

Remixing Entity Linking Evaluation Datasets for Focused Benchmarking

Jörg Waitelonis^a, Henrik Jürges^b and Harald Sack^c

^a *yovisto GmbH, August-Bebel-Str. 26-53, 14482 Potsdam, Germany*

E-mail: joerg@yovisto.com

^b *University of Potsdam, Am Neuen Palais 10, 14469 Potsdam, Germany*

E-mail: juerges@uni-potsdam.de

^c *FIZ Karlsruhe, Leibniz Institute for Information Infrastructure, Hermann-von-Helmholtz-Platz 1, 76344*

Eggenstein-Leopoldshafen, Germany

E-mail: harald.sack@fiz-karlsruhe.de

Editor(s): Axel-Cyrille Ngonga Ngomo, Institute for Applied Informatics, Leipzig, Germany; Irini Fundulaki, ICS-FORTH, Heraklion, Greece; Anastasia Krithara, National Center for Scientific Research “Demokritos”, Athens, Greece

Solicited review(s): Michelle Cheatham, Wright State University, Dayton, Ohio, USA; Ziqi Zhang, University of Sheffield, UK; Heiko Paulheim, University of Mannheim, Germany; One Anonymous Reviewer

Abstract. In recent years, named entity linking (NEL) tools were primarily developed in terms of a general approach, whereas today numerous tools are focusing on specific domains such as e. g. the mapping of persons and organizations only, or the annotation of locations or events in microposts. However, the available benchmark datasets necessary for the evaluation of NEL tools do not reflect this focalizing trend. We have analyzed the evaluation process applied in the NEL benchmarking framework GERBIL [37,30] and all its benchmark datasets. Based on these insights we have extended the GERBIL framework to enable a more fine grained evaluation and in depth analysis of the available benchmark datasets with respect to different emphases. This paper presents the implementation of an adaptive filter for arbitrary entities and customized benchmark creation as well as the automated determination of typical NEL benchmark dataset properties, such as the extent of content-related ambiguity and diversity. These properties are integrated on different levels, which also enables to tailor customized new datasets out of the existing ones by remixing documents based on desired emphases. Besides a new system library to enrich provided NIF [11] datasets with statistical information, best practices for dataset remixing are presented, and an in depth analysis of the performance of entity linking systems on special focus datasets is presented.

Keywords: Entity Linking, GERBIL, Evaluation, Benchmark

1. Introduction

Named entity linking (NEL) is the task of interconnecting natural language text fragments with entities in formal knowledge-bases with the purpose to e. g. help subsequent processing tools to cope with ambiguities of natural language. NEL has evolved to a fundamental requirement for a range of applications, such as (web-)search engines, e. g. by mapping the content of search queries to a knowledge-graph [32] or

to improve search rankings [39]. By linking textual content to formal knowledge-bases, exploratory search systems as well as content-based recommender systems greatly benefit from the underlying graph structures by leveraging semantic similarity and relatedness measures [35]. Likewise, social media and web monitoring systems benefit from NEL, e. g. by the identification of persons or companies in social media content as subject of observation or tracking. A general survey on current NEL systems has been provided in [31,16].

While the number of application scenarios for NEL is on the increase, likewise the number of different NEL approaches is evolving ranging from simple string matching techniques to complex optimization based on machine learning [26]. Most NEL approaches make use of a general solution strategy, however there is an uprising trend for specialized solutions. In [43] the authors demonstrate an approach focused on medical literature while [8] examine heritage texts with NEL. Other approaches are focused on specific entity types, such as e.g. [7], which is applied to the domain of art. Another interesting solution is [1], which can be utilized to build domain specific NEL tools. The approach of [41] extracts semantic information from mixed media types like scientific videos. This ongoing fragmentation of types of tasks aggravates the application of generic benchmarking frameworks for NEL optimization and comparison such as GERBIL [37,30] or NERD [28,27].

With GERBIL, a NEL tool optimized for the detection of person names only might be rather difficult to compare to other NEL tools of a more general focus or specialized for another topic. However, the benchmark datasets provided with GERBIL are annotated with all types of entities including organizations, events, etc. Therefore, by using these general typed benchmarks the overall achieved results with GERBIL might only be hard to compare since the assumed person-only NEL system would wrongly be punished with false negatives caused by non-person annotations contained in the benchmarks. The only valid way to achieve an objective evaluation would be to manually filter a dataset to only contain persons and upload it to GERBIL for the desired experiment. However, these experiments are not reproducible, because it is neither clear or standardized, how the applied filtering was carried out, nor is the newly created filtered dataset always publicly available for further experiments. Moreover, it is not desirable to manage a plethora of different versions of filtered datasets. As of now, GERBIL deploys 19 annotation systems and more than 20 datasets, whereas these numbers are subject to constant change. For a detailed overview on the systems and datasets provided by GERBIL we refer to the official version¹. Besides the already described problem, there are also more challenges faced by the GERBIL framework considering the recent development of new NEL approaches. For instance, it is highly desirable to

be able to quantify the ‘difficulty’ of NEL problems presented in the different evaluation datasets, as e.g. the average degree of ambiguity, the completeness of annotations, etc.

A first attempt to cope with this problem was made in [12] by manually compiling the Kore50² corpus with the goal to capture hard to disambiguate mentions of entities. Another problem arises with the quality of annotations as described in [15] and [38] including e.g. annotation redundancy, inter-annotation agreement, topicality according to the evolving knowledge bases, mention boundaries, as well as nested annotations. Especially, completeness and coverage of annotations are essential measures to assess those annotation tasks (A2KB cf. [37]) where also the entity mention detection contributes to the overall results.

Since no ‘all-in-one’ perfect dataset has emerged in the past, which covers all the aspects sufficiently well, it would be beneficial to measure and provide dataset characteristics on the document level to subsequently allow a recompilation of documents across different datasets according to predefined criteria into a customized corpus. For example, for the already mentioned person-only annotation system these measures would help to specifically select only those documents, which exhibit a significant number of person annotations providing a predefined level of ‘difficulty’. Remixing evaluation datasets on the document level leads to a better and more application specific focus of NEL tool evaluation while simultaneously ensuring reproducibility.

We have already introduced an extension of the GERBIL framework enabling a more fine grained evaluation and in depth analysis of the deployed benchmark datasets according to different emphases [40]. To achieve this, an adaptive filter for arbitrary entities has been introduced together with a system to automatically measure benchmark dataset properties. The implementation including a result visualization are integrated in the publicly available GERBIL framework.

In this paper, we present the following contributions: the work presented in [40] is brought up-to-date, consolidated, and furthermore extended with

- new additional dataset measures,
- a stand-alone library to enable customized remixing of datasets,

¹<http://aksw.org/Projects/GERBIL.html>

²<https://datahub.io/de/dataset/kore-50-nif-ner-corpus>

- a vocabulary to enrich NIF-based datasets with additional statistical information,
- a subset of available datasets has been reorganized to enable benchmarking according to the different dataset properties, and
- an in depth analysis of the performance of different systems on the reorganized datasets is presented.

The paper is structured as follows: after this introductory section, measures to characterize NEL datasets are introduced in Sect. 2. Sect. 3 explains the GERBIL integration as well as the stand-alone library in detail, while Sect. 4 elaborates on the most interesting properties on datasets we have determined so far and presents more insights on the systems performances on the reorganized and focused datasets. Finally, Sect. 5 concludes the paper with a summary of the presented work and an outlook on ongoing and future research.

2. Measuring NEL Dataset Characteristics

NEL datasets have already been analyzed to great extent. We consider these analyses to identify their potential shortcomings to be able to introduce characteristics and measures to establish more differentiated analyses. In [15] the basic characteristics of 9 NEL datasets were introduced including the number of documents, number of mentions, entity types, and number of NIL annotations. In [34] a more detailed view on the distribution of entity types was given including mapping coverage, entity candidate count, maximum recall, as well as entity popularity. The overlap among datasets was investigated in [38], they also introduced the new measures confusability, prominence and dominance as indicators for ambiguity, popularity, and difficulty.

In this paper, amongst others also a subset of the proposed characteristics has been integrated into the GERBIL benchmarking system. Compared to previous work, where either a theoretical only or an experimental only treatment of the problem was presented, this paper contributes a ready to use implementation by means of extending the GERBIL source code³ and also provides a publicly available on-line service⁴. Besides the implementation of filtering the benchmark datasets

according to the desired characteristics, the tool instantly updates and visualizes the per annotation system results including statistical summaries. The integration into GERBIL enables a standardized, consistent, extensible as well as reproducible way to analyze and measure dataset characteristics for NEL.

Building on that we also provide a stand-alone library⁵ that computes the proposed metrics directly on NIF datasets. Without limiting the generality of the forgoing, the following explanations refer to the annotation (A2KB) as well as disambiguation tasks (D2KB) of the GERBIL framework. D2KB is the task of disambiguation of a given entity mention against the knowledge base. With A2KB, first entity mentions have to be localized in the given input text before the subsequent disambiguation task is performed. Hence, for most implementations D2KB can be seen as a sub task of A2KB.

Before introducing the dataset characteristics one by one the terminology is presented.

A dataset D is a set of documents $d \in D$. We define a document as the tuple $d = (d_t, d_a)$ where d_t is the document text and $|d_t|$ is the number of words within the text of the document d . d_a is a set of annotations belonging to the document d and $|d_a|$ is the number of annotations for the document d .

An annotation $a \in d_a$ is defined as the tuple $a = (s, e, i, l)$. s is the surface form of a which can be located in the document text d_t with its character index i , indicating the begin of the annotation, and the text length l , indicating the number of characters the annotation encloses to the right of index i . The corresponding linked entity is denoted with e .

Furthermore, we define E as the infinite set of entities and S as the infinite set of surface forms such that they are supersets of all other sets of the form E^x and S^x . Moreover, we define E^D as the set of entities within the dataset D and S^D as the set of surface forms within D .

In the appendix of this paper a complete listing of the mathematical notation is given for overview purposes.

The hereafter defined measures might refer to different levels: dataset level, document level, and annotation (or entity) level. Table 1 contains an overview on which measure is considered at a specific level.

Some of the introduced measures are distinguished between micro and macro measurements [4]. Macro

³<https://github.com/santifa/gerbil/>

⁴<http://gerbil.s16a.org/>

⁵<https://github.com/santifa/hfts>

Table 1

Overview of the introduced measures and the according levels of reference, where **ds** stands for dataset level, **doc** for document level **an** for annotation level).

Measure	Level
Not annotated	ds
Density	ds, doc
Prominence	ds, doc, an
Maximum recall	ds
Likelihood of confusion	ds, doc, an
Dominance	ds
Types	ds, doc, an

measurement aggregates the average results of each single document. Regarding document length, all documents have the same influence on the aggregated result. In contrast, the micro measurement takes the results of each document into account as if they would belong to one single document, which consequently increases the influence of larger documents.

The formal definition is provided for both measurements for density, likelihood of confusion, dominance, and maximum recall. All other definitions are provided as macro measurement if not stated otherwise.

2.1 Number of Annotations

In general, the number of annotations is a measure to estimate the size of the disambiguation context. The average number of annotations for a dataset $na : D \rightarrow \mathbb{R}$ is defined as:

$$na(D) = \frac{\sum_{d \in D} |d_a|}{|D|} \quad (1)$$

2.2 Not Annotated Documents

Some of the available benchmark datasets even contain documents without any annotations at all. Documents without annotations might lead to an increase of false positives in the evaluation results and thereby might cause a loss of precision. The fraction of not annotated documents for a dataset $nad : D \rightarrow [0, 1]$ is defined as:

$$nad(D) = \frac{|\{d : |d_a| = 0\}|}{|D|} \quad (2)$$

Empty documents might be a problem for the annotation task (A2KB), but not for the disambiguation only task (D2KB), where empty document annotations are simply omitted in the processing.

2.3 Missing Annotations (Density)

Similar to not annotated documents, missing annotations in an otherwise annotated document might lead to a problem with the A2KB task. Annotation systems potentially identify these missing annotations, which are not confirmed in the available ground truth and thus are counted as false positives. It is not possible to determine the specific number of missing annotations without conducting an objective manual assessment of the entire ground truth data, which requires major effort. However, we propose to estimate this number by measuring an annotation density value which is the fraction of the number of annotations and the document text length. The $density : D \rightarrow [0, 1]$ is defined as:

$$density_{micro}(D) = \frac{\sum_{d \in D} \frac{|d_a|}{|d_t|}}{|D|} \quad (3)$$

$$density_{macro}(D) = \frac{\sum_{d \in D} |d_a|}{\sum_{d \in D} |d_t|}$$

If an annotation is spanning more than one word, it is only counted as one annotation.

2.4 Prominence (Popularity)

The assumption of [38] is, that an evaluation against a corpus with a tendency to focus strongly on prominent or popular entities may cause bias. Hence, NEL systems preferring popular entities potentially exhibit an increase in performance. To verify this, we have implemented two different measures on the annotation level. Similarly to [38], the prominence is estimated as PageRank [22] of entities, based on their underlying link graph in the knowledge base. Additionally, we also take into account Hub and Authorities (HITS) values as a complementary popularity related score. PageRank as well as HITS values were obtained from [25].

To classify annotations, documents, and datasets according to different levels of prominence of entities, the set of entities was partitioned as follows. PageRank (respectively HITS) underlies a power-law distribution (cf. Sect. 4.2.1), meaning that only a few entities exhibit a high PageRank and the majority of entities a lower PageRank (long-tail), cf. Fig 1. Highly prominent entities are then defined as the upper 10% of the top PageRank values. The subsequent 45% (i.e. 10% – 55%) define medium prominence and the lower 45% (i.e. 55% – 100%) low prominence.

It is important to mention that for a dataset with a stronger bias towards head entities, the entities of the

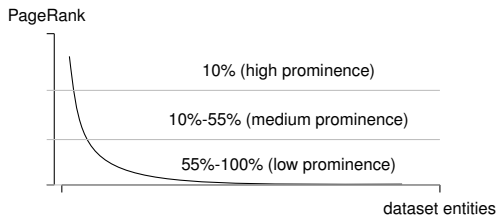


Fig. 1. Example partitioning for the PageRank.

middle or lower segment would then be in the higher segment for a dataset with a more even distribution. Thus, when working with multiple datasets, a global partitioning including all values of all entities is preferred.

For an arbitrary scoring algorithm P we can define the set of entities within a specific interval $a, b \in [0, 1]$ with $E_{a,b}^D : (P) \rightarrow E$ as:

$$E_{a,b}^D(P) = \{e \in E^D : a \leq P(e) \leq b\} \quad (4)$$

The resulting set contains all entities of a dataset that satisfies the given interval limits. A disadvantage of this approach is that entities, which do not have a score assigned, are not part of one of the resulting sets. Similarly the prominence can be determined using the HITS values or any other ranking score.

2.5 Likelihood of Confusion (Level of Ambiguity)

Since a surface form might denote multiple meanings as well as entities might be represented by different textual representatives the likelihood of confusion is a measure for the level of ambiguity for one surface form or entity. It was first proposed in [38] for surface forms. The authors pointed out that the true likelihood of confusion is always unknown due to a missing exhaustive collection of all named entities.

The likelihood of confusion needs some considerations beforehand. It can be determined for both sides of an annotation $a = (s, e, i, l)$. For a surface form s and the possible links to some entities E and for an entity e and the possible corresponding surface forms S .

We define a dictionary of an annotating system by W_E which is a mapping $W_E : S \rightarrow E$.

As shown in Fig. 2 the text `... Bruce ...` (lower box) has an annotation with ‘Bruce’ as surface form s . This surface form can be linked against different entities, i.e. they are homonyms, thus exhibiting the same writing but different meanings. As shown in the figure, an entity can belong to the dataset or is unknown to the

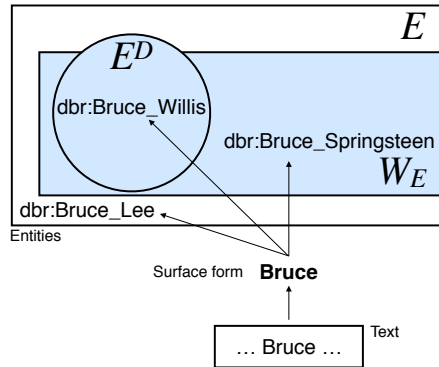


Fig. 2. The likelihood of confusion for a surface form is determined by the total number of possible entities known to some annotating system and a dataset $e^D \cup W_E$.

dataset but known to the annotating system. Also, the entity can be unknown to both sets.

For the other side we define a dictionary of the annotating systems W_S which is a mapping $W_S : E \rightarrow S$.

Fig. 3 shows the other side where the text annotation has `dbr:Bruce_Willis` as an entity. This entity can be linked against multiple possible surface forms which are synonyms. Again the surface form can be known to the dataset and the annotating system or unknown to one of them or both.

As already mentioned, a surface form s or an entity e can be placed within four possible locations:

1. Unknown to dictionary and dataset:
 $e \notin E^D \cup W_E$ or $s \notin S^D \cup W_S$
2. Only known to the dataset:
 $e \in E^D \setminus W_E$ or $s \in S^D \setminus W_S$
3. Only known to the dictionary:
 $e \in W_E \setminus E^D$ or $s \in W_S \setminus S^D$
4. Known to dictionary and dataset:
 $e \in E^D \cap W_E$ or $s \in S^D \cap W_S$

The example annotation system dictionaries W_E and W_S used for the experiments has been compiled from DBpedia entities’ labels, redirect labels, disambiguation labels, and `foaf:names`, if available.

For a dataset and a dictionary, the average likelihood of confusion is determined for surface forms $lc^{sf} : (D, W) \rightarrow \mathbb{R}^+$ with:

$$lc_{micro}^{sf}(D, W) = \frac{\sum_{d \in D} \frac{\sum_{a \in d_a} |W_E(s) \cup E^D(s)|}{|d_a|}}{|D|} \quad (5)$$

$$lc_{macro}^{sf}(D, W) = \frac{\sum_{s \in S^D} |W_E(s) \cup E^D(s)|}{|S^D|}$$

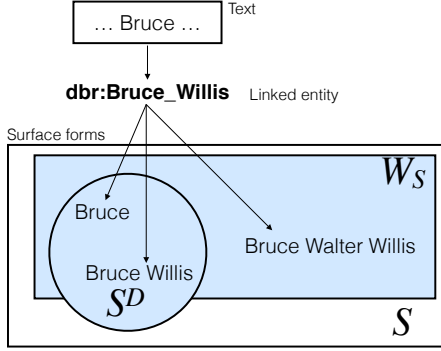


Fig. 3. The likelihood of confusion for an entity mention is the number of possible related surface forms shown in light blue.

The intuition is, the more entities exist per surface form, the larger is the likelihood of confusion lc^{sf} .

The average likelihood of confusion for entities $lc^e : (D, W) \rightarrow \mathbb{R}^+$ is:

$$lc_{micro}^e(D, W) = \frac{\sum_{d \in D} \frac{\sum_{a \in d_a} |W_S(e) \cup S^D(e)|}{|d_a|}}{|D|} \quad (6)$$

$$lc_{macro}^e(D, W) = \frac{\sum_{e \in E^D} |W_S(e) \cup S^D(e)|}{|E^D|}$$

Here the intuition is, the more surface forms exist per entity, the larger is the likelihood of confusion lc^e .

Again, an annotation within a dataset contains a surface form and an entity. For each side (surface form or entity) the likelihood of confusion is determined by counting the elements belonging to this particular side.

The measures should roughly indicate the difficulty distribution of a dataset.

2.6 Dominance (Level of diversity)

In [38] the dominance was introduced as a measure of how commonly a specific surface form is really meant for an entity with respect to other possible surface forms. A low dominance in a dataset leads to a low variance for an automated disambiguation system and to possible over-fitting. Similar to the likelihood of confusion, the true dominance remains unknown. Again, in addition to the work presented in [38] we estimate dominance for both sides of an annotation $a = (s, e, i, l)$: for the entities as well as surface forms. For an entire dataset and a dictionary, the average dominance is also determined in both directions.

For example the entity `dbr:Angelina_Jolie`, let there exist 4 different surface forms in the dataset,

while the dictionary provides overall 10 surface forms, which results in a 40% dominance of the entity `dbr:Angelina_Jolie` in the considered dataset. The dominance of an entity determines how many different surface forms of this entity are used in the dataset (synonyms).

As example for the other side, for the given surface form ‘Anna’ the dictionary provides 10 different entities, while the dataset only uses 2 entities for different mentions of the surface form ‘Anna’, which results in a 20% dominance of ‘Anna’ for the dataset under consideration. The dominance of a surface form determines how many different entities are used with this surface form in the dataset (homonyms). It indicates the variance or flexibility of the used vocabulary and expresses the dependency on context. Dominance indicates the expressiveness of the used dataset. An extensive one exhibits more diversity. The dominance of a dataset is closely related to the likelihood of confusion since it describes the coverage among the dataset and dictionary.

The average dominance for a dataset D is determined for all entities E^D with $dom^e : (W, D) \rightarrow \mathbb{R}^+$ and for surface forms S^D with $dom^{sf} : (W, D) \rightarrow \mathbb{R}^+$.

$$dom_{micro}^{sf}(D, W) = \frac{\sum_{d \in D} \frac{\sum_{a \in d_a} \frac{E^d(s)}{|W_E(s)|}}{|d_a|}}{|D|} \quad (7)$$

$$dom_{macro}^{sf}(D, W) = \frac{\sum_{s \in S^D} \frac{|E^D(s)|}{|W_E(s)|}}{|S^D|}$$

$$dom_{micro}^e(D, W) = \frac{\sum_{d \in D} \frac{\sum_{a \in d_a} \frac{S^d(e)}{|W_S(e)|}}{|d_a|}}{|D|} \quad (8)$$

$$dom_{macro}^e(D, W) = \frac{\sum_{e \in E^D} \frac{|S^D(e)|}{|W_S(e)|}}{|E^D|}$$

Since the actual dominance is unknown and the completeness of the applied dictionaries cannot be guaranteed, computed values above the nominal threshold of 1.0 are possible. These results refer to an incomplete dictionary, i.e. there are more patterns used in the dataset than the applied dictionary does contain. The subsequently described maximum recall takes care of this aspect.

2.7 Maximum Recall

Most of the NEL approaches apply dictionaries to look up possible entity candidates matching a given surface form. If the dictionary doesn't contain an appropriate mapping for the surface form the annotation system is unable to identify a possible entity candidate at all.

As Fig. 3 shows and as already mentioned before some parts of the dataset might not be contained within the dictionary. Surface forms not in the intersection are unlikely to be found by entity linking since the annotation systems are using dictionaries to look up potential relations. Therefore, an incomplete dictionary limits the performance of an NEL system since an unknown surface form will lead to a loss in precision. So the maximum recall can be seen as an artificial limit of a dataset.

To estimate the coverage of a mapping dictionary, the maximum recall measurement was introduced by [34].

For a dictionary W_S and a dataset the maximum recall is defined as $mr : (D, W) \rightarrow [0, 1]$:

$$mr_{micro}(D, W) = \frac{\sum_{d \in D} (1 - \frac{|S^d \setminus W_S|}{|S^d|})}{|D|} \quad (9)$$

$$mr_{macro}(D, W) = 1 - \frac{|S^D \setminus W_S|}{|S^D|}$$

2.8 Types

Since some NEL approaches might be focused on a specific domain or handle some entity categories in a different way, a filter has been implemented to distinguish dataset entities by their type. Besides the focus of NEL approaches in [38] it is also stated that types of entities may be differently difficult to disambiguate such as person names (esp. first names) might be more ambiguous and country names more or less unique. A type filter for some type T and E^T denoting the set of all entities for T is defined as $E^D : (T) \rightarrow E$:

$$E^D(T) = \{e \in E^D : e \in E^T\}. \quad (10)$$

Following these theoretical considerations, the extensions of the GERBIL framework and how the determined characteristics are exploited will be described in the subsequent sections.

3. Implementation

This section describes the implementation of the GERBIL extension and the standalone library. Furthermore, the vocabulary to integrate the calculated statistics in the NIF annotation model are explained in detail.

3.1. Extending GERBIL

Two new components have been implemented to extend the GERBIL framework: one component to filter and isolate subsets of the available datasets, and a second component to calculate aggregated statistics about the data (sub-)sets according to the newly introduced measures. It is important to mention that these filters and calculations can also be applied to newly uploaded datasets. Thus, the system can also be used to gain insights about any arbitrary 'non-official' datasets not yet part of the GERBIL framework. The implemented filter-cascade is of a generic type and can be adjusted via customized SPARQL queries. For example, to filter a dataset to only contain entities of type `foaf:Person` the following filter configuration has to be applied:

```
name=Filter Persons
service=http://dbpedia.org/sparql
query=select distinct ?v where {
  values ?v {##} .
  ?v rdf:type foaf:Person .
}
chunk=50
```

The `name` designates the filter in the GUI, `service` denotes an arbitrary SPARQL-endpoint, but also a local file encoded in RDF/Turtle can be specified to serve as the base RDF query dataset. The `query` is a SPARQL query that returns a list of entities to be kept in the filtered dataset. The `##` placeholder will be replaced with the specific entities of the dataset. To avoid the size limits for SPARQL queries, the `chunk` parameter can be specified to split the query automatically in several parts for the execution. Any number of filters can be specified to be included in the analysis. With the flexibility of configuring SPARQL-queries, filters of any complexity or depth can be specified.

To partition the datasets according to entity prominence (popularity) we have additionally implemented a filter to segment the datasets in three subsets containing the top 10%, 10% to 55%, and 55% to 100% of the

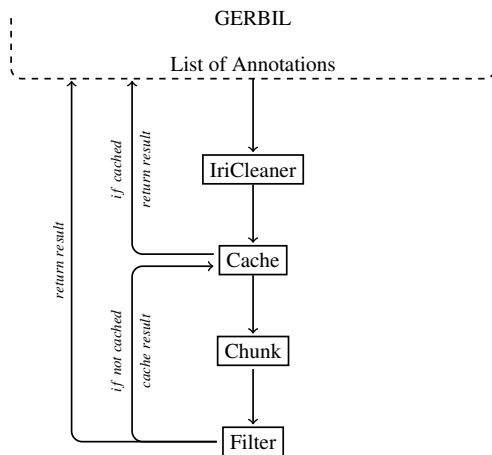


Fig. 4. Overview of the filter-cascade

entities. This segmentation is applied to PageRank as well as HITS values separately.

Fig. 4 shows a general overview of the filter cascade. The annotations produced by GERBIL are subsequently cleaned from invalid IRI's. If they are already cached the result is returned. Otherwise the set is chunked and passed to the defined filter.

Buttons have been added as new control elements to the A2KB, C2KB, and D2KB overview pages in GERBIL (cf. Fig. 5). The user now is able to choose between the classic view ‘no-filter’, the persons, places, organisations filter views, the PageRank/HITS top 10%, 10-55%, and 55-100% filter views, a comparison view, or a statistical overview. All implemented measures are visualized in GERBIL using HighCharts⁶. The existing charts are also replaced by the new chart API, since GERBIL was limited to only one single chart type. The comparison view enables the user to view two filters at the same time as well as the average for all annotation systems on a specific filter. The overview shows several statistics for all datasets, such as e. g., total number of types per filter, density, likelihood of confusion in average and total. A subset of these statistics is shown and discussed in section 4. The extended source code is publicly available at Github⁷. In addition, an online version of the system is available⁸.

Before discussing the dataset statistics as a result of the new GERBIL extension, the following section in-

⁶<http://www.highcharts.com/>

⁷<https://github.com/santifa/gerbil/>

⁸<http://gerbil.s16a.org/>

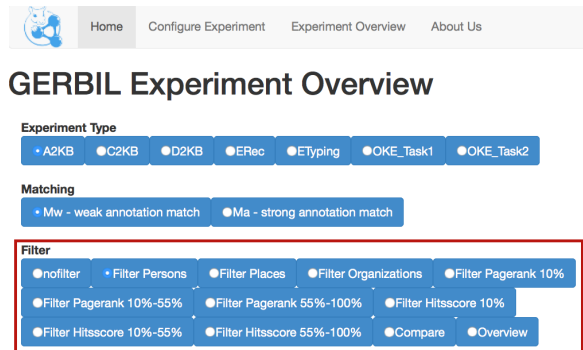


Fig. 5. New dataset filters for A2KB experiments in the GERBIL user interface.

roduces the stand-alone-library for statistics calculation as well as the new vocabulary.

3.2. Library and Vocabulary for Dataset Statistics

Following the considerations mentioned in the previous sections, the proposed measurements can also be calculated independently of GERBIL with a separate stand-alone library. The library consumes a NIF encoded input file, calculates the proposed statistics, and extends the NIF file with the newly determined information. A comprehensive documentation as well as the library source code is provided at Github⁹.

To serialize the calculated statistics generated by the GERBIL extension as well as by the library, a vocabulary has been defined with three layers to be integrated into the NIF model.

The first layer refers to an entity mention, respectively annotation, (e. g. NIF phrase) with its corresponding text fragment. The second layer addresses to the document (e. g. NIF context) that provides the text where the entity mentions are embedded. A third layer groups documents together to form a dataset. We introduce the `hfts:Dataset` class, which holds the documents with the `hfts:referenceDocuments` property. on the dataset level 13 properties have been introduced, which hold the measurements missing-annotation, density, maximum recall, dominance and likelihood of confusion on the dataset level. Some of them come with a micro as well as macro flavour while others are only computed once.

On the document level 6 new properties have been introduced to cover density, likelihood of confusion, and maximum recall. The likelihood of confusion,

⁹<https://github.com/santifa/hfts>

Table 2

Overview of the introduced properties and the corresponding measurements (**ds** stands for dataset level, **doc** for document level **an** for annotation level).

Measure	Property	Level
Not annotated	notAnnotated	ds
Density	microDensity	ds
	macroDensity	ds
Prominence	density	doc
	hits	an
Maximum recall	pagerank	an
	microMaxRecall	ds
Likelihood of confusion	macroMaxRecall	ds
	maxRecall	doc
	microAmbiguityEntities	ds
	macroAmbiguityEntities	ds
Dominance	ambiguityEntities	doc
	ambiguityEntity	an
	microAmbiguitySurfaceForms	ds
	macroAmbiguitySurfaceForms	ds
	ambiguitySurfaceForms	doc
	ambiguitySurfaceForm	an
Dominance	diversityEntities	ds
	diversitySurfaceForms	ds

prominence, and the types are also assigned on the entity mention level.

In Tab. 2 an overview over the introduced properties and their corresponding level is presented. Fig. 6 shows an excerpt of the extended Kore50 dataset for the new dataset class. One can see the new dataset statistics introduced by the RDF properties introduced by the *hfts:* prefix. In Fig. 7 an example for the document level is presented (*nif:Context*). Additionally to the existing NIF data the statistics have been serialized with the newly introduced *hfts:* properties. The entire definition and further documentation of the vocabulary is available at Github¹⁰.

Next, the possibility of remixing customized benchmark datasets will be explained including several examples.

3.3. Remixing Customized Datasets

The basic idea of remixing NEL benchmark datasets is to tailor new customized datasets from the exist-

¹⁰<https://raw.githubusercontent.com/santifa/hfts/master/ont/hfts.ttl##>>

```
<https://.../hfts/master/ont/nif-ext.ttl/kore50-nif>
a      hfts:Dataset ;
hfts:diversityEntities
      "0.0661871713645466"^^xsd:double ;
hfts:diversitySurfaceForms
      "0.08300283717687966"^^xsd:double ;
hfts:notAnnotatedProperty "0.0"^^xsd:double ;
hfts:referenceDocuments
      <http://.../KORE50.tar.gz/AIDA.tsv/CEL06#char=0,59> .
```

Fig. 6. An example of the new statistics properties on *dataset level* extending the KORE50 dataset.

```
<http://.../KORE50.tar.gz/AIDA.tsv/MUS03#char=0,97>
a      nif:RFC5147String , nif:String , nif:Context ;
nif:beginIndex "0"^^xsd:nonNegativeInteger ;
nif:endIndex "97"^^xsd:nonNegativeInteger ;
nif:isString "Three of the greatest ..."^^xsd:string ;
hfts:ambiguityEntities "17.0"^^xsd:double ;
hfts:ambiguitySurfaceForms "250.0"^^xsd:double ;
hfts:density "0.17647058823529413"^^xsd:double ;
hfts:maxRecall "1.0"^^xsd:double .
```

Fig. 7. An example of the new statistics properties on *document level* extending the KORE50 dataset.

ing ones by selecting documents based on desired emphases. This enables the compilation of focused benchmark datasets for NEL. For remixing it is proposed to store all analysed datasets in a single RDF triple store. This enables to quickly access the dataset documents via the SPARQL query language. In particular, SPARQL CONSTRUCT queries can be applied to select exactly those triples from the document annotations that meet a particular criteria, as e. g., popular persons, high possible maximum recall, places difficult to disambiguate, or any other arbitrary criteria, which can be expressed via SPARQL filter rules.

For this purpose, we introduce the basic query shown in Fig. 8. A CONSTRUCT statement creates RDF triples from document annotations meeting the filter requirement `maximumRecall >= 1.0`. This basic query utilizes the entire RDF induced graph and it might be useful to limit the number of documents that should be returned by the query. For this task, a subquery can be applied as shown in the second example in Fig. 9.

Another example is presented in Fig. 10. The SPARQL subselect chooses only documents that contain persons and aggregates their number. Subsequently, the CONSTRUCT statement selects documents that contain more than 4 persons with a maximum recall of at least 0.8.

To underline that any kind of filter can be applied, Fig. 11 shows a more specific example using a federated query to select only documents from the RDF graph with persons born before 1970. To achieve this,

```

# select document triples and annotation triples
CONSTRUCT {?doc ?dPredicate ?dObject .
           ?ann ?aPredicate ?aObject .}
WHERE {
# select all document triples
?ds hfts:referenceDocuments ?doc.
?doc ?dPredicate ?dObject .

# select all referenced annotations
?ann ?aPredicate ?aObject ;
    nif:referenceContext ?doc.

# use some filter condition
?doc hfts:maxRecall ?recall .
FILTER (xsd:double(?recall) >= 1.0) .
}

```

Fig. 8. Basic query that selects only documents with a maximum recall ≥ 1.0

```

# select document triples and annotation triples
CONSTRUCT {?doc ?dPredicate ?dObject .
           ?ann ?aPrediacte ?aObject .}
WHERE {
# get all document triples
?doc ?dPredicate ?dObject .

# limit the number of selected documents
{SELECT DISTINCT (?d AS ?doc)
 WHERE {
   ?ds hfts:referenceDocuments ?d.
   # use this instead of a global limit
   # to ensure only documents are limited
 } LIMIT 1
}
# select all referenced annotations
?ann ?aPredicate ?aObject ;
    nif:referenceContext ?doc.

# use some filter condition
}

```

Fig. 9. This query in addition limits the number of selected documents

the official DBpedia SPARQL endpoint is queried for additional information that is not present within the given benchmark datasets. More SPARQL examples can be found at Github¹¹.

For authoring arbitrary queries two aspects should be considered. First, many values of the proposed measurements are given as absolute values and are not always equally distributed across the datasets, documents, and annotations. Hence, it is necessary to investigate on the boundary values and value distribution before specifying a specific threshold. It is a subject of future work to normalize and harmonize the statistics adequately. Second, the proposed query examples are based on the document level. Therefore, if an annotation meets a requirement, the entire document together with all its annotations (which might not meet the requirement) is added to the result. Of course, queries can also be structured to only return the filtered anno-

¹¹<https://github.com/santifa/hfts/blob/master/Remix.md>

```

# document selection omitted
?doc hfts:maxRecall ?recall .

# use count for a later filter expression
{SELECT DISTINCT (?d AS ?doc) (COUNT(?a) AS ?aCount)
 WHERE {
   ?ds hfts:referenceDocuments ?d .
   # select matching entities
   ?a nif:referenceContext ?d ;
       itsrdf:taClassRef dbo:Person .
 } GROUP BY ?d LIMIT 100
}

# select referenced annotations omitted

# select only documents with more than three persons
# and a maximum recall of 0.8
FILTER(?aCount > 3) .
FILTER(xsd:double(?recall) >= 0.8) .

```

Fig. 10. Extract documents with a maximum recall of 0.8 and at least 4 person.

```

# construct block omitted
{SELECT DISTINCT (?d AS ?doc)
 WHERE {
   ?ds hfts:referenceDocuments ?d .
   # select matching entities
   ?a nif:referenceContext ?d ;
       itsrdf:taIdentRef ?ref ;
       itsrdf:taClassRef dbo:Person .

   # fetch data from another endpoint
   SERVICE <http://dbpedia.org/sparql> {
     ?ref dbo:birthDate ?date .
   }
   FILTER (?date <= xsd:date('1970-01-01')).
 }
}

```

Fig. 11. A SPARQL query that selects documents containing persons born before 1970 via additional data queried from the DBpedia SPARQL endpoint

tations, but this might lead to a missing annotation scenario that again might result in a drop of recall for the A2KB task.

Finally, the thereby newly created dataset can be uploaded to the GERBIL platform for a precisely tailored evaluation experiment.

4. Statistics and Results

This section presents the results of the execution of the proposed measures on the GERBIL datasets. Furthermore, an in depth overview on how to use the new library to partition the benchmarking datasets according to different criteria and to analyze the systems performances in much greater detail is presented.

4.1. GERBIL Datasets

The following datasets have been analyzed according to the characteristics introduced in Sect. 2: WES2015 [39], OKE2015 [21], DBpedia Spotlight [17],

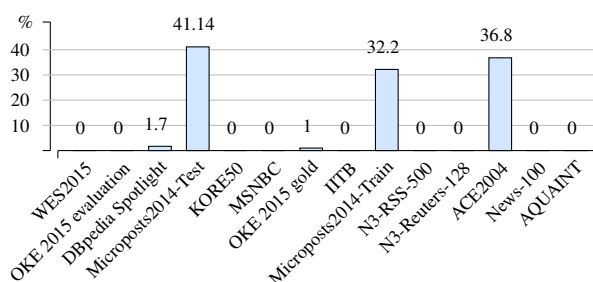


Fig. 12. Percentage of documents without annotations in the GERBIL datasets

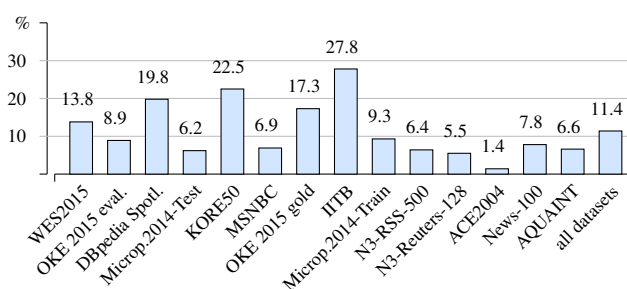


Fig. 13. Annotation density as relative number of annotations respective document length in words

KORE50 [12], MSNBC [5], IITB [14], RSS500 [29], Micropost2014 [2], Reuters128 [29], ACE2004 [19], AQUAINT [18], and NEWS-100 [29]. In this section, only the most significant results are presented. A complete listing of the achieved results is available online¹².

Fig. 12 shows the percentage of documents in the GERBIL datasets which were **not annotated**. Overall, there are 5 datasets that contain empty documents while 3 of them show a significant (i.e. >30%) number of empty documents. For A2KB tasks, these datasets might lead to an increased false positive rate and thus might lower the potentially achievable precision of an annotation system. Therefore, empty documents might be excluded from evaluation datasets to enable a sound evaluation. However, it should be noted that it is possible that these un-annotated documents are not actually mistakes but rather don't contain any entities.

Fig. 13 shows the **annotation density** of the GERBIL datasets as relative number of annotations with respect to document lengths in words. This serves as an estimation for potentially missing annotations, e. g. in

the IITB dataset 27.8% of all terms are annotated. If a dataset is annotated rather sparsely (low values), it is likely that the A2KB task will result in loss of precision, because the sparser the annotations the higher is the likelihood of potentially missing annotations (as it is shown in Sect. 4.2.7). Especially for NEL tools based on machine learning it should be considered, whether a sparsely annotated dataset is appropriate for the training task. Of course, this strongly depends on the according application. Nevertheless, it is arguable, if sparseness is problematic for A2KB, because all annotation systems are facing the same problem and the achieved results nevertheless might still be comparable.

Table 3 shows the **distribution of entity types and entity prominence** per dataset. A green (bold) label indicates the highest value and a red (italic) the lowest value in each category. Since not all entities can be linked with a type or affiliated with the ranking, the values for each partition do not necessarily sum up to 100%. For each dataset the percentage of entities per category is denoted, as e. g., of all the entities in the KORE50 dataset 47.1% are persons and 6.9% are places. In [34] it was demonstrated, there is a significant number of untyped entities in the DBpedia Spotlight and the KORE50 datasets. Therefore, an extra row for unspecified entities has been added to the table. The News-100 dataset exhibits the most unspecified entities because it is a German dataset and mostly contains annotations referring to the German DBpedia, but the analysis was based on the English DBpedia. The first partition (row 1–4) can be considered as an indicator of how specialized a dataset is. Thus, e. g., for the evaluation of an annotation system with focus on persons, the KORE50 dataset with 45.1% of person annotations might be better suited than the IITB dataset with only 2.4% of person annotations. The second and third partition (PageRank and HITS) show the entities categorized according to their popularity. It can be observed that many datasets are slightly unbalanced towards popular entities. A well balanced dataset should exhibit a relation of 10%, 45%, 45% among the three subset categories.

Fig. 14 shows the **average likelihood of confusion** to correctly disambiguate an entity or a surface form for several datasets. The blue bar (left) indicates the average number of surface forms that can be assigned to an entity, i. e. it refers to surface forms per entity, respectively synonyms. The red/hatched bar (right) shows the average number of entities that can

¹²<http://gerbil.s16a.org/>

Table 3
Percentage of entities by entity type and entity popularity per dataset

	WES 2015	OKE 2015 eval	DBpedia Spotl.	Microp. 2014 Test	KORE50	MSNBC	OKE 2015 gold	IITB	Microp. 2014 Train	N3-RSS-500	N3-Reuters-128	ACE2004	News-100	AQUAINT	all datasets
Persons	18.4	30.3	3.0	16.6	45.1	27.2	29.3	2.4	16.2	15.9	6.5	6.5	1.0	7.8	16.16
Org.	3.4	11.1	3.0	9.0	16.0	9.0	18.3	2.0	13.8	10.5	20.7	20.3	0.6	0.5	9.9
Places	9.4	14.0	8.2	8.9	6.9	17.5	14.5	3.5	14.2	7.2	17.2	35.0	0.3	24.5	13.0
unspecified	68.8	44.6	85.1	65.5	32	46.3	37.9	92.1	55.8	66.4	55.6	38.2	98.1	55.5	60.14
PageRank															
10%	27.9	24.4	30.0	21.3	28.5	28.5	24.9	14.8	26.0	14.3	18.8	22.2	0.4	25.3	22.0
10%-55%	48.9	39.5	47.6	49.8	48.6	32.2	0.3	29.8	45.8	23.0	31.4	37.6	1.1	43.7	34.2
55%-100%	22.5	16.6	19.7	28.0	19.4	24.8	7.7	15.0	25.6	11.1	19.0	15.1	0.7	23.9	17.8
HITS															
10%	28.4	21.1	32.4	31.4	27.8	29.8	26.9	12.3	32.9	18.3	19.0	28.4	0.4	29.0	24.2
10%-55%	12.9	12.4	18.2	14.4	20.8	22.8	0.3	12.2	13.6	7.3	9.1	11.4	1.4	45.3	14.4
55%-100%	58.0	47.0	48.2	51.8	47.2	32.1	50.2	35.2	50.6	23.2	40.6	15.3	0.4	18.9	37.1

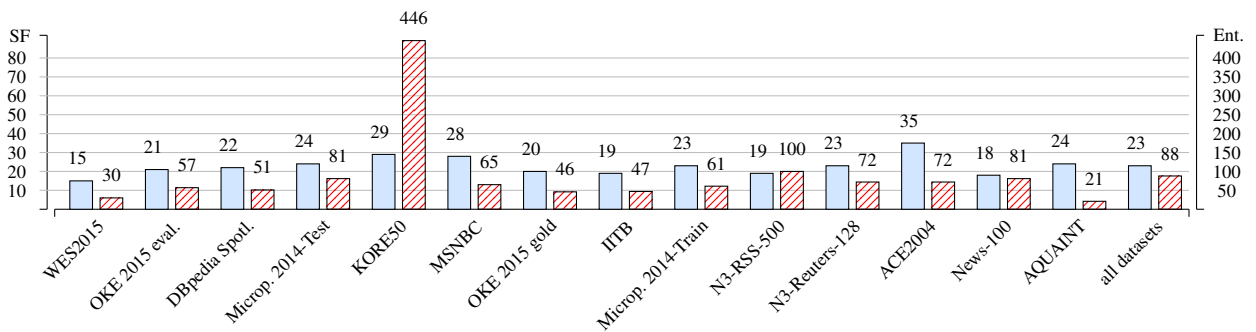


Fig. 14. Average number of surface forms (SF) per entity (blue, left) and average number of entities per surface form (red/hatched, right) indicating the likelihood of confusion for each dataset

be assigned to a surface form, i. e. it refers to entities per surface form, respectively homonyms. The figure shows clearly that KORE50 uses surface forms with a high number of potential entity candidates, i. e. it contains a large number of homonyms. Since this dataset is focused on persons it is not surprising that surface forms representing first names, such as e. g. ‘Chris’ or ‘Steve’, can be associated with a large number of corresponding entity candidates. KORE50 was compiled

with the aim to capture hard to disambiguate mentions of entities, which is confirmed by these observations. ACE2004 exposes the highest average number of surface forms for possible entities (35), i. e. it contains many synonyms.

In Section 4.2.2 a correlation analysis between the likelihoods of confusion for entities and surface forms with precision and recall is presented.

Fig. 15 shows the **average dominance of entities and surface forms** in percent. The red/hatched bars show the *average dominance of entities*. The dominance of an entity expresses the relation between an entity’s surface forms used in the dataset with respect to all its existing surface forms in the dictionary. Referring to Fig. 15, the KORE50 dataset uses only 9% of the surface forms that are provided in the dictionary. This indicates also how well the dataset’s surface forms are covered by the dictionary’s surface forms.

On the other hand, the blue bars show the *average dominance of surface forms*. The dominance of a surface form expresses the relation of how many entities are using this surface form in the considered dataset and the overall number of entities in the dictionary using this surface form.

Referring to Fig. 15, the KORE50 dataset in which many persons are annotated uses only 7% of the possible entities for the contained surface forms. In average, entities are represented in the WES2015 dataset with 21% of their surface forms.

Since the datasets with a high likelihood of confusion have a low dominance, it is arguable that these two measures express somehow the contrary. For example, the KORE50 dataset has a high likelihood of confusion for surface forms with 446 entities for one surface form on the average. This means that for a high dominance each surface form is represented by more than 400 entities within this dataset. Such a high dominance means also that a high coverage of surface forms (dominance of entities) or entities (dominance of surface forms) is present. For example, in the WES2015 dataset, which is focused on blog posts on rather specific topics, many rare entities (i.e. entities with a low popularity) with many different notations are used resulting in a likelihood of confusion of 15 surface forms for an entity on the average. The average dominance of entities is quite high with 21%, since the likelihood of confusion is low and topic specific blog posts often vary the surface forms for an entity to enrich the spirit-fulness of the text. This is commonly known from articles or essays, where the author usually tries to minimize frequent repetitions of surface form by varying the surface form for the entity under consideration to avoid monotony and to make the article more interesting to read. It might be concluded that a high dominance covers the diversity of natural language more precisely and therefore could be considered a means to prevent overfitting.

The News-100 dataset shows an anomaly in the dominance of entities, which is larger than 100%. The

reason for that is that the dataset contains a large number of entities from the German DBpedia. For these entities a surface form cannot be found in the dictionary (which was generated from the English DBpedia). That means, there are more surface forms present in the dataset than in the dictionary, which results in a dominance value larger than 100%.

This section has introduced and discussed the results of the statistical dataset analysis. Based on the information embedded in the NIF dataset files, a customized reorganisation of datasets can be accomplished as explained in the following section.

4.2. Insights from Remixing Datasets

To gain more insights on the interplay of annotation systems performance and the introduced dataset characteristics, this section describes how the datasets are reorganized to determine each system’s performance with focus on a given measure.

The approach is to first combine the datasets into one large dataset and then divide it into partitions. Each partition contains only those annotations or documents that lie in a specified interval of values of one of the proposed measures. For this purpose and to insert the statistical data into the NIF document the proposed library has been applied. Subsequently, the entire dataset was stored in an RDF triple store. With the SPARQL queries proposed in the previous sections, each partition was constructed and stored in a separate NIF document, which was submitted to the official GERBIL service to acquire the results.

For the conducted experiments the following public and GERBIL ‘shipped’ datasets have been used: DBpedia Spotlight, KORE50, Reuters128, RSS500, ACE2004, IITB, MSNBC, News100, AQUAINT. Other available datasets were either not publicly available or not in the NIF format.

Since the official GERBIL service was used to conduct the experiments, the therewith provided systems are included in the experiments. Unfortunately, not all systems returned consistent results due to too many errors or insufficient availability. However, if sufficient results could be provided, the system was included in the analysis.

The following annotation systems provided by GERBIL have been used: AGDISTIS [36], AIDA [13], Babelfy [20], DBpedia Spotlight [17], Dexter [3], Entityclassifier.eu [6], FOX [33], Kea [42], WAT [23] and PBOH [9].

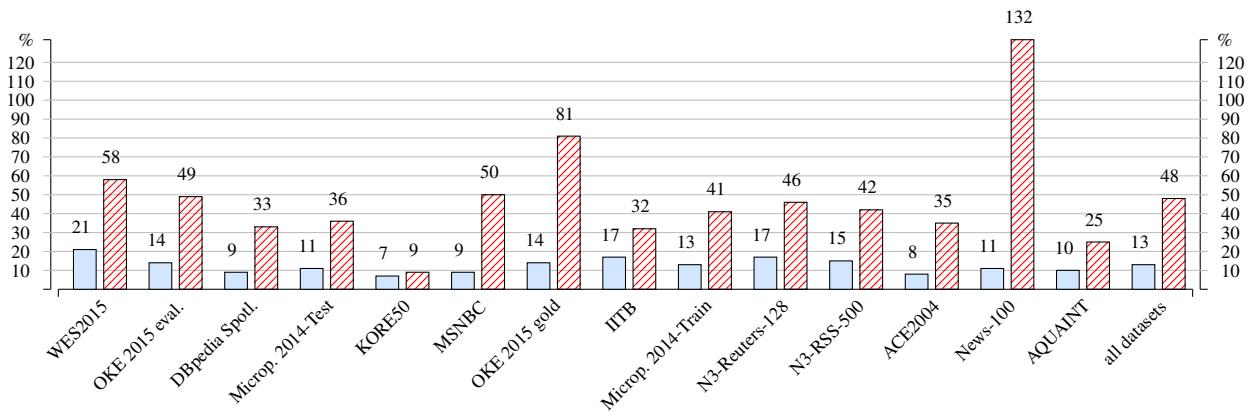


Fig. 15. Average dominance for surface forms (blue) and entities (red/hatched) per dataset

The measures used in the subsequent experiments are the measures currently supported by the library (i.e. likelihood of confusion, HITS, PageRank, density, and numbers of annotations). In general, both the A2KB as well as D2KB types of experiments, might be applied. For likelihood of confusion, HITS and PageRank only D2KB is provided because these are characteristics of the annotations. Number of annotations as estimation for the size of the disambiguation context is used with A2KB and D2KB types of tasks, density as characteristic of documents is used with A2KB only. All data as well as the achieved results can be found online¹³

4.2.1. Value distribution and partitioning

Fig. 16 presents the distribution of the data values over all datasets. In total, the dataset contains 16,821 annotation in 1043 documents. The figure shows a distribution chart for each measure. On the charts, the x-axis shows the number of annotations (for confusions, HITS, PageRank) or documents (for density and number of annotations). The y-axis shows the absolute values of the measures. Each of the charts approximate a power-law distribution, i.e. only a few items exhibit large values and many items smaller values. For HITS and PageRank only 14,372 items are available, because for 2,449 entities no HITS or PageRank value could be determined.

We have decided to apply a decile partitioning. It seems a reasonably well choice to indicate low, medium, large as well as the boundary values. When partitioning on the item values an uneven distribution of values over the partitions occurs because of the

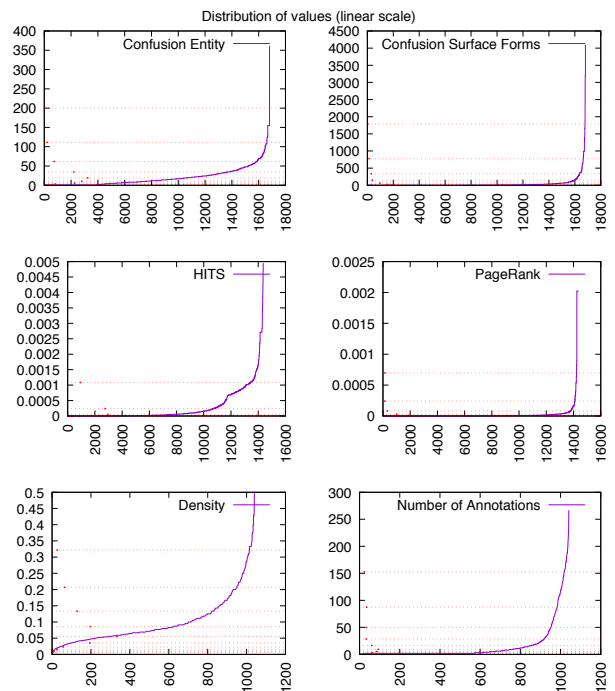


Fig. 16. Distribution of values (linear scale).

power-law, i.e. the first partition would contain a very large disproportionate number of items and the last partition only a very few items. To achieve a more even distribution a logarithmic scaling on the values is applied as shown in Fig. 17. The red horizontal dashed lines indicate the partition boundaries. Table 4 shows for each measure the threshold values (thr) for the partition boundaries as well as the number of items per partition (qty). For HITS and PageRank an additional partition was introduced to also include the items with-

¹³<https://github.com/santifa/hfts/blob/master/Results.md>

Table 4
Partitioning thresholds (log-based) and annotation/document quantities (this table is best viewed in color).

Part.	Conf. Surf.		Conf. Ent.		PageRank		HITS		Num. Anno.		Density	
	thr	qty	thr	qty	thr	qty	thr	qty	thr	qty	thr	qty
0	<2	8143	<2	3946	unspec.	2449	unspec.	2449	<2	20	<0.009	4
1	5	1368	3	599	<1.39E-07	3211	<5.77E-09	2456	3	595	0.015	10
2	12	1893	6	812	4.03E-07	1341	2.63E-08	19	5	63	0.023	26
3	28	2026	11	2256	1.17E-06	1504	1.20E-07	200	9	86	0.035	58
4	64	1581	19	2802	3.39E-06	2072	5.48E-07	446	16	93	0.055	194
5	147	963	34	3245	9.85E-06	2753	2.50E-06	819	29	61	0.086	333
6	338	382	62	2204	2.86E-05	1869	1.14E-05	1474	50	33	0.133	197
7	777	297	111	744	8.29E-05	1010	5.21E-05	2314	87	33	0.207	129
8	1786	128	200	203	2.40E-04	331	2.38E-04	2960	153	35	0.322	65
9	4105	40	361	10	6.98E-04	135	0.001	2744	267	24	0.500	27
10					0.002	146	0.005	940				

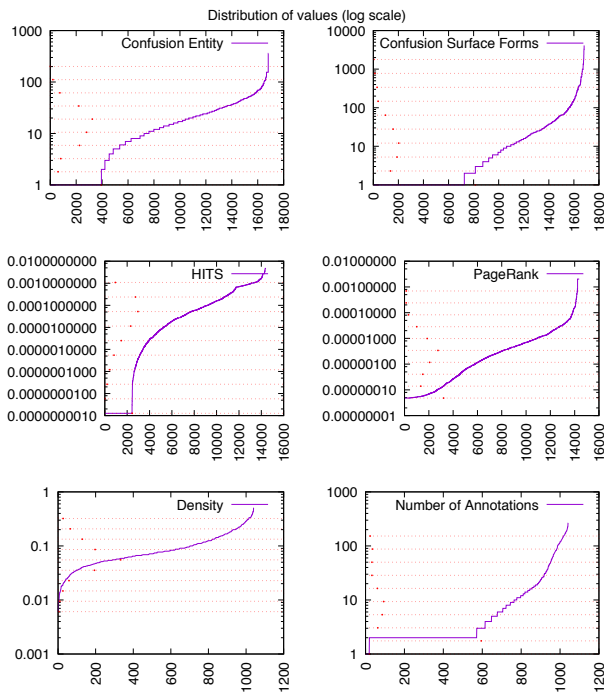


Fig. 17. Distribution of values (log scale).

out a value (unspec.). Each threshold is meant as the upper boundary of the partition, thus the lower boundary is the threshold of the previous partition. The color coding in the background of the cells will be explained later.

4.2.2. Likelihood of confusion of surface forms

Fig. 18 shows the experimental results of each system for the likelihood of confusion of surface forms. Each graph shows the partitions (x-axis), as well as

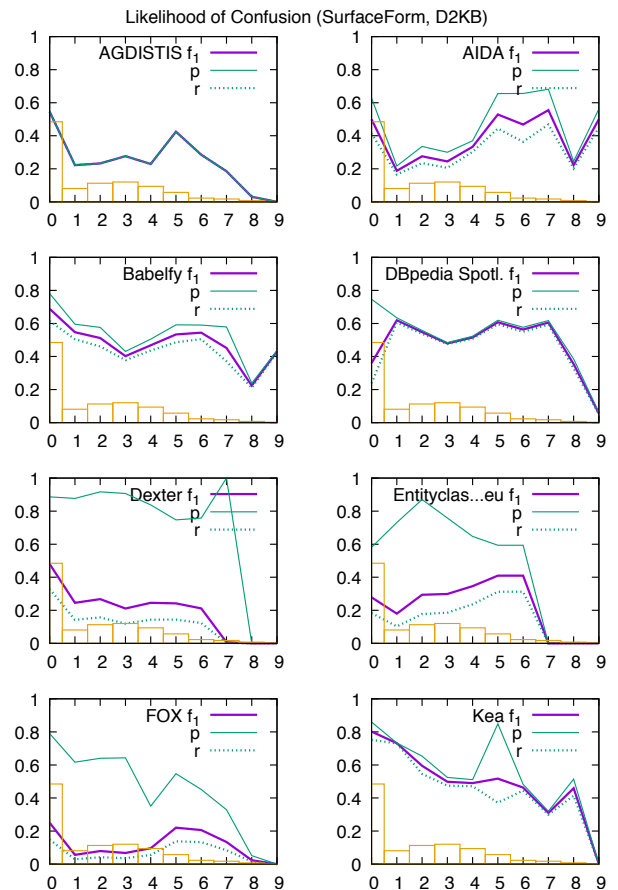


Fig. 18. Likelihood of confusion for surface forms (D2KB)

the determined F_1 -measure (f_1), precision (p), and recall (r) for each partition. In the background the relative sizes of the partitions are indicated with boxes (see Tab. 4 for specific values).

The likelihood of confusion for surface forms describes the number of entities mapping to one particular surface form. For an annotation in the dataset, a confusion of 30 signifies that 30 possible entities for that surface form exist (homonymy).

The leftmost partition (0) contains lower values, thus annotations contain surface forms with fewer numbers of entities mapping to them and therefore a lower likelihood of confusion. Typical are for example surface forms mentioning full names, as e.g., ‘Britney Spears’, ‘Northwest Airlines’, or ‘JavaScript’. The rightmost partition (9) shows larger values. It is expected that the annotations in the right partitions are more difficult to disambiguate since they exhibit a larger likelihood of confusion. The first partition contains almost half of all values, indicating that for almost half of the annotations only one entity maps to the surface form. For the second to sixth partition a reasonable even distribution is given. Considering Tab. 4, only 40 items are in the rightmost partition. These include the names Allen, Bill, Bob, Carlos, David, Davis, Eric, Jan, John, Johnson, Jones, Karl, Kim, Lee, Martin, Mary, Miller, Paul, Robert, Ryan, Steve, Taylor, and Thomas.

This experiment was applied as a disambiguation task (D2KB)¹⁴. However, the Entityclassifier.eu system did not provide results for partitions 7,8, and 9 (set to zero). WAT and PBOH created too many errors and have been excluded in this experiment.

To interpret the figures in general, the presented graphs show a trend from the upper left to the lower right, meaning that the systems performance decreases with growing likelihood of confusion. Many systems, except AIDA and Babelfy, fail with surface forms having more than ca. 1,700 entities mapping to (8th partition and above). Entityclassifier.eu, Dexter, and FOX show a very strong focus on precision, at the expense of recall, as we can also see in the further experiments.

It can be concluded that the fewer entities are mapping to a particular surface form, the easier seems the disambiguation task. For surface forms with more than 1,700 potential entity candidates the reliability of the disambiguation might drop dramatically.

4.2.3. Likelihood of confusion of entities

Fig. 19 shows the experimental results of each system for the likelihood of confusion of entities. The graphs are presented in the same way as for the previ-

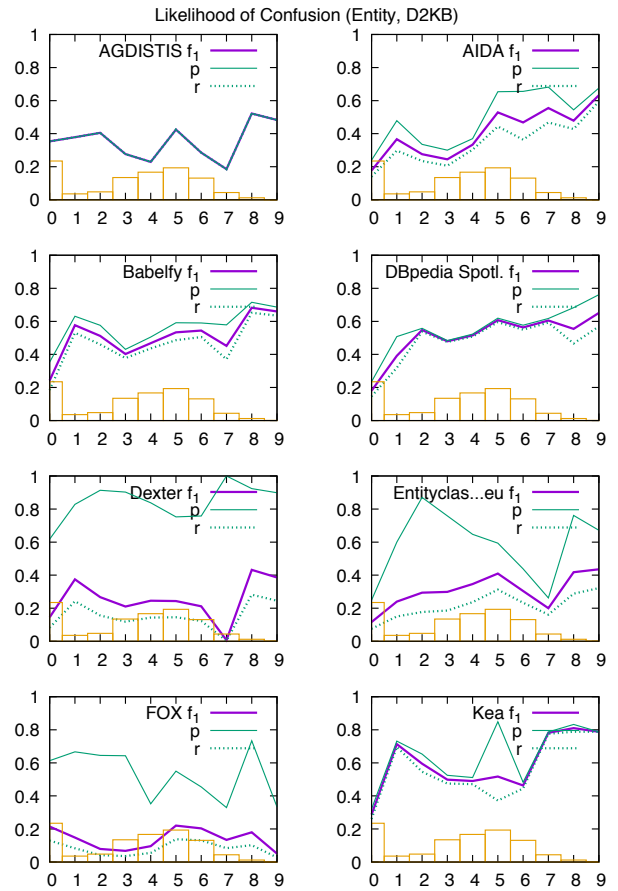


Fig. 19. Likelihood of confusion for entities (D2KB)

ous measure. The likelihood of confusion for entities describes to how many surface forms the entity of an annotation is mapping to. For an annotation, a confusion of 30 means that 29 surface forms besides the one within the annotation share the same entity.

The leftmost partition (0) contains lower values, thus annotations with entities mapping to only one surface form. The rightmost partition (9) contain annotations with entities mapping to more than 361 surface forms e.g. dbp:United_States. The number of items across the partitions is more evenly distributed than for the previous measure.

This experiment was applied as disambiguation task (D2KB)¹⁵. All participating systems except WAT and PBOH returned valid results, Entityclassifier.eu returned several faulty results.

¹⁴<http://gerbil.aksw.org/gerbil/experiment?id=201712060006>

¹⁵<http://gerbil.aksw.org/gerbil/experiment?id=201712050002>

In general, there is an upward trend, i.e., the more surface forms are available for an entity, the better it is. However, almost all systems have in common, that the performance drops rather abruptly on the first partition (0) compared to the second partition (1). A closer look on the partition data revealed that a large share of the entities in partition 0 are resources originating from Wikipedia redirect and disambiguation pages (e.g. `dbp:Diesel`, `dbp:Thermoelectricity`). Typically, these resources only map to a single surface form, which is why they occur in partition 0. Assumably, the systems are not annotating redirect and disambiguation resources, since they prefer to use the main resource and not resources directing to it. Some datasets show a drop at partition 7, but the partition data does not show obvious anomalies. Since we only can access the performance values provided by the GERBIL experiments, and therefore cannot access the actual annotations systems results, it is impossible to further investigate on that now.

Overall, it can be concluded that the more surface forms an entity is mapping to, the better the systems performances are. Furthermore, the datasets containing a larger number of redirect and disambiguation resources can bias the systems performances. Future work will repeat this analysis without bias to gain insights about, how well the systems really perform on the first partition.

4.2.4. PageRank

Fig. 20 shows the systems performances on the popularity estimation via PageRank values. Now, an additional partition is included in the graphs, which is located left (partition 0) showing the results on the 2,449 annotations, where no PageRank was given. For all other partitions, the PageRank values increase from left to right. Thus, popular entities can be found on the right hand. The distribution of values across the partitions is reasonable even.

The experiments were conducted as D2KB task¹⁶. With the exception of Entityclassifier.eu and FOX, all systems returned error free results. For the time of the execution of these experiments, also the WAT system was available. PBOH was not available.

In the graph a general uprising trend can be observed, i.e. popular entities are better disambiguated than unpopular entities, but with the exception of

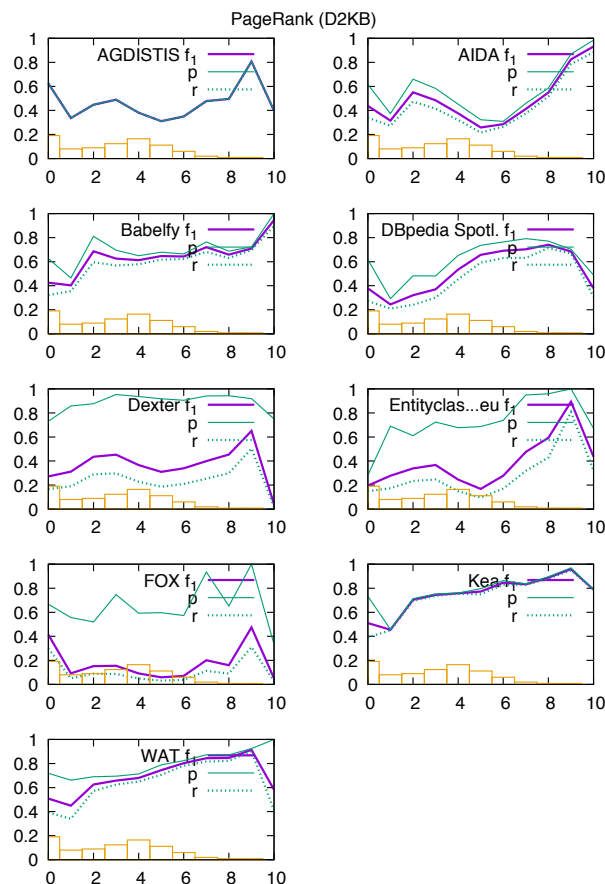


Fig. 20. Results for PageRank (D2KB)

AIDA and Babelify, all systems struggle with extremely popular entities (partition 10). A view in the data revealed that the 146 annotations only refer to the 4 entities `dbp:Germany`, `dbp:United_States`, `dbp:Americas` and `dbp:Animal`. It might be that some of the effect comes from the confusion of `dbp:United_States` and `dbp:Americas`. Therefore, partition 10 might not be sufficiently representative. The entities with the largest PageRanks (e.g. from partition 8) mostly refer to countries and popular locations as well as to the entity `dbp:Insect`.

In conclusion, a positive correlation (>0.7) between the PageRank values and the systems performances can be observed. It seems likely that popular entities are used much more frequently, while being described via many varying surface forms.

4.2.5. HITS

Similarly to PageRank, HITS values were not provided for all entities, thus partition 0 contains the annotations with unspecified values (see Fig. 21). For the

¹⁶<http://gerbil.aksw.org/gerbil/experiment?id=201712060001>

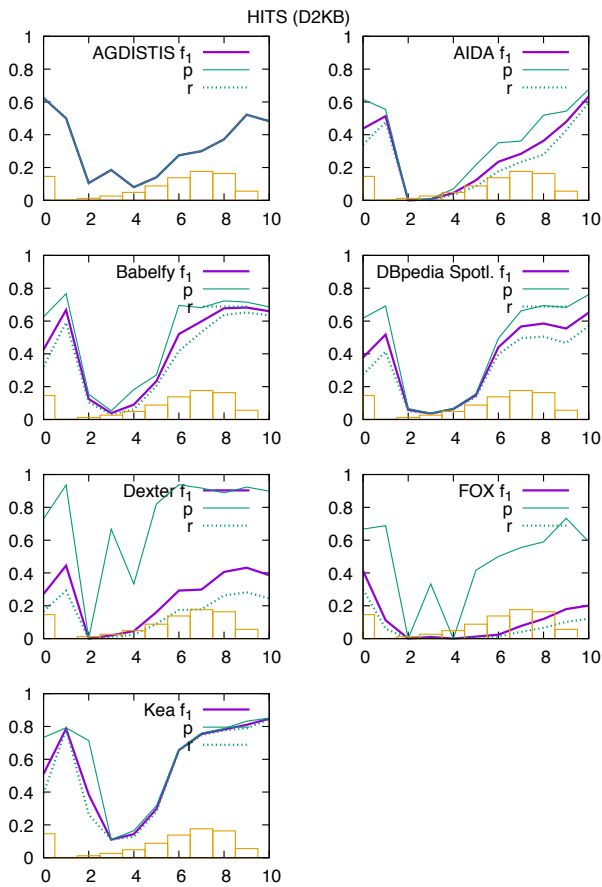


Fig. 21. Results for HITS (D2KB)

other partitions the HITS values are increasing from left to right. According to Tab. 4, partition 2 contains only very few annotations (19). The other partitions contain a more representative number of items.

Again, the experiments were conducted as D2KB tasks¹⁷. However, the Entityclassifier.eu, WAT, and PBOH produced too many faulty results and had to be excluded from the evaluation.

The HITS analysis reveals that for very low values (partition 1) and higher values (partition 6 and upwards) the systems provide better results than for the medium values (partitions 2-5). There is a weak correlation among HITS and the systems performances (>0.4). This could be interpreted as with increasing partition number there are more entities with higher popularity, which might cause better disambiguation results.

¹⁷<http://gerbil.aksw.org/gerbil/experiment?id=201712060011>

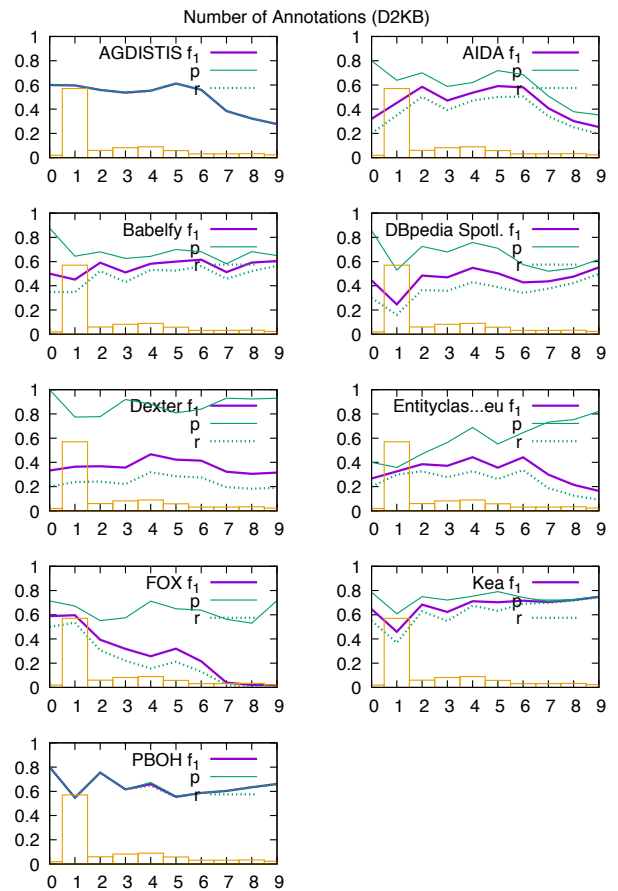


Fig. 22. Results for Number of Annotations (D2KB)

4.2.6. Number Of Annotations

Fig. 22 and 23 show the results for the number of annotations measure. This measure is not to be interpreted as a quality of the annotations but of the documents. Tab. 4 shows that more than half (595) of the 1,043 documents contain exactly 3 annotations, indicated by partition 1. Only 20 documents contain fewer annotations (partition 0). The number of annotations also corresponds to the size of the ‘disambiguation context’.

For this measure both experiment types D2KB¹⁸ (Fig. 22) and A2KB¹⁹ (Fig. 23) were conducted. For the A2KB task, the AGDISTIS system was not available, because it is only capable of D2KB tasks. For the period of D2KB experiments also the PBOH system was available. Entityclassifier.eu produced several er-

¹⁸<http://gerbil.aksw.org/gerbil/experiment?id=201711280011>

¹⁹<http://gerbil.aksw.org/gerbil/experiment?id=201711280030>

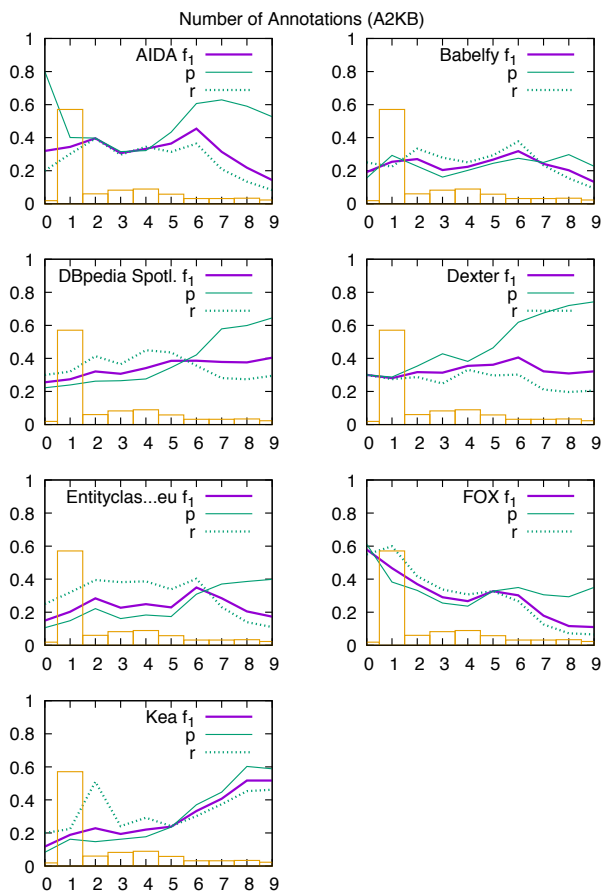


Fig. 23. Results for Number of Annotations (A2KB)

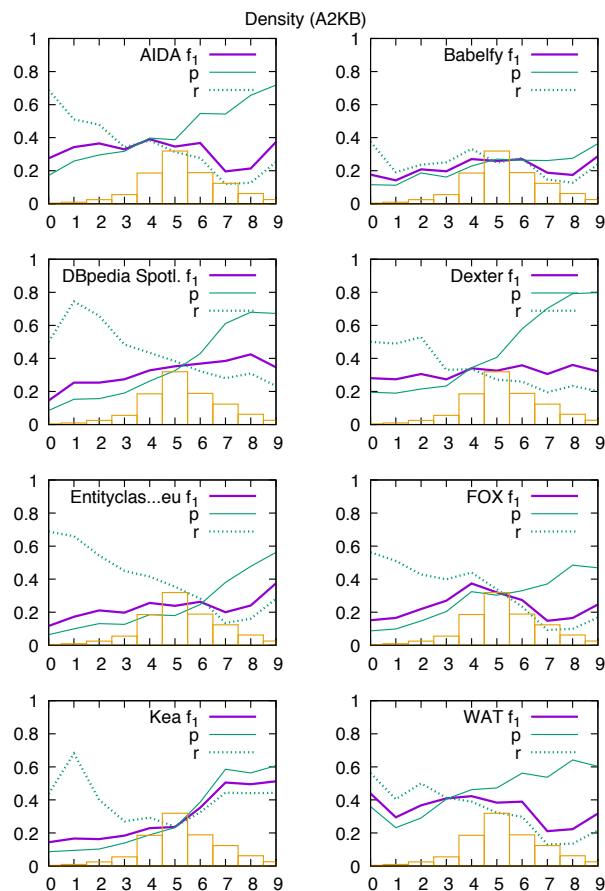


Fig. 24. Results for Density (A2KB)

rors, but overall, the results seem to be valid. WAT was not available.

In Fig. 22 (D2KB) it can be observed that some systems are not robust against growing context size, as e. g., AGDISTIS, AIDA, Entityclassifier.eu, and FOX. The other systems exhibit a more or less constant behaviour. The annotation tasks (A2KB) presented in Fig. 22 confirm this observation. Almost every system increases precision with growing context sizes, but on the expense of recall. This drifting apart occurs between the 4th and 6th partition (16 to 50 annotations per document). KEA seems to strongly benefit from increasing context sizes, while FOX benefits from smaller context sizes.

4.2.7. Density

The results for the density measure are presented in Fig. 24. Density also is a quality of the documents and not of their annotations. Low density (left hand partitions) signifies that a longer document has only a few annotations. High density (right hand partitions) on the

other hand signifies that a document contains many annotations relative to its length.

For density the experiments were conducted as A2KB tasks²⁰. All participating systems except PBOH provided valid results.

From the presented graphs it can be observed that the systems perform on low dense documents with high recall, but comparably low precision. On the other hand, dense documents are annotated with higher precision, but lower recall. While Babelify performs more or less evenly distributed, KEA seems to also maintain recall with denser documents. The break even point between precision and recall is located between the 4th and 6th partition (density between 0.055 and 0.133).

The density only estimates the number of missing annotations and the correlation between this metric and precision and recall supports this to some extent.

²⁰<http://gerbil.aks.org/gerbil/experiment?id=201712050010>

However, it is important to also take into account the reasons for sparsity. Sparsity in the annotations can also stem from a specific combination of a knowledge base and documents. Very domain specific documents with little coverage in the knowledge base will often be sparsely annotated, even if the annotation is complete with respect to the knowledge base. This limits this metric's utility. It would be interesting to assess whether the density can be put in relation to the dominance of entities and surfaces forms in order to reduce domain and knowledge base dependencies.

4.2.8. General Results

Table 5 shows the achieved micro- f_1 results of the systems for the D2KB task. The top row indicates the original GERBIL results²¹ (No Filter). Top results are indicated in green (bold) and the lowest results in red (italic). Each row shows the results for the dataset filtered according to a specific criteria. The second column shows the number of remaining annotations in the dataset after filtering. The penultimate column shows the average of the systems, the last column the Pearson correlation of the current row to the first row. Unfortunately the WAT system did not produce usable results and had to be excluded.

For persons²², organizations²³ and places²⁴ the results achieved by the systems are rather similar, but do not perfectly correlate to the baseline (first row). For persons and organizations PBOH seems to be the best system. KEA produces the best results for places and for the entities not falling into these categories (others). The others category strongly correlates with the baseline.

The next 2 rows separate annotations into a dataset containing entities with `itsrdf:taClassRef` statement (with Classes²⁵) and without (without Classes²⁶). The first dataset correlates very strongly to the baseline. For the annotations without class assignment the correlation is not so clear, furthermore the annotation performance was comparably low.

²¹<http://gerbil.aksw.org/gerbil/experiment?id=201711230013>

²²<http://gerbil.aksw.org/gerbil/experiment?id=201711280013>

²³<http://gerbil.aksw.org/gerbil/experiment?id=2017112800143>

²⁴<http://gerbil.aksw.org/gerbil/experiment?id=201711280015>

²⁵<http://gerbil.aksw.org/gerbil/experiment?id=201711280028>

²⁶<http://gerbil.aksw.org/gerbil/experiment?id=201711280020>

Another filtering was performed by filtering entities according to class membership of typical classes of the three different domains: Music²⁷, Science²⁸, and Movie/TV²⁹. In every domain a different system performed best. The Pearson value for Music indicates a lower correlation.

The last four rows show datasets filtered according to thresholds of the proposed measures. For the first, we removed the first and last decile partition to avoid bias caused by disambiguation and redirect resources, too popular and unpopular entities, entities without information about PageRank and HITS, extremely short and large contexts, extreme homonyms and synonyms (likelihood of confusion). Furthermore, the density was restricted to a moderate level around the break even points between precision and recall to avoid major bias caused by extreme strong and low density. The filtered dataset is denoted as the 'low skew' dataset³⁰. The dataset contains 765 annotations in 118 documents. Considering Tab. 4, a grey cell background indicates that this partition was *not* included in the 'low skew' dataset.

From all these restrictions, all annotations have been filtered, which fall into the intersection of the opposite filters, denoted as the 'high skew' dataset³¹ (grey cells of Tab. 4). This results in only 66 annotations in 22 documents.

Tab. 5 shows that the results for the 'low skew' dataset are overall better than for the 'high skew' dataset. But surprisingly, 3 systems (KEA, AGDISTIS, Dexter) perform with larger f-measure than on the 'low skew' dataset. With a larger value of 0.898 the Pearson value suggests a slightly better correlation with the baseline for the 'low skew' dataset than for the 'high skew' dataset with 0.866.

For the 'high skew' dataset, 66 annotations might not be very representative, but applying all the restrictions resulted in this rather small dataset. To increase the size we attempted to relax the restrictions

²⁷<http://gerbil.aksw.org/gerbil/experiment?id=201712060008> <http://gerbil.aksw.org/gerbil/experiment?id=201712110000>

²⁸<http://gerbil.aksw.org/gerbil/experiment?id=201712060009> <http://gerbil.aksw.org/gerbil/experiment?id=201712110001>

²⁹<http://gerbil.aksw.org/gerbil/experiment?id=201712060007>

³⁰<http://gerbil.aksw.org/gerbil/experiment?id=201712100002>

³¹<http://gerbil.aksw.org/gerbil/experiment?id=201712100003>

slightly and created the ‘medium skew’ dataset³². For this dataset the filters have not been applied all at once. According to a ‘leave-one-out’ principle, for each of the 6 filters (see header of Tab. 4) a datasets was created. The dataset for a particular filter was restricted only to the *other* filters. Finally, a join has been applied to the sets resulting in 595 annotations. The results presented in Tab. 5 show that the systems performed in a similar manner compared to the ‘high skew’ results. Four systems performed slightly better, the results of three diminished. Unfortunately we were not able to produce results for Dexter and Entityclassifier.eu. Thus, the correlation quotient is not informative. In general, the results of both ‘high skew’ datasets should be treated with caution since they are created on purpose from outliers and very likely contain bias. The consequence is, that the results are not trustful, and a system performing well on the ‘high skew’ datasets, e.g. KEA, must not necessarily perform well overall. However, we can see, that PBOH performs best on the ‘low skew’ dataset, and this seems to be an objective and reliable result.

The last two remixed datasets are derived from the ‘low skew’ dataset. The first one was compiled with the intent to include only annotations, which are comparably ‘easy’ to disambiguate³³. The other one includes annotations which are considered more ‘difficult’ to resolve³⁴. Considering Tab. 4 the green, orange, and white partitions belong to the easy dataset, the red, orange and white partitions belong to the difficult datasets. Thus, the easy dataset preferably contains annotations with more popular entities and lower likelihood of confusion of entities and surface forms. For the difficult dataset annotations with unpopular entities and higher likelihood of confusion of entities and surface forms are considered. We did not further restrict the number of annotations and density values compared to the ‘low skew’ dataset, because the resulting datasets would have been too small.

KEA performed well on the dataset that was considered easier, but not on the difficult dataset where PBOH is ahead of all other systems. The average numbers of the easy and difficult datasets suggest that ex-

pectations have been fulfilled. The dataset considered more difficult to solve in fact is more difficult to solve and the easy dataset easier to solve than others. The results for the difficult dataset only slightly correlate with the overall results.

4.2.9. Measures of remixed datasets

For a further detailed view on the data, the characteristics of the remixed datasets have been calculated and are presented in Fig. 25, 26, and 27.

Fig. 25 shows the density values of the remixed datasets. Since the datasets are filtered on annotation level and therewith some annotations were not included, it is to be expected that the density values are overall smaller compared to the unfiltered datasets (see Fig. 13). For the experiments conducted as D2KB tasks, the density does not influence the results. For A2KB tasks it might be more useful to remix on document level instead of annotation level.

In Fig. 26 the likelihood of confusions are presented. As expected, the difficult dataset contains a larger average number of entities per surface form, indicating more homonyms compared to the easy dataset. Furthermore, the number of surface forms per entity is smaller for the difficult compared to the easy dataset, indicating a smaller number of synonyms. We might conclude, that items of the Science category are more difficult to disambiguate than items of the Place category. The ‘high skew’ category almost only contains one item per surface form, respectively entity. Revising the data revealed that with the filtering of this category (cf. Tab. 4 grey background cells) partition 9 for the likelihoods of confusion has been completely cleared out by the other restrictions (PageRank, HITS, etc.). Thus, it seems that there exist some dependencies between the measures.

In Fig. 27 the dominance of entities and surface forms is presented. In general, for the remixed datasets the dominance of entities is larger than for the original datasets. This is to be expected, because by filtering out annotations in the dataset (reducing S^D and E^D), the remaining entities are gaining more dominance.

4.2.10. Dataset coverage

To also observe the distribution of the origin datasets over the remixed datasets the following analysis is performed. Tab. 6 shows the coverage of the origin datasets (rows) and the remixed datasets (columns). The first data row shows the number of annotations in the remixed datasets. Column ‘Complete’ corresponds to the join of all origin datasets. The origin datasets are described row by row. Each origin dataset item

³²<http://gerbil.aksw.org/gerbil/experiment?id=201804090002>

³³<http://gerbil.aksw.org/gerbil/experiment?id=201712120003> <http://gerbil.aksw.org/gerbil/experiment?id=201804020001>

³⁴<http://gerbil.aksw.org/gerbil/experiment?id=201712120004>

Table 5
Micro-f₁ results of D2KB systems for different remixed datasets.

	IAI	Babelfy	Spotl.	Dexter	Ent.cl.	Fox	KEA	AGDI	AIDA	PBOH	AVG	Pearson
No Filter	16821	0.572	0.485	0.349	0.285	0.167	0.704	0.407	0.374	0.625	0.441	
Person	1556	0.830	0.407	0.506	0.505	0.268	0.795	0.645	0.756	0.839	0.617	0.779
Organization	1084	0.731	0.530	0.519	0.487	0.325	0.732	0.675	0.756	0.838	0.621	0.796
Places	1477	0.702	0.643	0.643	0.695	0.257	0.866	0.693	0.809	0.856	0.685	0.763
other	12931	0.512	0.467	0.265	0.164	0.113	0.651	0.333	0.259	0.561	0.369	0.987
with Classes	11306	0.658	0.560	0.410	0.342	0.129	0.807	0.406	0.425	0.742	0.498	0.992
without Classes	5515	0.381	0.324	0.212	0.168	0.235	0.467	0.413	0.277	0.385	0.318	0.829
Music	525	0.545	0.449	0.560	0.511	0.189	0.704	0.582	0.684	0.656	0.542	0.693
Science	225	0.797	0.574	0.364	0.259	0.136	0.778	0.307	0.451	0.756	0.491	0.953
Movie/TV	305	0.631	0.367	0.406	0.379	0.239	0.618	0.477	0.515	0.688	0.480	0.871
Low skew	765	0.617	0.614	0.327	0.428	0.144	0.646	0.361	0.500	0.694	0.481	0.898
High skew	66	0.517	0.234	0.489	0.208	0.029	0.760	0.364	0.415	0.621	0.404	0.866
Medium skew	595	0.567	0.297	err	err	0.183	0.681	0.366	0.344	0.603	0.435	0.932
Easy	235	0.716	0.769	0.647	0.654	0.625	0.811	0.566	0.630	0.809	0.692	0.795
Difficult	98	0.601	0.421	0.070	0.126	0.065	0.194	0.071	0.552	0.622	0.302	0.538

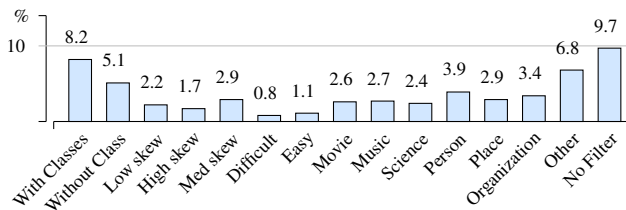


Fig. 25. Annotation density as relative number of annotations respective document length in words

contains three lines of numbers. The first line shows the number of annotations covered in the columnwise dataset, e.g. the KORE50 dataset contains 144 annotations, whereas 3 of them also belong to the ‘low skew’ dataset. The second line just shows the relative numbers, e.g. the KORE50 dataset contributes 0.86% to the ‘Complete’ set of annotations and 0.39% to the ‘low skew’ dataset. The third line relates the number of annotations of the column to the size of the origin dataset, meaning that e.g. 2.08% of the KORE50 dataset also belong to the ‘low skew’ dataset. Special aspects are highlighted through bold font.

It is observable, that the IITB dataset contributes almost two thirds to the entire experiments, which also leads to a large coverage over the remixed datasets. 4.99% of its annotations fall into the ‘med skew’ category. IITB and KORE50 seemingly are the most ‘high skew’ datasets. But, the number of ‘high skew’ annota-

tions overall is considerably low, so that it can be said, that there is no origin dataset which might suffer too much skewness.

On the other side, we see MSNBC and AQUAINT as considerably low skewed datasets. Over 20% of their annotations fall in that category.

With 11.97% of annotations AQUAINT has the largest fraction of easy annotations. The dataset with the largest relative number of difficult annotations is MSNBC with 5.23%. Surprisingly, the KORE50 dataset does not contribute to the difficult dataset at all, which contradicts KORE50’s creation intention.

In summary, the share of ‘high skew’ elements overall is rather small. There is no dataset that should be excluded in further evaluation experiments because it is completely ‘out of order’.

5. Conclusion

In this paper an extension of the GERBIL framework has been introduced to enable a more fine grained evaluation of NEL systems.

According to the predefined entity types, the KORE-50 benchmark dataset contains the most persons, N3-Reuters-500 the most organizations, and ACE2004 the most places. The IITB dataset on the other hand contains almost no persons, organizations, or places. According to the PageRank algorithm the DBpedia

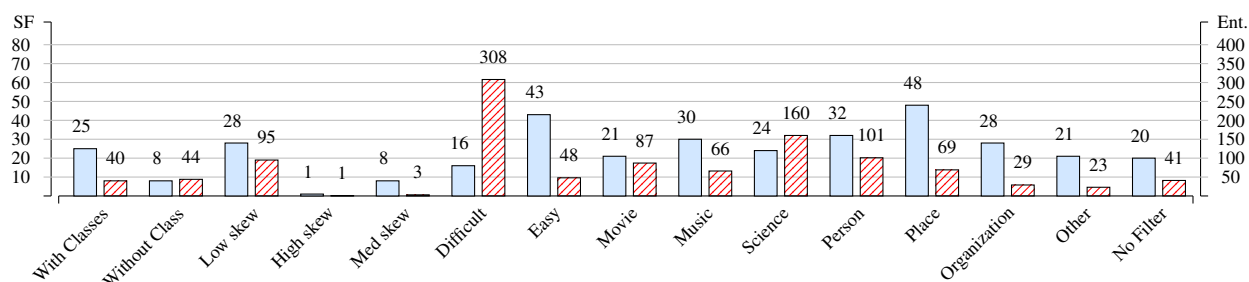


Fig. 26. Average number of surface forms (SF) per entity (blue, left) and average number of entities per surface form (red/hatched, right) indicating the likelihood of confusion for each dataset

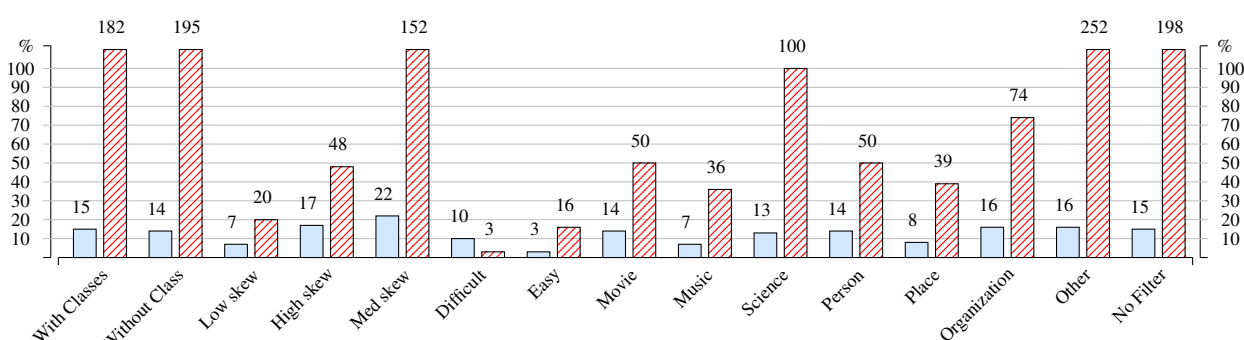


Fig. 27. Average dominance for surface forms (blue) and entities (red/hatched) per dataset

Spotlight dataset contains the most prominent entities, while the Micropost 2014 Test dataset contains the most entities with medium and low prominence. N3-RSS contains the fewest popular and OKE 2015 gold standard the fewest medium and low prominence entities. The HITS value showed a more diverse picture with Micropost 2014 Train containing the most popular entities, MSNBC with the most medium prominence entities, and WES2015 with the most low prominence entities. On the other hand, IITB contains the fewest high prominence entities and OKE 2015 gold standard follows with the fewest medium prominence entities. N3-RSS-500 contains the fewest low prominence entities.

A stand-alone library has been introduced to enrich documents encoded in the NIF format with additional meta information. This enables researchers to remix existing NIF-based datasets according to their needs in a reproducible manner.

An exhaustive example was presented, on how to use the library to reorganize datasets according to the measures introduced earlier. Therefore, datasets were combined and partitioned to determine and visualize for each system correlations between a dataset prop-

erty and the system's performance. It was ascertained that systems fail with homonyms with a likelihood of confusion beyond ca. 1,700 entities mapping to the surface form. From the analysis on entities' likelihood of confusions, it was confirmed that redirect and disambiguation resources strongly bias the overall results. However, the overall performance increases the more surface forms an entity is mapping to. It was also shown that the PageRank of entities correlates with the systems performance, but only up to a certain threshold. Interestingly, for the HITS measure the systems produced poor results on low to medium, but very good results on very low and larger values. It was further shown that not all systems are robust against a rising number of annotations in a text to disambiguate. Many systems tend to suffer loss of recall with larger numbers of items to disambiguate. While FOX greatly performs on smaller contexts, KEA benefits from larger numbers of annotations in a context. Finally, the density measure shows that text with rather few annotations can promote recall and demote precision very unevenly.

Furthermore, an overall comparison of different filtered datasets was given including a focus on specific

Table 6
Coverage of origin datasets and remixed datasets

	Complete	Low skew	High skew	Med skew	Easy	Difficult
	16821	765	66	595	235	98
DBp. Spotl.	330 1.96%	14 1.83%	0 0.00%	3 0.50%	4 1.70%	0 0.00%
		4.24%	0.00%	0.91%	1.21%	0.00%
KORE50	144 0.86%	3 0.39%	2 3.03%	7 1.18%	1 0.43%	0 0.00%
		2.08%	1.39%	4.86%	0.69%	0.00%
MSNBC	650 3.86%	137 17.91%	0 0.00%	5 0.84%	16 6.81%	34 34.69%
		21.08%	0.00%	0.77%	2.46%	5.23%
IITB	11182 66.48%	357 46.67%	63 95.45%	558 93.78%	100 42.55%	42 42.86%
		3.19%	0.56%	4.99%	0.89%	0.38%
N3-RSS500	1000 5.94%	0 0.00%	1 1.52%	12 2.02%	0 0.00%	0 0.00%
		0.00%	0.10%	1.20%	0.00%	0.00%
N3-Reuters-128	880 5.23%	89 11.63%	0 0.00%	5 0.84%	26 11.06%	11 11.22%
		10.11%	0.00%	0.57%	2.95%	1.25%
ACE2004	253 1.50%	3 0.39%	0 0.00%	4 0.67%	1 0.43%	0 0.00%
		1.19%	0.00%	1.58%	0.40%	0.00%
News-100	1655 9.84%	1 0.13%	0 0.00%	1 0.17%	0 0.00%	1 1.02%
		0.06%	0.00%	0.06%	0.00%	0.06%
AQUAINT	727 4.32%	161 21.05%	0 0.00%	0 0.00%	87 37.02%	10 10.20%
		22.15%	0.00%	0.00%	11.97%	1.38%

domains, as e. g., persons, organizations, places, music, science, movies/tv. Although KEA and PBOH perform well in the majority of cases, they are not necessarily the best performing systems. Babelfy greatly performs on the science domain, thus, there are domain and dataset structure specific preferences across the systems. Therefore, it is of major importance always to take into account the characteristics of datasets for entity linking benchmarks.

It is impossible to define how a perfect ‘one for all’ dataset should look like. However, we attempted to compile at least one dataset that is almost free of the apparent biasing factors ascertained from the proposed measures. To determine the ‘difficulty’ of a dataset, the confusion and popularity measures seem to be appropriate measures, but only in combination with moderate size of context and balanced density. Extreme out-

liers should be avoided as possible. Also redirect and disambiguation resources distort the result very much.

From the remixing we have learned, that there are in fact domain differences in the performance of the annotating systems. The systems have their peculiarities according to the introduced measures and there are differences in the quality of datasets. But, we cannot find evidence, that the datasets under consideration contain a harmful number of inappropriate annotations.

Further biasing factors identified in the datasets are NIL (notInWiki) annotations and the mixture of language versions of DBpedia, as for example caused by including the News-100 dataset. Both should be taken into account in further versions of this work. Unfortunately, the applied online annotation systems were not always available. Moreover, it is not clear what the current development state of the systems is or how many systems exist that are not connected to GERBIL, which

might also worthwhile to be included in further analysis.

Ongoing research is focused on the implementation of additional measures, such as e. g. those introduced by [10,24] and the annotation systems' performance breakdown should also include the dominance and maximum recall measures. More datasets such as WES2015 and the Microposts series should be included in future versions.

Also, we would like to introduce difficulty levels for datasets along with new properties for annotation, which might be useful for further remixing, as e. g. a distinction of the NEL annotation for common and proper nouns, or the dependency on temporal context. The inter-systems agreement might also be a valuable measure to be included into an evaluation.

The results of this work as well as the provided source code and the public online service enable to improve further benchmarks, to optimize systems for a unprecedented level of detail, and the results enable to find the right tool or method for the desired annotation task.

In summary, the evaluation at a finer granular level allows a better understanding of the NEL process and also promotes the development of improved NEL systems.

Appendix

A. Mathematical notation

D	Dataset, a set of annotated documents
$ D $	Number of documents in D
$d = (d_t, d_a)$	A document $d \in D$
d_t	Text of document d
$ d_t $	Number of words within the text of document d
d_a	Annotations in document d
$ d_a $	Number of annotations in d
e	An entity $e \in E$
s	A surface form $s \in S$
$a = (s, e, i, l)$	An annotation with surface form s , entity e , text index i , and length l
E	A set of entities
E^D	Set of entities in dataset D
E^d	Set of entities in a document d
$E(s)$	Set of entities for the surface form s
W_E	A mapping (dictionary) from surface forms to entities $W_E : S \rightarrow E$ of an annotation system
$W_E(s)$	Set of entities in dictionary W_E for surface form s
S	Set of surface forms
S^D	Set of surface forms in dataset D
S^d	Set of surface forms in document d
$S(e)$	Set of surface forms for the entity e
W_S	A mapping (dictionary) from entities to surface forms $W_S : E \rightarrow S$ of an annotation system
$W_S(e)$	Set of surface forms in the dictionary W_S for entity e
P	Arbitrary scoring algorithm (e.g. PageRank, HITS) to estimate popularity

B. Formula overview

Average number of annotations:

$$na(D) = \frac{\sum_{d \in D} |d_a|}{|D|} \quad (11)$$

Average number of not annotated documents:

$$nad(D) = \frac{|\{d : |d_a| = 0\}|}{|D|} \quad (12)$$

Density of a document d :

$$density(d) = \frac{|d_a|}{|d_t|} \quad (13)$$

Density of dataset D :

$$density_{micro}(D) = \frac{\sum_{d \in D} density(d)}{|D|} \quad (14)$$

$$density_{macro}(D) = \frac{\sum_{d \in D} |d_a|}{\sum_{d \in D} |d_t|}$$

Set of entities with prominence in interval $[a, b]$ for a scoring algorithm P :

$$E_{a,b}^D(P) = \{e \in E^D : a \leq P(e) \leq b\} \quad (15)$$

Average likelihood of confusion for all surface forms of dataset D

$$lc_{micro}^{sf}(D, W) = \frac{\sum_{d \in D} \frac{\sum_{a \in d_a} |W_E(s) \cup E^D(s)|}{|d_a|}}{|D|} \quad (16)$$

$$lc_{macro}^{sf}(D, W) = \frac{\sum_{s \in S^D} |W_E(s) \cup E^D(s)|}{|S^D|}$$

Average likelihood of confusion for all entities of dataset D :

$$lc_{micro}^e(D, W) = \frac{\sum_{d \in D} \frac{\sum_{a \in d_a} |W_S(e) \cup S^D(e)|}{|d_a|}}{|D|} \quad (17)$$

$$lc_{macro}^e(D, W) = \frac{\sum_{e \in E^D} |W_S(e) \cup S^D(e)|}{|E^D|}$$

Dominance of surface forms:

$$dom_{micro}^{sf}(D, W) = \frac{\sum_{d \in D} \frac{\sum_{a \in d_a} \frac{E^d(s)}{|W_E(s)|}}{|d_a|}}{|D|} \quad (18)$$

$$dom_{macro}^{sf}(D, W) = \frac{\sum_{s \in S^D} \frac{|E^D(s)|}{|W_E(s)|}}{|S^D|}$$

Dominance of entities :

$$dom_{micro}^e(D, W) = \frac{\sum_{d \in D} \frac{\sum_{a \in d_a} \frac{S^d(e)}{|W_S(e)|}}{|d_a|}}{|D|} \quad (19)$$

$$dom_{macro}^e(D, W) = \frac{\sum_{e \in E^D} \frac{|S^D(e)|}{|W_S(e)|}}{|E^D|}$$

Maximum recall:

$$mr_{micro}(S, W) = \frac{\sum_{d \in D} (1 - \frac{|S^d \setminus W_S|}{|S^d|})}{|D|} \quad (20)$$

$$mr_{macro}(S, W) = 1 - \frac{|S^D \setminus W_S|}{|S^D|}$$

Set of entities in dataset D with type T where E^T is the set of all entities with type T :

$$E^D(T) = \{e \in E^D : e \in E^T\} \quad (21)$$

References

- [1] S. Bhatia and A. Jain. Context Sensitive Entity Linking of Search Queries in Enterprise Knowledge Graphs. In *Proceedings of the International Semantic Web Conference 2016 (ESWC2016)*, pages 50–54. Springer, Cham, 2016. https://doi.org/10.1007/978-3-319-47602-5_11.
- [2] A. E. Cano, G. Rizzo, A. Varga, M. Rowe, M. Stankovic, and A.-S. Dadzie. Making Sense of Microposts:(# microposts2014) Named Entity Extraction & Linking Challenge. In *CEUR Workshop Proceedings*, volume 1141, pages 54–60, 2014.
- [3] D. Ceccarelli, C. Lucchese, S. Orlando, R. Perego, and S. Trani. Dexter: an Open Source Framework for Entity Linking. In *Proceedings of the 6th International Workshop on Exploiting Semantic Annotations in Information Retrieval*, pages 17–20. ACM, 2013, <https://doi.org/10.1145/2513204.2513212>.
- [4] M. Cornolti, P. Ferragina, and M. Ciaramita. A framework for benchmarking entity-annotation systems. In *Proceedings of the 22nd International Conference on World Wide Web (WWW2013)*, pages 249–260. ACM, 2013, <https://doi.org/10.1145/2488388.2488411>.
- [5] S. Cucerzan. Large-Scale Named Entity Disambiguation Based on Wikipedia Data. In *Proceedings of the 2007 Joint Conference on Empirical Methods in Natural Language Processing and Computational Natural Language Learning (EMNLP-CoNLL)*, pages 708–716. ACL, 2007, <http://www.aclweb.org/anthology/D/D07/D07-1074>.
- [6] M. Dojchinovski and T. Kliegr. Entityclassifier.eu: Real-time Classification of Entities in Text with Wikipedia. In *Joint Eu-*

- ropean Conference on Machine Learning and Knowledge Discovery in Databases, volume 8190 of *LNCS*, pages 654–658. Springer, Berlin, Heidelberg, 2013, https://doi.org/10.1007/978-3-642-40994-3_48.
- [7] M. Dragoni, E. Cabrio, S. Tonelli, and S. Villata. Enriching a Small Artwork Collection Through Semantic Linking. In *Proceedings of the International Semantic Web Conference (ESWC 2016)*, volume 9678 of *LNCS*, pages 724–740. Springer, Cham, 2016, https://doi.org/10.1007/978-3-319-34129-3_44.
- [8] F. Frontini, C. Brando, and J.-G. Ganascia. Semantic Web Based Named Entity Linking for Digital Humanities and Heritage Texts. In *1st International Workshop Semantic Web for Scientific Heritage at the 12th ESWC 2015 Conference*, volume 1364 of *CEUR-WS*, pages 77–88, Portorož, Slovenia, June 2015.
- [9] O.-E. Ganea, M. Ganea, A. Lucchi, C. Eickhoff, and T. Hofmann. Probabilistic Bag-Of-Hyperlinks Model for Entity Linking. In *Proceedings of the 25th International Conference on World Wide Web (WWW '16)*, pages 927–938. ACM, 2016, <https://doi.org/10.1145/2872427.2882988>.
- [10] B. Hachey, J. Nothman, and W. Radford. Cheap and easy entity evaluation. In *Proceedings of the 52nd Annual Meeting of the Association for Computational Linguistics (Short Papers)*, pages 464–469. ACL, 2014.
- [11] S. Hellmann, J. Lehmann, S. Auer, and M. Brümmer. Integrating NLP using linked data. In *International Semantic Web Conference (ISWC2013)*, volume 8219 of *LNCS*, pages 98–113. Springer, Berlin, Heidelberg, 2013, https://doi.org/10.1007/978-3-642-41338-4_7.
- [12] J. Hoffart, S. Seufert, D. B. Nguyen, M. Theobald, and G. Weikum. KORE: Keyphrase Overlap Relatedness for Entity Disambiguation. In *21st ACM International Conference on Information and Knowledge Management*, pages 545–554. ACM, New York, NY, USA, 2012, <https://doi.org/10.1145/2396761.2396832>.
- [13] J. Hoffart, M. A. Yosef, I. Bordino, H. Fürstenau, M. Pinkal, M. Spaniol, B. Taneva, S. Thater, and G. Weikum. Robust Disambiguation of Named Entities in Text. In *Proceedings of the Conference on Empirical Methods in Natural Language Processing*, pages 782–792. ACL, 2011.
- [14] S. Kulkarni, A. Singh, G. Ramakrishnan, and S. Chakrabarti. Collective Annotation of Wikipedia Entities in Web Text. In *Proceedings of the 15th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pages 457–466. ACM, New York, NY, USA, 2009, <https://doi.org/10.1145/1557019.1557073>.
- [15] X. Ling, S. Singh, and D. S. Weld. Design Challenges for Entity Linking. *Transactions of the Association for Computational Linguistics*, 3:315–328, 2015.
- [16] J. L. Martinez-Rodriguez, A. Hogan, and I. Lopez-Arevalo. Information Extraction meets the Semantic Web: A Survey. *Semantic Web*, (accepted 2018; to appear) <http://www.semantic-web-journal.net/system/files/swj1909.pdf>.
- [17] P. N. Mendes, M. Jakob, A. García-Silva, and C. Bizer. Dbpedia spotlight: Shedding Light on the Web of Documents. In *Proceedings of the 7th International Conference on Semantic Systems (I-Semantic 2011)*, pages 1–8. ACM, New York, NY, USA, 2011, <https://doi.org/10.1145/2063518.2063519>.
- [18] D. Milne and I. H. Witten. Learning to Link with Wikipedia. In *Proceedings of the 17th ACM International Conference on Information and Knowledge Management*, pages 509–518. ACM, New York, NY, USA, 2008, <https://doi.org/10.1145/1458082.1458150>.
- [19] A. Mitchell, S. Strassel, S. Huang, and R. Zakhary. ACE 2004 Multilingual Training Corpus LDC2005T09. Web download, Linguistic Data Consortium, Philadelphia, 2005, <https://catalog.ldc.upenn.edu/ldc2005t09>.
- [20] A. Moro, A. Raganato, and R. Navigli. Entity Linking meets Word Sense Disambiguation: a Unified Approach. *Transactions of the Association for Computational Linguistics*, 2:231–244, 2014.
- [21] A. G. Nuzzolese, A. L. Gentile, V. Presutti, A. Gangemi, D. Garigliotti, and R. Navigli. Open Knowledge Extraction Challenge. In *Semantic Web Evaluation Challenge*, volume 548 of *CCIS*, pages 3–15. Springer, Cham, 2015, https://doi.org/10.1007/978-3-319-25518-7_1.
- [22] L. Page, S. Brin, R. Motwani, and T. Winograd. The PageRank Citation Ranking: Bringing Order to the Web. *Stanford InfoLab*, 1999.
- [23] F. Piccinno and P. Ferragina. From TagME to WAT: a new Entity Annotator. In *Proceedings of the 1st International Workshop on Entity Recognition & Disambiguation (ERD2014)*, pages 55–62. ACM, New York, NY, USA, 2014, <https://doi.org/10.1145/2633211.2634350>.
- [24] S. Pradhan, X. L. an Marta Recasens, E. H. Hovy, V. Ng, and M. Strube. Scoring Coreference Partitions of Predicted Mentions: A Reference Implementation. In *52nd Annual Meeting of the Association for Computational Linguistics*, pages 30–35. ACL, 2014, <https://doi.org/10.3115/v1/P14-2006>.
- [25] D. Reddy, M. Knuth, and H. Sack. DBpedia GraphMeasures. Hasso Plattner Institute, Potsdam, July 2014, <http://s16a.org/node/6>.
- [26] G. Rizzo, A. E. C. Basave, B. Pereira, and A. Varga. Making Sense of Microposts (#microposts2015) Named Entity Recognition and Linking (NEEL) Challenge. In *5th Workshop on Making Sense of Microposts at 24th Int. World Wide Web Conference*, volume 1395 of *CEUR-WS*, pages 44–53, 2015.
- [27] G. Rizzo and R. Troncy. NERD: A framework for unifying named entity recognition and disambiguation extraction tools. In *Proceedings of the Demonstrations at the 13th Conference of the European Chapter of the Association for Computational Linguistics*, pages 73–76. ACL, Stroudsburg, PA, USA, 2012.
- [28] G. Rizzo, M. van Erp, and R. Troncy. Benchmarking the Extraction and Disambiguation of Named Entities on the Semantic Web. In *9th Int. Conf. on Language Resources and Evaluation (LREC2014)*. European Language Resources Association (ELRA), 2014.
- [29] M. Röder, R. Usbeck, S. Hellmann, D. Gerber, and A. Both. N³-A Collection of Datasets for Named Entity Recognition and Disambiguation in the NLP Interchange Format. In *Proceedings of the 9th International Conference on Language Resources and Evaluation (LREC2014)*. European Language Resources Association (ELRA), 2014.
- [30] M. Röder, R. Usbeck, and A. N. Ngomo. GERBIL - benchmarking named entity recognition and linking consistently. *Semantic Web*, 9(5):605–625, 2018, <https://doi.org/10.3233/SW-170286>.
- [31] W. Shen, J. Wang, and J. Han. Entity Linking with a Knowl-

- edge Base: Issues, Techniques, and Solutions. *IEEE Transactions on Knowledge and Data Engineering*, 27(2):443–460, Feb 2015, <https://doi.org/10.1109/TKDE.2014.2327028>.
- [32] A. Singhal. Introducing the knowledge graph: things, not strings. *Official Google Blog*, May, 2012.
- [33] R. Speck and A.-C. N. Ngomo. Named Entity Recognition Using FOX. In *Proceedings of the ISWC 2014 Posters & Demonstrations Track within the 13th International Semantic Web Conference (ISWC 2014)*, volume 1272 of *CEUR-WS*, pages 85–88, 2014.
- [34] N. Steinmetz, M. Knuth, and H. Sack. Statistical Analyses of Named Entity Disambiguation Benchmarks. In *Proceedings of NLP & DBpedia 2013 workshop at 12th International Semantic Web Conference*, volume 1064 of *CEUR-WS*, 2013.
- [35] T. Tietz, J. Waitelonis, J. Jäger, and H. Sack. Smart Media Navigator: Visualizing Recommendations based on Linked Data. In *Proceedings of the Industry Track at the International Semantic Web Conference 2014 Co-located with the 13th International Semantic Web Conference (ISWC 2014)*, volume 1383 of *CEUR-WS*, 2014.
- [36] R. Usbeck, A.-C. N. Ngomo, M. Röder, D. Gerber, S. A. Coelho, S. Auer, and A. Both. AGDISTIS - Graph-Based Disambiguation of Named Entities using Linked Data. In *Proceedings of the 2014 International Semantic Web Conference (ISWC2014)*, volume 8796 of *LNCS*, pages 457–471. Springer, Cham, 2014, https://doi.org/10.1007/978-3-319-11964-9_29.
- [37] R. Usbeck, M. Röder, A.-C. Ngonga Ngomo, C. Baron, A. Both, M. Brümmer, D. Ceccarelli, M. Cornolti, D. Cherix, B. Eickmann, P. Ferragina, C. Lemke, A. Moro, R. Navigli, F. Piccinno, G. Rizzo, H. Sack, R. Speck, R. Troncy, J. Waitelonis, and L. Wesemann. GERBIL – General Entity Annotation Benchmark Framework. In *Proceedings of the 24th International Conference on World Wide Web (WWW2015)*, pages 1133–1143. International World Wide Web Conferences Steering Committee, 2015, <https://doi.org/10.1145/2736277.2741626>.
- [38] M. van Erp, P. Mendes, H. Paulheim, F. Ilievski, J. Plu, G. Rizzo, and J. Waitelonis. Evaluating Entity Linking: An Analysis of Current Benchmark Datasets and a Roadmap for doing a better Job. In *Proceedings of the 10th International Conference on Language Resources and Evaluation (LREC 2016)*, Paris, France, May 2016. ELRA.
- [39] J. Waitelonis, C. Exeler, and H. Sack. Linked Data Enabled Generalized Vector Space Model to Improve Document Retrieval. In *Proceedings of the Third NLP&DBpedia Workshop (NLP & DBpedia 2015) co-located with the 14th International Semantic Web Conference 2015 (ISWC 2015)*, volume 1581 of *CEUR-WS*, pages 33–44, 2015.
- [40] J. Waitelonis, H. Jürges, and H. Sack. Don’t Compare Apples to Oranges: Extending GERBIL for a Fine Grained NEL Evaluation. In *Proceedings of the 12th International Conference on Semantic Systems (SEMANTiCS 2016)*, pages 65–72. ACM, New York, NY, USA, 2016, <https://doi.org/10.1145/2993318.2993334>.
- [41] J. Waitelonis, M. Plank, and H. Sack. TIBIAV-Portal: Integrating Automatically Generated Video Annotations into the Web of Data. In *International Conference on Theory and Practice of Digital Libraries (TPDL2016)*, volume 9819 of *LNCS*, pages 429–433. Springer, Cham, 2016, https://doi.org/10.1007/978-3-319-43997-6_37.
- [42] J. Waitelonis and H. Sack. Named Entity Linking in #Tweets with KEA. In *Proceedings of the 6th Workshop on ‘Making Sense of Microposts’ co-located with the 25th International World Wide Web Conference (WWW 2016)*, volume 1691 of *CEUR-WS*, pages 61–63, 2016.
- [43] J. G. Zheng, D. Howsmon, B. Zhang, J. Hahn, D. McGuinness, J. Hendler, and H. Ji. Entity linking for biomedical literature. *BMC Medical Informatics and Decision Making*, 15(1):S4, May 2015, <https://doi.org/10.1186/1472-6947-15-S1-S4>.