1) entity metadata (the entity ID and some Wikidata internal fields), and the (2) revision metadata (the revision ID and when and by whom it was authored) of the Wikidata revision it is based upon, as well as sets of (3) *RDF triple deletions* and (4) *RDF triple additions* in comparison to the previous revision of the respective entity. As hinted at in Sect. 2, we thus decide

¹⁰ Taken from <https://www.wikidata.org/wiki/Special:EntityData/Q42.ttl>. Documentation at https://www.mediawiki.org/wiki/Wikibase/Indexing/RDF_Dump_Format.

```

1  $\Delta_{\text{Global}} \leftarrow$  empty sequence
2 download pages-meta-history dump of Wikidata
3 foreach entity  $e_i$  in dump do
4    $\Delta_{\text{Entity}} \leftarrow$  empty sequence
5    $r_{\text{RDF-prev}} \leftarrow \{\}$ 
6   for  $j = 1$  to  $n_i$  do
7      $r_{\text{Meta}} \leftarrow$  take revision metadata of  $r_{i,j}$  from dump
8      $r_{\text{JSON}} \leftarrow$  take JSON blob of  $r_{i,j}$  from dump
9
10     $r_{\text{RDF}} \leftarrow$  serialize  $r_{\text{JSON}}$  as RDF graph
11     $r_{\text{Del}} \leftarrow r_{\text{RDF-prev}} \setminus r_{\text{RDF}}$ 
12     $r_{\text{Add}} \leftarrow r_{\text{RDF}} \setminus r_{\text{RDF-prev}}$ 
13     $r_{\text{RDF-prev}} \leftarrow r_{\text{RDF}}$ 
14
15    append incremental revision  $(e_i, r_{\text{Meta}}, r_{\text{Del}}, r_{\text{Add}})$  to  $\Delta_{\text{Entity}}$ 
16
17  output  $\Delta_{\text{Entity}}$  as entity-stream variant of entity  $e_i$ 
18  append all elements of  $\Delta_{\text{Entity}}$  to  $\Delta_{\text{Global}}$ 
19
20 sort  $\Delta_{\text{Global}}$  after ascending revision IDs  $\text{id}(r_{i,j})$  across entities
21 output  $\Delta_{\text{Global}}$  as global-stream variant

```

Algorithm 1. Construction of the Wikidated 1.0 dataset.

to define change between Wikidata revisions via the difference in triples between their RDF serializations. This allows for more straightforward dataset modeling and consumption, as opposed to defining change for each of the different aspects of the Wikidata data model (fingerprint, site links, and statements with qualifiers, references, and ranks).

Wikidated 1.0 consists of two complementary variants of the same data:

1. The *global-stream* variant consists of all incremental revisions across all entities sorted in chronological order.
2. The *entity-streams* variant contains a separate stream of incremental revisions for each entity of Wikidata.

The former can be used for global analysis, e. g., for analyzing the number or style of revisions in a specific time period, whereas the latter is useful for entity-centered analysis, e. g., when one is only interested in a subset of all entities or when the aim is to directly compare consecutive revisions of the same entity.

Algorithm 1 outlines the steps of creating Wikidated 1.0. First, we download a full Wikidata dump¹¹ (Line 2). Next, we iterate over all entities in the dump (Line 3). For each entity, we iterate over all of its revisions in the order they were created in (Line 6). Because the dump files store entities and revisions

¹¹ Specifically, Wikidated 1.0 is based on the 20210601-pages-meta-history dump, the history of Wikidata from its inception on 30 October 2012 until June 2021. The ID of the last revision is 1433475551, which was authored on 2 June 2021 at 05:35:58.

in exactly this order, this amounts to linearly traversing the dump files. For each revision, we first extract revision metadata and the JSON blob of revision contents from the dump files (Lines 7 and 8). We then serialize the revision contents as an RDF graph (Line 9), and compute the sets of RDF triple deletions and additions (Lines 10 and 11) compared to the RDF graph of the previous revision, for which we maintain a helper variable (Lines 5 and 12). Having now transformed the revision into its incremental counterpart, we append it to a sequence of all incremental revisions of the entity Δ_{Entity} (Line 13, initialization in Line 4). After all revisions of an entity have been iterated over, we output Δ_{Entity} as the entity-stream variant of that entity (Line 14). Finally, we maintain a sequence of all incremental revisions across all entities Δ_{Global} (Lines 1 and 15). After sorting all revisions in it globally (Line 16), we also output Δ_{Global} as the global-stream variant (Line 17).

3.2 Implementation

We provide a Python API for browsing and iterating the dataset without having to know how the dataset is stored on disk. Internally, both variants are stored as `gzip`-compressed text files in JSON Lines format, i. e., each line is a JSON object encoding one incremental revision. For the entity-streams variant, the files for all entities are packaged in a `tar` archive. While approaches for storing differences between RDF graphs in RDF itself exist [3,20], we hereby opt for a less Semantic-Web-oriented distribution format, because we feel that it allows for easier consumption by most users. We are open to releasing the dataset in RDF later (based on community demand).

The file size of the global-stream variant of Wikidated 1.0 is 239 GiB whereas the `tar` archive for the entity-streams variant is 329 GiB (both `gzip`-compressed). In contrast to this, the official Wikidata dump of non-incremental revisions that Wikidated 1.0 is built from is available in the two compression formats `bz2` and `7z` with a size of 1 040 GiB and 339 GiB, respectively. While Wikidated 1.0’s smallest variant is thus only 71% the size of the official dump’s smallest format, one might have expected an even larger reduction in size due to the usage of incremental revisions that do not repeat all statements from earlier revisions. We suspect that this benefit is offset by using the RDF serialization, which is more verbose than the JSON blobs of the official dumps, that the non-incremental revisions are more compressible through their repeated statements, and by the inferior compression of `gzip` compared to `7z`. We are therefore looking into publishing our dataset in additional compression formats in the future, but have opted for the universally-available and streamable `gzip` format for the first release.

For serializing Wikidata revisions as RDF graphs (Line 9), we use Wikidata Toolkit¹². Because it does not provide a way to serialize Wikidata revisions that are redirects, we encode these via the `owl:sameAs` predicate, following the choice of the Wikidata Query Service¹³ [15]. Additionally, we discard all RDF triples

¹² Available at https://www.mediawiki.org/wiki/Wikidata_Toolkit.

¹³ Available at <https://query.wikidata.org/>.

from Wikidata Toolkit’s output that are not directly related to the entity at hand. The triples discarded in this manner contain ontological information about Wikidata concepts such as items and properties. If needed, these can always be reconstructed from context and they never change between revisions.

For computing set differences between RDF graphs (Lines 10 and 11), we can avoid the difficult task of finding an optimal mapping between blank nodes. In the case of Wikidata’s RDF serialization, blank nodes are only used to encode the special `some` values, and each blank node never occurs in more than one RDF triple. Because of this, there is an efficient way to determine the set difference between RDF serializations of two Wikidata revisions: two RDF triples can be treated as equal if their non-blank-node components are equal.

Last, assembling the incremental revision streams Δ_{Global} and Δ_{Entity} (Lines 15 and 16) is not as straightforward as presented, because both would quickly exceed available main memory. Instead, we directly append any incremental revisions (Line 13) to the target file without keeping them in memory. To merge all Δ_{Entity} streams into a sorted Δ_{Global} , we use a hierarchical multiway merge.

3.3 Limitations

The main limitation of Wikidated 1.0 is that it does not contain any record of deleted entities or revisions¹⁴ because these are not recorded in Wikidata’s revision history dumps, as explained in Sect. 2.1. However, entity deletions are comparatively rare in Wikidata, since the common case of duplicate entities is addressed through merges, i. e., redirects of one entity to another, which—unlike deletions—are recorded in Wikidated 1.0. Additionally, statements of other entities may still target deleted entities, so partial history of them is recorded.

Multiple implementations of RDF serializations of Wikidata revisions exist. The one in Wikidata Toolkit was used to construct Wikidated 1.0. Both the RDF exports of individual Wikidata entities and the official Wikidata RDF dumps use two other slightly different implementations. In practice, the differences are minimal and mostly amount to how certain metadata is encoded—the important parts, i. e., facts about entities, are encoded identically in all implementations.

By design, Wikidated 1.0 only contains the revision history of entities. As a consequence the “meta level” of Wikidata is not part of the dataset. Among other things, this includes the talk pages of all entities, where editors discuss aspects such as how certain content should be modeled or what is in scope for Wikidata, or the help pages, which document how to use Wikidata. While this plain text data is part of the Wikidata dumps it is not RDF serializable.

In May 2018, three new entity types have been added to the Wikidata data model to model lexicographical data: lexemes, forms, and senses¹⁵ [16]. While these are part of the Wikidata dumps and a RDF serialization for them has been defined [16], it has not yet been implemented in Wikidata Toolkit and lexicographical data is thus not part of this first release of our dataset.

¹⁴ In Sect. 5, we review the work of Shenoy et al. [29], which describes an approach that is able to obtain some information about deleted entities (from monthly dumps).

¹⁵ Documentation at https://www.wikidata.org/wiki/Wikidata:Lexicographical_data.

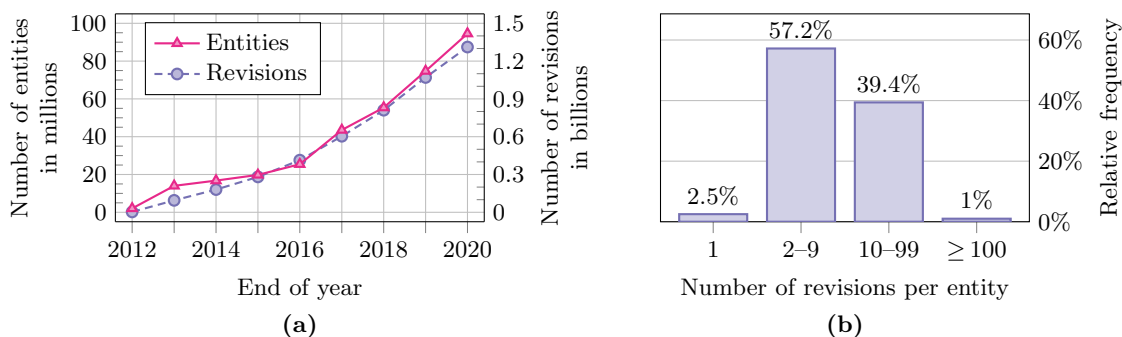


Fig. 1. (a) Number of entities and revisions at the end of each year. The totals by 2 June 2021 (end of dataset’s time range) are 96.6 million entities and 1.4 billion revisions, respectively. (b) Histogram of number of revisions per entity (mean 14.60, standard deviation 24.94, median 7).

4 Wikidated 1.0 Dataset Characteristics

In this section, we present statistical characteristics of the Wikidated 1.0 dataset in order to provide context for any research work building upon it. Much of the analysis also applies to Wikidata given that Wikidated 1.0 is a direct representation of it.

We plot the number of entities and revisions over time in Fig. 1(a)¹⁶. For both, we observe a mostly linear growth, with a slightly stronger incline since 2016, which follows the integration of Freebase into Wikidata in the latter half of 2015 [32]. Interestingly, this means that the ratio between both frequencies is roughly constant at about 14 revisions per entity.

Fig. 1(b) shows the number of revisions per entity in more detail. Roughly half of all entities have fewer than 10 revisions, and the majority have more than one. 99% of entities have less than 100 revisions, whereas some of the remaining entities can have significantly more.

The time between consecutive revisions of the same entity is visualized in Fig. 2(a). 30% of revisions are being authored within less than a minute since the previous revision. We suspect that this stems either from heavy activity on the most popular entities, or from editors performing multiple related changes directly after another, such as changing two related statements or reverting erroneous edits. For 60% of revisions, the time since the previous revision is less than a month. For less than 4% of revisions, that time is more than a year. Coupled with the previous figure’s data of only 2.5% of entities having exactly one revision, it stands to reason that most entities in Wikidata are edited somewhat frequently.

¹⁶ In all figures of this section, we treat revisions that are redirects as regular revisions (being represented by exactly one `owl:sameAs` RDF triple). As a consequence, we also count entities whose latest revision is a redirect as regular entities.

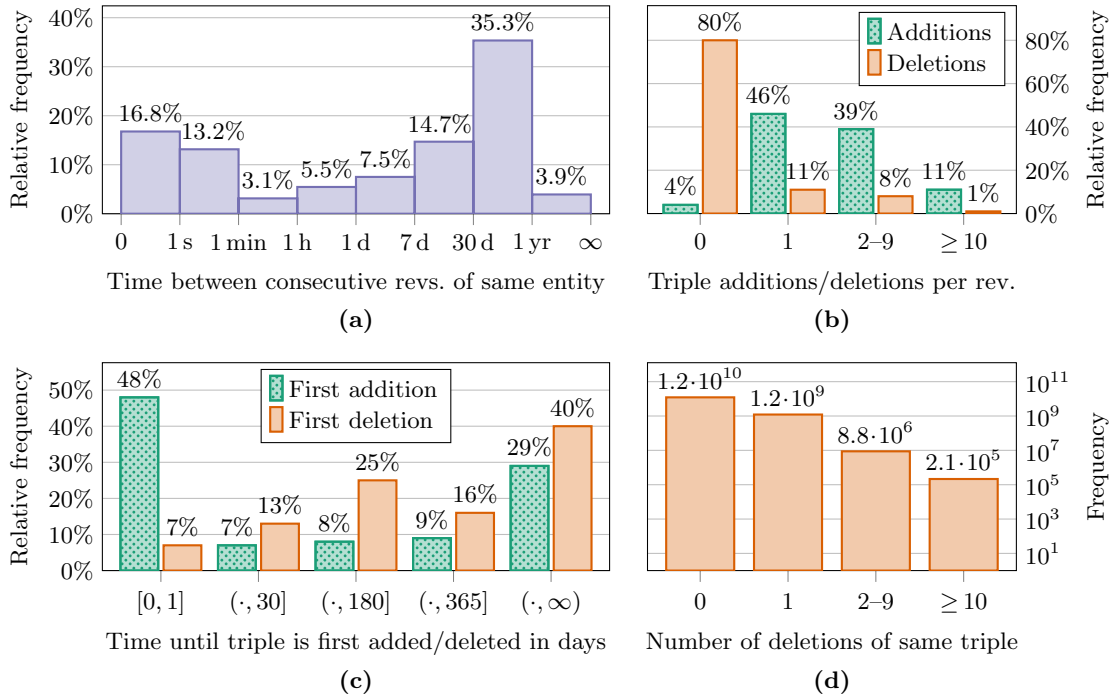


Fig. 2. (a) Histogram of time between consecutive revisions of the same entity (mean 69.98 days, standard deviation 149.39 days, median 12 days). (b) Histogram of number of RDF triple additions per revision (mean 9.56, standard deviation 38.60, median 2) and deletions per revision (mean 0.88, standard deviation 14.73, median 0). (c) Histogram of time until a RDF triple is first added/deleted. For added triples, the time since the creation of the entity is measured (mean 363.64 days, standard deviation 626.19 days, median 8 days). For deleted triples, the time since the triple has been added is measured (mean 396.84 days, standard deviation 450.04 days, median 236 days). (d) Histogram of number of deletions of same RDF triple (mean 0.09, standard deviation 0.36, median 0).

However, this does not imply that most entities are checked by humans with some frequency, as these changes could also have been made by automated bots.

Getting closer to the contents of revisions, we look at the number of additions and deletions of RDF triples per revision in Fig. 2(b). Note that there is no one-to-one correspondence between the number of Wikidata statements and the number of RDF triples. For example, a single site link is expressed in more than one triple (compare Lines 8 to 9 of Listing 1). Additions are much more common than deletions with 80% of revisions not featuring any triple deletions. Since 89% of revisions contain less than 10 triple additions, we conclude that most revisions constitute atomic changes and that the case in which multiple statements of a single entity change is much rarer.

In Fig. 2(c), we show the time until a RDF triple is first added or deleted. Approximately half of triples are added within less than a day since the the

creation of its entity. Deletions take far longer: more than half of all deleted triple are deleted more than half a year after they had originally been added. Besides changes to the Wikidata schema—like the deletion of properties—that potentially entail (semi-) automated changes to otherwise unchanged entities, we see two competing explanations for these late deletions: it might simply take a while until facts in the real world change and Wikidata can only update its record of them once they do, or alternatively, Wikidata might take a while to detect incorrect knowledge for the less popular entities. A more detailed look into classifying the types and causes of changes will therefore be necessary for further investigation.

Last, Fig. 2(d) visualizes repeated deletions of the same RDF triple. Unsurprisingly, the vast majority of triples added to Wikidata are never deleted. Slightly less than 10% of triples are deleted exactly one time; 4% of which are added back into Wikidata again afterwards (not shown in figure). Even though only less than 1% of triples are deleted from Wikidata more than once, a few of these are deleted very many times. For example, around 52 thousand triples are deleted and added to Wikidata more than 100 times (not shown in figure). We suspect heavy edit wars—potentially between bots—as the main cause for this.

To summarize, we have quantified how Wikidata changes over time on a macro level through analyzing statistical characteristics of Wikidated 1.0, which demonstrates its fitness as a dataset for evolving knowledge graph research.

5 Related Work

Based on their naming, Wikidata’s incremental dumps¹⁷ may seem to address the exact same problem as Wikidated 1.0. These dumps are published every 24 hours and contain all revisions authored since the last dump. However, like the full dumps discussed in Sect. 2.1, each revision is stored in its full state and no obvious way exists to identify what changed from one revision to the next. Additionally, incremental dumps older than a few months are routinely taken offline. Because of this, they do not offer a way to trace the full edit history since Wikidata’s inception like Wikidated 1.0 does. Their main use case is to keep live services operating on Wikidata’s contents up-to-date.

Much closer to our setting is the history query service¹⁸ [20]. It consists of a SPARQL [33] endpoint that allows users to query for Wikidata revision differences—similarly to Wikidated 1.0. The paper’s main contributions are on how to express revision additions and deletions in a RDF data model and how to index them for efficient query answering. On the other hand, the paper does not discuss *how* revision additions and deletions are computed, does not discuss any limitations, and does not provide a stable, downloadable dataset. It is therefore not suitable as an environment for reproducible evaluations.

¹⁷ Available at <https://dumps.wikimedia.org/other/incr/wikidatawiki/>.

¹⁸ Available at <https://wdhq.wmflabs.org/>, however only displaying a “502 Bad Gateway” error during the time of writing (June to October 2021). Documentation at https://www.wikidata.org/wiki/Wikidata:History_Query_Service.

In contemporary work, Shenoy et al. [29] follow an alternative approach¹⁹ for studying changes in Wikidata over time. In contrast to our approach of parsing a single dump of Wikidata’s full revision history, they utilize monthly dumps of Wikidata’s current state, i. e., dumps that only contain the most recent revision of each entity for the respective month. They then analyze which statements were deleted and added from one month to the next. In comparison to Wikidated 1.0, which records statement changes at the revision level and thus includes revision metadata such as the exact point in time when a statement was deleted/added, their approach thus only aggregates all changes to an entity per month. This aggregation implies the inability to record phenomena such as the frequency of revisions to entities per month or changes that are reverted within the same month. The upside of their approach is that records of deleted entities, which are purged from the full revision history dump, are still available in those monthly dumps that were created before the entity was deleted. Their data thus provides a useful addition to Wikidated 1.0.

The CorHist dataset [19] is another dataset build from Wikidata’s edit history. However, it limits itself to recording constraint-related data. Wikidata constraints are similar to database integrity constraints and used to aid Wikidata editors in finding erroneous data. Specifically, the CorHist dataset records past constraint violations and their corrections. Wikidated 1.0, in comparison, records all statement changes (including constraint violations, albeit in a different format than CorHist and only implicitly) including revision metadata, such as when and by whom a revision was authored, which makes it a more complete resource.

Other research that studies Wikidata’s evolution includes Sarasua et al. [27], which studies the engagement of Wikidata’s editors over time; Piscopo et al. [22], which evaluates the quality of provenance information in Wikidata; and Piscopo and Simperl [23], which investigates the relation between different types of editors and their impact on the Wikidata ontology over time.

6 Conclusion

We have presented Wikidated 1.0, a dataset containing Wikidata’s revision history as incremental revisions, i. e., sets of deletions and additions of RDF triples. To the best of our knowledge, it constitutes the first large evolving knowledge graph dataset of its kind. We foresee applications both from the Wikidata community for studying how Wikidata changed over time, as well as from the wider knowledge graph community for evaluating techniques over evolving knowledge graphs, such as incremental knowledge graph embeddings or updatable indexing structures for efficient query answering.

Besides releasing the dataset and the accompanying codebase¹, in this paper we have documented our methodology for creating Wikidated 1.0, and discussed its implementation and limitations, the biggest one being the omission of deleted entities, which are not contained in the openly available revision history dumps

¹⁹ Data and analysis scripts available at https://w3id.org/wd_quality.

of Wikidata. Additionally, we have presented statistical characteristics of our dataset, and compared it to related work.

6.1 Future Work

In the future, we plan to release new versions and additional variants of our dataset. In particular, we aim to establish a release cadence for publishing a new Wikidated version on the most recent Wikidata dump in regular time intervals, based on community uptake. Additionally, we are working on extracting subsets of Wikidated 1.0 that only contain RDF serializations of simple statements, i. e., statements with qualifiers and references removed, or “thruthy” statements, i. e., statements with the highest rank—similar to the existing equally-named variant of the official Wikidata dumps. Last, we are thinking about additional ways of reducing the large file size of Wikidated 1.0, such as switching to better compression formats, subsampling the dataset, or aggregating deletions and additions of all revisions in a fixed time frame (e. g., an hour, day, or week).

Further future work includes deeper analysis of the editing dynamics of Wikidata recorded in Wikidated 1.0, such as detecting updates from sets of triple deletions and additions, or classifying the type and source of changes; and consideration of Wikidata’s more recent lexicographic data, the first step towards this would be to implement a RDF serialization of it in Wikidata Toolkit. We invite the Wikimedia Foundation specifically to consider also releasing official Wikidata dumps in an incremental format—such as the one used for Wikidated 1.0—in order to save bandwidth and storage space for users. Additionally, we would welcome any way to access and integrate revision data of deleted entities and revisions into Wikidated. While the raw data itself contains sensible personal information and copyrighted material that is unavailable to the general public for good reason, the metadata of deleted entities and revisions, e. g., number of deleted statements, is by itself interesting to us. For instance, it could be published by replacing all literal values in deleted revisions with generated ones, thus only preserving the data’s graph structure.

Acknowledgments Lukas Schmelzeisen was supported by the German Research Foundation (DFG) via grant agreement number STA 572/18-1 (Open Argument Mining). Corina Dima was supported by the German Federal Ministry for Economic Affairs and Energy (BMWi) via grant agreement number 01MK20008F (Service-Meister). Thanks to Raphael Menges for suggesting the name “Wikidated”.

References

1. Ahn, J., Im, D.H., Eom, J.H., Zong, N., Kim, H.G.: G-Diff: A Grouping Algorithm for RDF Change Detection on MapReduce. In: JIST. Lecture Notes in Computer Science, vol. 8943, pp. 230–235. Springer (2014)

2. Balazevic, I., Allen, C., Hospedales, T.M.: TuckER: Tensor Factorization for Knowledge Graph Completion. In: EMNLP/IJCNLP (1). pp. 5184–5193. Association for Computational Linguistics (2019)
3. Berners-Lee, T., Connolly, D.: Delta: an ontology for the distribution of differences between RDF graphs. Tech. rep., W3C (2004), <https://www.w3.org/DesignIssues/lncs04/Diff.pdf>
4. Bhowmik, R., de Melo, G.: Explainable Link Prediction for Emerging Entities in Knowledge Graphs. In: ISWC (1). Lecture Notes in Computer Science, vol. 12506, pp. 39–55. Springer (2020)
5. Cyganiak, R., Wood, D., Lanthaler, M.: RDF 1.1 Concepts and Abstract Syntax. W3C Recommendation, W3C (2014), <https://www.w3.org/TR/2014/REC-rdf11-concepts-20140225/>
6. Daruna, A., Gupta, M., Sridharan, M., Chernova, S.: Continual Learning of Knowledge Graph Embeddings. *IEEE Robot. Autom Lett* **6**(2), 1128–1135 (2021)
7. Dettmers, T., Minervini, P., Stenetorp, P., Riedel, S.: Convolutional 2D Knowledge Graph Embeddings. In: AAAI. pp. 1811–1818. AAAI Press (2018)
8. Erxleben, F., Günther, M., Krötzsch, M., Mendez, J., Vrandečić, D.: Introducing Wikidata to the Linked Data Web. In: International Semantic Web Conference (1). Lecture Notes in Computer Science, vol. 8796, pp. 50–65. Springer (2014)
9. Hee, Q., Chen, B.C., Agarwal, D.: Building The LinkedIn Knowledge Graph (2016), <https://engineering.linkedin.com/blog/2016/10/building-the-linkedin-knowledge-graph>
10. Hernández, D., Hogan, A., Krötzsch, M.: Reifying RDF: What Works Well With Wikidata? In: SSWS@ISWC. CEUR Workshop Proceedings, vol. 1457, pp. 32–47. CEUR-WS.org (2015)
11. Hogan, A., Blomqvist, E., Cochez, M., d’Amato, C., de Melo, G., Gutiérrez, C., Kirrane, S., Labra Gayo, J.E., Navigli, R., Neumaier, S., Polleres, A., Ngonga Ngomo, A.C., Rashid, S.M., Rula, A., Schmelzeisen, L., Sequeda, J.F., Staab, S., Zimmermann, A.: Knowledge Graphs. *ACM Comput. Surv.* **54**(4), 71:1–71:37 (2021)
12. Huang, X., Zhang, J., Li, D., Li, P.: Knowledge Graph Embedding Based Question Answering. In: WSDM. pp. 105–113. ACM (2019)
13. Lantzaki, C., Papadakos, P., Analyti, A., Tzitzikas, Y.: Radius-aware approximate blank node matching using signatures. *Knowl. Inf. Syst.* **50**(2), 505–542 (2017)
14. Liu, J., Zhang, Q., Fu, L., Wang, X., Lu, S.: Evolving Knowledge Graphs. In: INFOCOM. pp. 2260–2268. IEEE (2019)
15. Malyshev, S., Krötzsch, M., González, L., Gonsior, J., Bielefeldt, A.: Getting the Most Out of Wikidata: Semantic Technology Usage in Wikipedia’s Knowledge Graph. In: International Semantic Web Conference (2). Lecture Notes in Computer Science, vol. 11137, pp. 376–394. Springer (2018)
16. Nielsen, F.Å.: Lexemes in Wikidata: 2020 status. In: LDL@LREC. pp. 82–86. European Language Resources Association (2020)
17. Noy, N.F., Gao, Y., Jain, A., Narayanan, A., Patterson, A., Taylor, J.: Industry-scale Knowledge Graphs: Lessons and Challenges. *ACM Queue Tomorrows Comput. Today* **17**(2), 20 (2019)
18. Papavasileiou, V., Flouris, G., Fundulaki, I., Kotzinos, D., Christophides, V.: High-level change detection in RDF(S) KBs. *ACM Trans. Database Syst.* **38**(1), 1:1–1:42 (2013)
19. Pellissier Tanon, T., Bourgaux, C., Suchanek, F.M.: Learning How to Correct a Knowledge Base from the Edit History. In: WWW. pp. 1465–1475. ACM (2019)

20. Pellissier Tanon, T., Suchanek, F.M.: Querying the Edit History of Wikidata. In: ESWC (Satellite Events). Lecture Notes in Computer Science, vol. 11762, pp. 161–166. Springer (2019)
21. Pernischová, R., Dell’Aglío, D., Horridge, M., Baumgartner, M., Bernstein, A.: Toward Predicting Impact of Changes in Evolving Knowledge Graphs. In: ISWC Satellites. CEUR Workshop Proceedings, vol. 2456, pp. 137–140. CEUR-WS.org (2019)
22. Piscopo, A., Kaffee, L.A., Phethean, C., Simperl, E.: Provenance Information in a Collaborative Knowledge Graph: An Evaluation of Wikidata External References. In: ISWC (1). Lecture Notes in Computer Science, vol. 10587, pp. 542–558. Springer (2017)
23. Piscopo, A., Simperl, E.: Who Models the World?: Collaborative Ontology Creation and User Roles in Wikidata. *Proc ACM Hum Comput Interact* **2**(CSCW), 141:1–141:18 (2018)
24. Pittman, R.J.: Cracking the Code on Conversational Commerce (2017), <https://www.ebayinc.com/stories/news/cracking-the-code-on-conversational-commerce/>
25. Pomp, A., Kraus, V., Poth, L., Meisen, T.: Semantic Concept Recommendation for Continuously Evolving Knowledge Graphs. In: ICEIS (Revised Selected Papers). Lecture Notes in Business Information Processing, vol. 378, pp. 361–385. Springer (2019)
26. Rebele, T., Suchanek, F.M., Hoffart, J., Biega, J., Kuzey, E., Weikum, G.: YAGO: A Multilingual Knowledge Base from Wikipedia, Wordnet, and Geonames. In: International Semantic Web Conference (2). Lecture Notes in Computer Science, vol. 9982, pp. 177–185 (2016)
27. Sarasua, C., Checco, A., Demartini, G., Difallah, D.E., Feldman, M., Pintscher, L.: The Evolution of Power and Standard Wikidata Editors: Comparing Editing Behavior over Time to Predict Lifespan and Volume of Edits. *Comput Support Coop. Work* **28**(5), 843–882 (2019)
28. Schreiber, G., Raimond, Y.: RDF 1.1 Primer. W3C Note, W3C (2014), <https://www.w3.org/TR/2014/NOTE-rdf11-primer-20140624/>
29. Shenoy, K., Ilievski, F., Garijo, D., Schwabe, D., Szekely, P.A.: A Study of the Quality of Wikidata. *CoRR* **abs/2107.00156** (2021)
30. Singhal, A.: Introducing the Knowledge Graph: things, not strings (2012), <https://www.blog.google/products/search/introducing-knowledge-graph-things-not/>
31. Sun, Z., Deng, Z.H., Nie, J.Y., Tang, J.: RotatE: Knowledge Graph Embedding by Relational Rotation in Complex Space. In: ICLR (Poster). OpenReview.net (2019)
32. Tanon, T.P., Vrandečić, D., Schaffert, S., Steiner, T., Pintscher, L.: From Freebase to Wikidata: The Great Migration. In: WWW. pp. 1419–1428. ACM (2016)
33. The W3C SPARQL Working Group: SPARQL 1.1 Overview. W3C Recommendation, W3C (2013), <https://www.w3.org/TR/2013/REC-sparql11-overview-20130321/>
34. Tummarello, G., Morbidoni, C., Bachmann-Gmür, R., Erling, O.: RDFSyc: Efficient Remote Synchronization of RDF Models. In: ISWC/ASWC. Lecture Notes in Computer Science, vol. 4825, pp. 537–551. Springer (2007)
35. Vrandečić, D., Krötzsch, M.: Wikidata: A Free Collaborative Knowledgebase. *Commun. ACM* **57**(10), 78–85 (2014)
36. Wang, Q., Mao, Z., Wang, B., Guo, L.: Knowledge Graph Embedding: A Survey of Approaches and Applications. *IEEE Trans. Knowl. Data Eng.* **29**(12), 2724–2743 (2017)
37. Wu, T., Khan, A., Gao, H., Li, C.: Efficiently Embedding Dynamic Knowledge Graphs. *CoRR* **abs/1910.06708** (2019)