

# Remixing Entity Linking Evaluation Datasets for Focused Benchmarking

Jörg Waitelonis<sup>a</sup>, Henrik Jürges<sup>b</sup> and Harald Sack<sup>c</sup>

<sup>a</sup> *Hasso-Plattner-Institute, Prof.-Dr.-Helmert-Str. 2-3, 14482 Potsdam, Germany*

*E-mail: joerg.waitelonis@hpi.de*

<sup>b</sup> *University of Potsdam, Am Neuen Palais 10, 14469 Potsdam, Germany*

*E-mail: juerges@uni-potsdam.de*

<sup>c</sup> *FIZ Karlsruhe, Leibniz Institute for Information Infrastructure, Hermann-von-Helmholtz-Platz 1, 76344*

*Eggenstein-Leopoldshafen, Germany*

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

**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 [17] 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. The implemented system as well as an adapted result visualization has been integrated in the publicly available GERBIL framework. In addition, a new system library to enrich provided NIF [3] datasets with statistical information including best practices for dataset remixing are 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 [14] or to improve search rankings [19]. 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 or relatedness measures

[16]. Likewise, social media and web monitoring systems benefit from NEL, for e. g. by the identification of persons or companies in social media content as subject of observation or tracking. A general survey on NEL systems is given by Chen et al. [13].

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 via machine learning [9]. Most NEL approaches usually follow a generic solution strategy, but there is an uprising trend for many systems to focus on the solution of rather specific tasks only, e. g. by the restriction to a specific domain of interest, document-, or entity type. This ongoing fragmentation of types

of tasks aggravates the application of generic benchmarking frameworks for NEL optimization and comparison such as GERBIL [17,12] or NERD [11,10]. With GERBIL, a NEL tool optimized for the detection of person names only is rather difficult to compare to other NEL tools with a more general focus. 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 are not comparable since the assumed person-only NEL annotator 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 14 annotators and 17 datasets, whereas these numbers are subject to constant change. For a detailed overview on annotators and datasets provided by GERBIL we refer to the official version<sup>1</sup>.

Besides the already described problem, there are 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.

A first attempt was made by Hoffart et al. [4] by manually compiling the Kore50<sup>2</sup> corpus aiming to capture hard to disambiguate mentions of entities. Another problem arises with the quality of annotations as described by [5] and [18] including e. g. annotation redundancy, inter-annotation agreement, topicality according to the evolving knowledge-bases, mention boundaries and nested annotations. Especially, completeness and coverage of annotations are essential measures to assess the annotation tasks (A2KB cf. [17]) where the entity mention detection contributes to the overall results.

Since no 'all-in-one' perfect data-set has emerged in the past, which covers all the aspects sufficiently

---

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

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

well, it would be beneficial to measure and provide dataset characteristics on document level to subsequently allow a re-compilation of documents across different datasets according to predefined criteria into a customized corpus. E. g. for the already mentioned person-only annotator these measures would help to specifically select only those documents, which exhibit a significant amount of person annotations providing a specific level of 'difficulty'. Remixing evaluation datasets on 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 deep analysis of the deployed benchmark datasets according to different emphases [20]. 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, the work presented in [20] is brought up-to-date and consolidated. Furthermore it is extended with new additional dataset measures, a stand-alone library to enable customized remixing of datasets, as well as a vocabulary to enrich NIF-based datasets with additional statistical information. Finally, we introduce best practices and examples to remix new datasets matching customizable criteria.

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 results we have achieved and how they can be exploited for the recompilation of customized new 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. Ling et al. [5] have introduced the basic characteristics of nine NEL datasets including the number of documents, number of mentions, entity types, number of NIL annotations. Steinmetz et al. [15] went one

step further with a more detailed view on the distribution of entity types including mapping coverage, entity candidate count, maximum recall, and entity popularity. Erp et al. [18] investigated on the overlap among datasets and introduced as 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 code<sup>3</sup> and also provides a publicly available on-line service<sup>4</sup>. Besides the implementation of filtering the benchmark datasets according to the desired characteristics, the system instantly updates and visualizes the per annotator results together with statistical summaries. The integration in 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 library<sup>5</sup> 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.

To enable a more differentiated NEL evaluation, the following characteristics are introduced with the purpose to perform analysis on dataset, document, as well as entity mention level.

### 2.1 Not Annotated Documents

Some of the available benchmark datasets also contain documents without any annotations at all. Documents without annotations lead to an increase of false positives in the evaluations and thereby cause a loss of precision. For a dataset  $\mathcal{D}$ , documents  $t \in \mathcal{D}$  and the set of annotations  $a(t)$  within  $t$ , the relative number of

documents without any annotation at all  $e: \mathcal{D} \rightarrow (0, 1)$  is determined as:

$$e(\mathcal{D}) = \frac{|\{t \in \mathcal{D} | a(t) = \emptyset\}|}{|\mathcal{D}|} \quad (1)$$

Empty documents are 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.

### 2.2 Missing Annotations (Density)

Similar to not annotated documents, missing annotations in an otherwise annotated document lead to a problem with the A2KB task. Annotators might 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 ground truth data, which requires major effort. However, we propose to estimate this number by measuring an annotation density value as the relation between the number of annotations in the ground truth  $a(t)$  and the overall document length  $len(t)$ , determined as the number of words, with  $d: \mathcal{D} \rightarrow (0, 1)$ :

$$d(\mathcal{D}) = \frac{\sum_{t \in \mathcal{D}} |a(t)|}{\sum_{t \in \mathcal{D}} len(t)} \quad (2)$$

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

### 2.3 Prominence (Popularity)

The assumption of [18] is, that evaluation against a corpus with a tendency to focus strongly on prominent or popular entities may cause problems. Hence, NEL systems preferring popular entities might exhibit an increase in performance. To verify this, we have implemented two different measures on entity level. Similarly to [18], the prominence is estimated as PageRank [6] 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 [8].

To evaluate annotators according to different levels of prominence of entities, the set of entities was partitioned as follows. A power-law distribution of the PageRank (respectively HITS) values over all entities is assumed, meaning that only a few entities exhibit a high PageRank and many entities a lower PageRank (long-tail), cf. Fig 1. Highly prominent entities are then

<sup>3</sup><https://github.com/santifa/gerbil/>

<sup>4</sup><http://gerbil.s16a.org/>

<sup>5</sup><https://github.com/santifa/hfts>

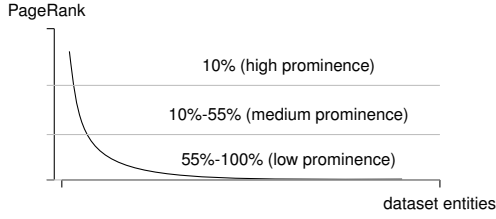


Fig. 1. Example partitioning for the PageRank.

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.

For a dataset, the relative amount of entities for every category is determined with  $p: (D, P) \rightarrow [0, 1]$  using the PageRank  $P \in \mathcal{P}$  and the interval  $a, b \in \mathbb{R}$  where  $e$  refers to a single entity:

$$p(D, P) = \frac{|\{e \in D | e \in P, a \leq e \leq b\}|}{|e \in D|} \quad (3)$$

The resulting set contains all entities of a dataset that satisfies the given interval limits. Similarly the prominence can be determined using the HITS values or any other ranking score.

#### 2.4 Likelihood of Confusion (Level of Ambiguity)

Since a surface form can have multiple meanings and entities can have multiple 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 [18] 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.

An example is given in the following two figures. In Fig. 2 a document with text fragment ... *Bruce* ... that contains an entity mention is shown (lower box). The surface form 'Bruce' of the entity mention can be linked to different possible entities, i.e. they are homonyms, thus exhibiting the same writing but different meanings. The overall set of all possible entities for a surface form is  $\mathcal{V}_{sf}$ . The dictionary known to the annotator  $\mathcal{W}$  is a subset of  $\mathcal{V}_{sf}$ . The dictionary known to the dataset containing the document is  $\mathcal{D}$ , also a subset of  $\mathcal{V}_{sf}$ . The likelihood of confusion for the surface form 'Bruce' is then determined by the cardinality of the union of the known entities  $\mathcal{D} \cup \mathcal{W}$ . The larger the cardinality, the higher is the likelihood of confusion.

In Fig. 3 a document with text fragment ... *Bruce* ... that contains an entity mention linking to the entity

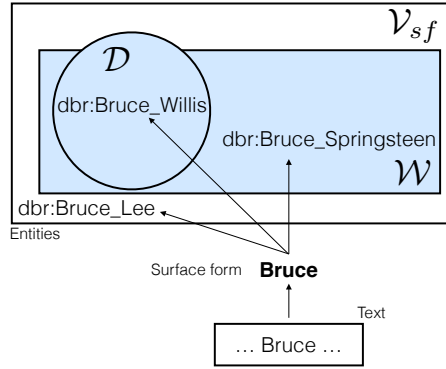


Fig. 2. The likelihood of confusion for a surface form is determined by the total number of possible entities known to the annotating system and the dataset  $\mathcal{D} \cup \mathcal{W}$ .

`dbr:Bruce_Willis` is shown. This entity could also be mapped from multiple other surface forms (synonyms). The overall set of all possible surface forms is  $\mathcal{V}_e$  (outer lower box). The annotator knows only  $\mathcal{W}$  (inner lower box), a subset of  $\mathcal{V}_e$ , and the dataset under consideration only contains  $\mathcal{D}$ , also a subset of  $\mathcal{V}_e$ . *Bruce* and *Bruce Willis* are both surface forms used within the dataset for the entity 'Bruce'. However, the annotation system provides *Bruce Walter Willis* as another additional possible surface form for this entity. The likelihood of confusion for an entity is then determined by the cardinality of the union of the known surface forms  $\mathcal{D} \cup \mathcal{W}$ .

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

1. Unknown to dictionary and dataset
2. Only known to the dictionary
3. Only known to the dataset
4. Known to dictionary and dataset

The annotator system dictionary  $\mathcal{W}$  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  $\mathcal{W}$ , the average likelihood of confusion is determined for surface forms  $S$  with  $c_{sf}: (\mathcal{W}, S) \rightarrow \mathbb{R}^+$  and entities  $E$  with  $c_e: (\mathcal{W}, S) \rightarrow \mathbb{R}^+$ :

$$c_{sf}(\mathcal{W}, S) = \frac{\sum_{s \in S} |\{e | s \in W\}|}{|S|} \quad (4)$$

$$c_e(\mathcal{W}, E) = \frac{\sum_{e \in E} |\{s | e \in W\}|}{|E|} \quad (5)$$

Since the dictionary is a multi set, the term  $\{x | y \in W\}$  refers to the set containing all elements  $x$  that are

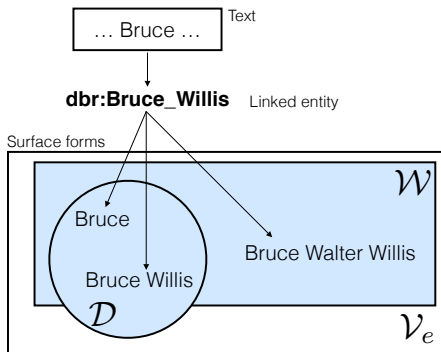


Fig. 3. The likelihood of confusion for an entity mention is the number of possible related surface forms shown in light blue.

referenced by a search variable  $y$ . The likelihood of confusion gives only a rough overview of how difficult it might be to correctly disambiguate the entities and surface forms contained in the dataset.  $c_{sf}(W, S)$  can also be seen as an indicator for homonyms, and  $c_e(W, E)$  as an indicator of synonyms.

### 2.5 Dominance (Level of diversity)

Erp et al. introduced the dominance as a measure of how commonly a specific surface form is really meant for an entity with respect to other possible surface forms [18]. 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 and an approximation of the dominance is computed based on the same dictionary. In addition to the work in [18] we estimate dominance for both sides the entity as well as the surface form side. For an entire dataset and a dictionary, the average dominance is determined in both directions.

In the one direction the amount of surface forms used for one specific entity in the dataset  $e(D)$  is divided by the amount of possible surface forms referencing that entity in the dictionary  $e(W)$ . For example, for the entity `dbp: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 `dbp:Angelina_Jolie` in the considered dataset. Again, the dominance of an entity determines how many different surface forms of this entity are used in the dataset (synonyms). Dominance indicates the expressiveness of the used vocabulary. An extensive vocabulary exhibits more diversity and is more appropriate to avoid over-fitting.

In the other direction we divide the amount of all entities for one specific surface form used within the dataset  $s(D)$  by the possible number  $s(W)$  referenced in the dictionary. For example, for the given surface form 'Anna' the dictionary provides 10 different entities, while the dataset only uses 2 entities for different mentions with surface form 'Anna', which results in a 20% dominance of 'Anna' for the dataset under consideration. Again, 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.

As shown in Fig. 2 and Fig. 3 the likelihood of confusion is determined by counting everything that is contained in the dictionary and the dataset. The dominance is related to this since it describes the coverage among the dataset and dictionary.

The average dominance for an entire dataset is computed over all entities  $e \in D$  with  $dom_e: (W, D) \rightarrow \mathbb{R}^+$  and surface forms  $s \in D$  with  $dom_{sf}: (W, D) \rightarrow \mathbb{R}^+$ :

$$dom_e(W, D) = \frac{\sum_{e \in E} \frac{e(D)}{e(W)}}{|e \in D|} \quad (6)$$

$$dom_{sf}(W, D) = \frac{\sum_{s \in E} \frac{s(D)}{s(W)}}{|s \in D|} \quad (7)$$

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 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 annotator 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 annotators 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.

To estimate the coverage of a mapping dictionary, the maximum recall measurement was introduced by [15]. For a dictionary  $W$  and an entire dataset  $D$  the maximum recall is defined as the intersection of entity mentions in the dataset and the dictionary for a given surface form  $S$  where  $max\_recall : (\mathcal{W}, \mathcal{D}) \rightarrow (0, 1)$ :

$$max\_recall(W, D) = \frac{|\{s | s \in D : s \in W\}|}{|s \in D|}. \quad (8)$$

## 2.6 Types

Since different NEL approaches focus on different categories of entities, we have implemented a filter to analyze the following DBpedia entity types separately: person, places, organizations, and others. Besides the focus of NEL approaches Erp et al. also stated that types of entities may be differently difficult to disambiguate such as person names might be more ambiguous and country names more or less unique [18]. For a dataset  $\mathcal{D}$ , the relative amount of entities of a specific type  $\mathcal{T}$  is determined by  $t : (\mathcal{D}, \mathcal{T}) \rightarrow (0, 1)$ :

$$t(\mathcal{D}, \mathcal{T}) = \frac{|\{e | e \in \mathcal{D} : e \in \mathcal{T}\}|}{|e \in \mathcal{D}|}. \quad (9)$$

## 2.8 Micro and Macro Measurement

In accordance to Cornolti et al. [1], we distinguish between micro and macro measurements for the following introduced measures: density, likelihood of confusion and maximum recall. Macro measurement aggregates the average of 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 were one single document, which consequently increases the influence of larger documents.

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 will be explained in detail.

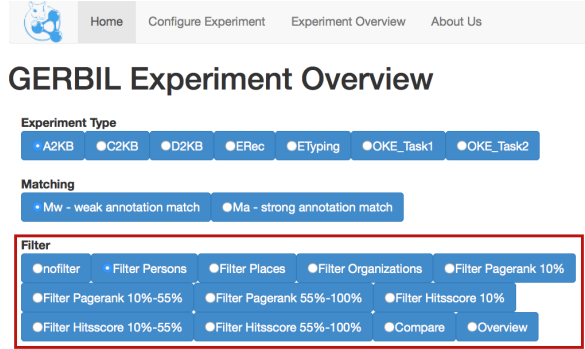


Fig. 4. New dataset filters for A2KB experiments in the GERBIL user interface.

### 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 another 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 arbitrary 'non-official' datasets. The implemented filter-cascade is of a generic type and can be adjusted via customized SPARQL queries. E. g. 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 entities. This segmentation is applied to PageRank as well as HITS values separately.

Buttons have been added as new control elements to the A2KB, C2KB, and D2KB overview pages in GERBIL (cf. Fig. 4). The user can now 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 HighCharts<sup>6</sup>. The existing charts are also replaced by the new chart API, since GERBIL was limited to only one chart type. The comparison view enables the user to view two filters at the same time as well as the average for all annotators on a specific filter. The overview shows several statistics for all datasets, such as e. g., total amount of types per filter, density, likelihood of confusion in average and total. The extended source code is publicly available at Github<sup>7</sup>. In addition, an online version is available<sup>8</sup>.

Before discussing the dataset statistics as a result of the new GERBIL extension the next section introduces 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 Github<sup>9</sup>.

To serialize the calculated statistics generated by the GERBIL extension as well as by the library, a vocabulary was 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

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

Table 1

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

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 e. g. holds the documents with the `hfts:referenceDocuments` property. On dataset level, there are 13 properties introduced, which hold the measurements missing-annotation, density, maximum recall, dominance and likelihood of confusion on dataset level. Some of them come with a micro and macro flavour while others are only computed once.

On document level six new properties were introduced to cover density, likelihood of confusion and maximum recall. The likelihood of confusion, prominence and the types are also assigned on entity mention level.

In Tab. 1 an overview over the introduced properties and their corresponding level is given. Fig. 5 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. 6 an example for the document level is presented (`nif:Context`). Besides with the existing NIF vocabulary the statistics are serialized with

<sup>6</sup><http://www.highcharts.com/>

<sup>7</sup><https://github.com/santifa/gerbil/>

<sup>8</sup><http://gerbil.s16a.org/>

<sup>9</sup><https://github.com/santifa/hfts>



```
<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. 5. 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. 6. An example of the new statistics properties on *document level* extending the KORE50 dataset.

the newly introduced *hfts:* properties. The entire definition and further documentation of the vocabulary can be found at Github<sup>10</sup>.

The following section introduces a selection of the most interesting results we have achieved so far as well as from the library as also from the GERBIL integration. Building on these insights, the dataset remixing is explained with various examples.

#### 4. Dataset Statistics and Remixing

Before it is introduced how documents from different datasets can be combined into customized new datasets, the quantification of dataset characteristics is presented.

##### 4.1. GERBIL Generated Results

The datasets and annotators have been analyzed according to the characteristics introduced in Sect. 2. In this section, only the most significant results are presented. A complete listing of the achieved results is available online<sup>11</sup>.

Fig. 7 shows the **percentage of empty documents** in a dataset. Overall, there are six datasets that contain

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

<sup>11</sup><http://gerbil.s16a.org/>

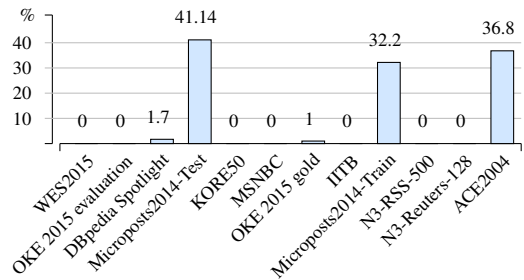


Fig. 7. Percentage of documents without annotations in a dataset

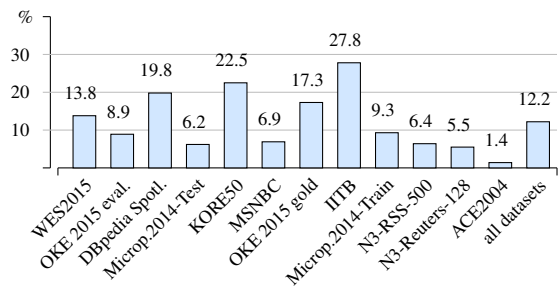


Fig. 8. Annotation density as relative number of annotations respective document length in words

empty documents while four of these show a significant (>30%) amount of empty documents. For A2KB tasks, these datasets will lead to an increased false positive rate and thus will lower the potentially achievable precision of an annotator. Therefore, empty documents should be excluded from evaluation datasets for a sound evaluation.

Fig. 8 shows the **annotation density** of the 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. In order to find evidence for this correlation, we have determined the Pearson correlation between density and achieved precision with a result of 0.7, which supports our original assumption. Especially for NEL tools based on machine learning it is of importance, if a sparsely annotated dataset is appropriate for the training task. Of course, this strongly depends on the application. Nevertheless, it is arguable, if sparseness is problematic for A2KB, because all annotators are facing the same



problem and the achieved results might still be comparable.

Table 2 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. As Steinmetz et al. have demonstrated there are many untyped entities in the DBpedia Spotlight and the KORE50 datasets. Therefore, an extra row for unspecified entities has been added to the table. 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 annotator 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.

Table 3 shows the achieved micro- $f_1$  **results of the annotators** for the D2KB task, partitioned in the same way as for table 2. The top row indicates the original GERBIL results (No Filter). Top results are indicated in green (bold) and the lowest results in red (italic). Each row indicates an entity restriction either by entity type or by entity popularity measure (PageRank and HITS) being applied before evaluation. For persons, organizations and places the results achieved by the annotators are rather similar, except for FOX, Entityclassifier.eu, and Dexter, which have achieved significantly lower scores. DBpedia Spotlight performs best for places and KEA for persons as well as for organizations. Developers are free to build on these results to optimize their systems accordingly. For example, the KEA system could be improved by investigating on why places are not sufficiently well recognized and linked.

The second and third partition of the table shows the results achieved for entities of different popularity according to PageRank and HITS. The subsets for high, medium and low popularity show that all annotators achieve rather similar results for each subset. This ob-

servations is further supported by the average PageRank and HITS values denoted in the last column. There is no significant difference in the achieved results for popular entities vs. less popular entities. More detailed results for the complete set of filter and dataset combinations as well as the results for the A2KB tasks can be obtained online.

Fig. 9 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 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 also confirmed by these numbers. ACE2004 exposes the highest average number of surface forms for possible entities (35), i. e. it contains many synonyms.

To measure a correlation between likelihoods of confusion for entities and surface forms with precision and recall, the following Pearson correlation values have been determined: entity-recall = 0.156, entity-precision = -0.858, surface-recall = 0.126, and surface-precision = -0.351. Carefully speaking, these results indicate negative correlations for precision, which was expected. Therefore, the more potential candidate entities exist for each surface form in a dataset (homonyms), the lower is the achieved precision. Likewise, the more different surface forms exist for the entities in a dataset (synonyms), the lower is precision.

With regard to recall, only a very slight positive correlation can be observed, which does not allow to draw a clear conclusion. Furthermore, the stated values do not include the KORE50 dataset, which was excluded as outlier since it exposes a very large number of homonyms within a rather small dataset only.

Fig. 10 shows the **average dominance of entities and surface forms** in percent. The blue 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

	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	all datasets
Persons	18.4	30.3	3.0	16.6	<b>45.1</b>	27.2	29.3	<i>2.4</i>	16.2	15.9	6.5	6.5	18.1
Org.	3.4	11.1	3.0	9.0	16.0	9.0	18.3	<i>2.0</i>	13.8	10.5	<b>20.7</b>	20.3	11.4
Places	9.4	14.0	8.2	8.9	6.9	17.5	14.5	<i>3.5</i>	14.2	7.2	17.2	<b>35.0</b>	13.0
unspecified	68.8	44.6	85.1	65.5	<b>32</b>	46.3	37.9	<i>92.1</i>	55.8	66.4	55.6	38.2	57.4
PageRank 10%	27.9	24.4	<b>30.0</b>	21.3	28.5	28.5	24.9	14.8	26.0	<i>14.3</i>	18.8	22.2	23.5
PageRank 10%-55%	48.9	39.5	47.6	<b>49.8</b>	48.6	32.2	<i>0.3</i>	29.8	45.8	23.0	31.4	37.6	36.2
PageRank 55%-100%	22.5	16.6	19.7	<b>28.0</b>	19.4	24.8	<i>7.7</i>	15.0	25.6	11.1	19.0	15.1	18.7
HITS 10%	28.4	21.1	32.4	31.4	27.8	29.8	26.9	<i>12.3</i>	<b>32.9</b>	18.3	19.0	28.4	25.7
HITS 10%-55%	12.9	12.4	18.2	14.4	20.8	<b>22.8</b>	<i>0.3</i>	12.2	13.6	7.3	9.1	11.4	13.0
HITS 55%-100%	<b>58.0</b>	47.0	48.2	51.8	47.2	32.1	50.2	35.2	50.6	<i>23.2</i>	40.6	15.3	41.6

Table 2

Percentage of entities by entity type and entity popularity per dataset

	Babelify	DBpedia Spotl.	Dexter	Entityclassifier.eu	FOX	KEA	TagMe 2	WAT	AGDISTIS	average
No Filter	0.53	0.56	0.39	0.33	<i>0.32</i>	<b>0.63</b>	0.59	0.58	0.52	0.49
Persons	0.81	0.69	0.53	0.57	<i>0.44</i>	<b>0.84</b>	0.77	0.80	0.74	0.69
Org.	0.71	0.83	0.65	0.75	<i>0.55</i>	<b>0.88</b>	0.79	0.80	0.77	0.75
Places	0.77	<b>0.82</b>	0.57	0.55	<i>0.54</i>	0.78	0.81	0.80	0.75	0.71
PageRank 10%	0.68	0.76	0.50	0.48	<i>0.39</i>	<b>0.79</b>	0.74	0.75	0.63	0.64
PageRank 10%-55%	0.69	0.75	0.50	0.50	<i>0.40</i>	<b>0.80</b>	0.75	0.74	0.62	0.64
PageRank 55%-100%	0.72	0.70	0.48	0.46	<i>0.36</i>	<b>0.81</b>	0.74	0.75	0.63	0.63
HITS 10%	0.67	0.78	0.48	0.48	<i>0.40</i>	<b>0.82</b>	0.74	0.74	0.62	0.64
HITS 10%-55%	0.69	0.74	0.51	0.52	<i>0.40</i>	<b>0.79</b>	0.75	0.75	0.64	0.64
HITS 55%-100%	0.68	0.69	0.48	0.47	<i>0.36</i>	<b>0.79</b>	0.74	0.73	0.61	0.62

Table 3

Micro-f<sub>1</sub> results of D2KB annotators by entity type and entity popularity per dataset

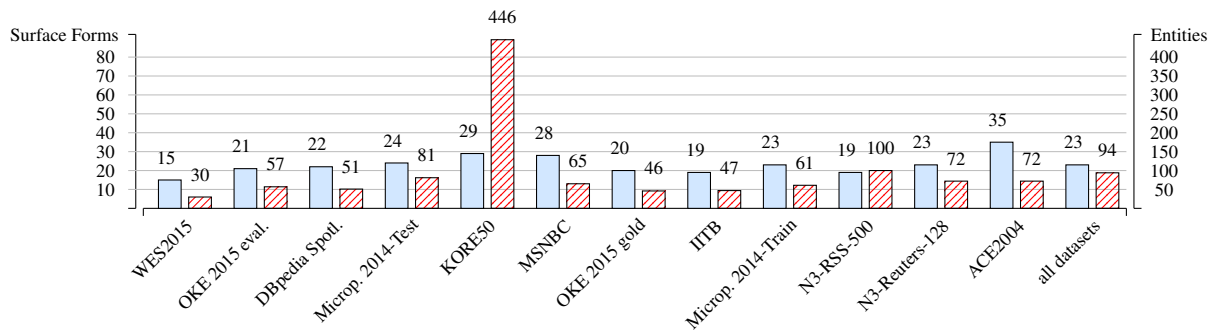


Fig. 9. Average number of surface forms per entity (blue, left) and average number of entities per surface form (red/hatched, right) indicating the likelihood of confusion for each dataset

existing surface forms in the dictionary. Referring to Fig. 10, the ACE2004 dataset uses only 8% of the surface forms existing in the dictionary. It indicates also how well the dataset's surface forms are covered by the dictionary's surface forms.

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

Referring to Fig. 10, the KORE50 dataset in which many persons are annotated uses only 9% of the possible entities for the contained surface forms. In average, entities are represented in the WES2015 dataset with 21% of their surface forms.

To verify a potential correlation between dominance for entities and surface forms with precision and recall, the following Pearson correlation values have been determined: entity-precision = 0.056, entity-recall = 0.063, surface-precision = -0.095, surface-recall = 0.674. Only the surface-recall relation shows a potential positive correlation. That means, to enable improved recall, the surface form dominance of the datasets has to be increased. Again, because of the diversity of the datasets and only scarcely available data points, these numbers are rather vague and only enable a tentative insight.

Since the datasets with a high likelihood of confusion have a low dominance, it is arguable that these two measures are somehow contrary. E. g. 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 the 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. E. g. 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 are ideal to vary the surface forms for an entity. This is commonly known from articles or essays, where the author usually tries to minimize surface form repetitions 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.

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

#### 4.2. Remixing Customized Datasets

The basic idea of remixing NEL benchmark datasets is to tailor new customized datasets from the existing 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 one single triple store. This enables to quickly access the dataset documents via the SPARQL query language. In particular, SPARQL CONSTRUCT queries can be used to select exactly those triples from the document annotations that meet a particular criteria, as e. g., popular per-

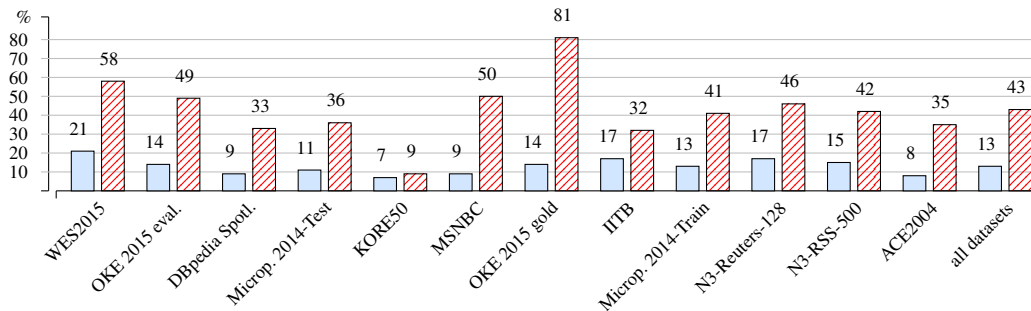


Fig. 10. Average dominance for surface forms (blue) and entities (red) per dataset

sons, high possible maximum recall, difficult places, or any other arbitrary criteria, which can be specified via SPARQL filter rules.

Therefore, we introduce the basic query shown in Fig. 11. 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 amount 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. 12.

Another example is given in Fig. 13. The SPARQL subselect chooses documents that contain persons and aggregates their count. Afterwards the CONSTRUCT statement selects only documents that contain more than four persons and with a maximum recall of at least 0.8.

To underline that any kind of filter can be applied, Fig. 14 shows a more exotic example using a federated query to select only documents from the RDF graph with persons born before 1970. To achieve this, the official DBpedia SPARQL endpoint is queried for additional information that is not present within the given benchmark datasets. More examples can be found at Github<sup>12</sup>.

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 subject of future work to normalize and harmonize the statistics adequately. Second, the proposed query examples are

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

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

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

Fig. 11. Basic query selecting only documents with a maximumRecall >= 1.0

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

  # limit the amount 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. 12. This query limits the number of selected documents.

based on document level. That means that 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 annotations, but this can lead to a missing annotation scenario which might result in a drop of recall in the A2KB task.

<sup>12</sup><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. 13. Extract documents with a maximum recall of 0.8 and at least four 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. 14. A query which selects documents containing persons born before 1970. This is archived by retrieving additional data from the DBpedia SPARQL endpoint.

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

## 5. Conclusion

In this paper an extension of the GERBIL framework has been introduced to enable a more fine grained evaluation of NEL annotators. It was shown that not all of the available datasets are equally suitable for the A2KB task. According to our evaluation, the best suited datasets for the A2KB task are WES2015, OKE 2015 evaluation, DBpedia Spotlight, KORE50 and IITB. We have also shown that the original general assumption about popular entities and disambiguation does not hold for the considered datasets and annotators, i.e. popular entities are not easier to disambiguate than less popular entities. The average scores achieved by each annotator for different levels of entity popularity are almost identical.

According to our predefined entity categories, the KORE50 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 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. As a result, users might chose the best suited annotator for specific texts according to the properties of the considered texts.

We have documented that some of the presented measures directly correlate to precision and recall. Since there are only a few data points available and the datasets exhibit a strong variation, the correlation numbers should be considered with caution. Also we have introduced a stand-alone library 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.

Ongoing research is focused on the implementation of additional measures, such as e. g. those introduced by [2,7] including the analysis of NIL-annotations. Future work will additionally include the creation of a user-friendly web front-end to ease the use of remixing datasets. Also we would like to introduce difficulty levels for datasets along with new properties for annotation, which might be useful for further remixing, e. g. a distinction of the NEL annotation for common and proper nouns.

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

In summary, evaluation on a more diverse as well as fine granular level will enable a better understanding of

the NEL process and likewise fosters the development of improved NEL annotators.

## References

- [1] 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*, pages 249–260. ACM, 2013.
- [2] B. Hachey, J. Nothman, and W. Radford. Cheap and easy entity evaluation. In *52nd Annual Meeting of the Association for Computational Linguistics*, pages 464–469. ACL, 2014.
- [3] S. Hellmann, J. Lehmann, S. Auer, and M. Brümmer. Integrating NLP using linked data. In *International Semantic Web Conference*, pages 98–113. Springer, 2013.
- [4] J. Hoffart, S. Seufert, D. B. Nguyen, M. Theobald, and G. Weikum. Kore: Keyphrase overlap relatedness for entity disambiguation. In *21st ACM Int. Conf. on Information and Knowledge Management*, pages 545–554, New York, NY, USA, 2012. ACM.
- [5] X. Ling, S. Singh, and D. S. Weld. Design Challenges for Entity Linking. *Transactions of the Association for Computational Linguistics*, 3:315–28, 2015.
- [6] L. Page, S. Brin, R. Motwani, and T. Winograd. The PageRank Citation Ranking: Bringing Order to the Web. *Stanford InfoLab*, 1999.
- [7] S. Pradhan, X. Luo, M. 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.
- [8] D. Reddy, M. Knuth, and H. Sack. DBpedia GraphMeasures. Hasso Plattner Institute, Potsdam, July 2014, <http://s16a.org/node/6>.
- [9] 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.
- [10] G. Rizzo and R. Troncy. NERD: A framework for unifying named entity recognition and disambiguation web extraction tools, Eurecom 3677, Avignon, France, 2012.
- [11] 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*. ELRA, 2014.
- [12] M. Röder, R. Usbeck, and A.-C. Ngonga Ngomo. Gerbil’s new stunts: Semantic annotation benchmarking improved. Technical report, Leipzig University, 2016.
- [13] W. Shen, J. Wang, and J. Han. Entity linking with a knowledge base: Issues, techniques, and solutions. *IEEE Transactions on Knowledge and Data Engineering*, 27(2):443–460, Feb 2015.
- [14] A. Singhal. Introducing the knowledge graph: things, not strings. *Official Google Blog*, May, 2012.
- [15] N. Steinmetz, M. Knuth, and H. Sack. Statistical analyses of named entity disambiguation benchmarks. In *Proc. of NLP & DBpedia 2013 workshop at 12th Int. Semantic Web Conference*. CEUR-WS, 2013.
- [16] T. Tietz, J. Waitelonis, J. Jäger, and H. Sack. Smart Media Navigator: Visualizing recommendations based on Linked Data. In *13th Int. Semantic Web Conference, Industry Track*, pages 48–51, 2014.
- [17] R. Usbeck et al. GERBIL – general entity annotation benchmark framework. In *24th World Wide Web Conf.* ACM, 2015.
- [18] 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 *Proc. of the 10th Int. Conf. on Language Resources and Evaluation (LREC 2016)*, Paris, France, May 2016. European Language Resources Association (ELRA).
- [19] J. Waitelonis, C. Exeler, and H. Sack. Linked Data Enabled Generalized Vector Space Model to Improve Document Retrieval. In *NLP & DBpedia 2015 workshop at 14th Int. Semantic Web Conf.*, volume 1581, pages 33–44. CEUR-WS, 2015.
- [20] 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, New York, NY, USA, 2016. ACM.