



City Research Online

City, University of London Institutional Repository

Citation: Abdelmageed, N., Chen, J., Cutrona, V., Efthymiou, V., Hassanzadeh, O., Hulsebos, M., Jimenez-Ruiz, E., Sequeda, J. & Srinivas, K. (2022). Results of SemTab 2022. Semantic Web Challenge on Tabular Data to Knowledge Graph Matching, 3320, ISSN 1613-0073

This is the published version of the paper.

This version of the publication may differ from the final published version.

Permanent repository link: <https://openaccess.city.ac.uk/id/eprint/29744/>

Link to published version:

Copyright: City Research Online aims to make research outputs of City, University of London available to a wider audience. Copyright and Moral Rights remain with the author(s) and/or copyright holders. URLs from City Research Online may be freely distributed and linked to.

Reuse: Copies of full items can be used for personal research or study, educational, or not-for-profit purposes without prior permission or charge. Provided that the authors, title and full bibliographic details are credited, a hyperlink and/or URL is given for the original metadata page and the content is not changed in any way.

City Research Online:

<http://openaccess.city.ac.uk/>

publications@city.ac.uk

Results of SemTab 2022

Nora Abdelmageed^{1,*}, Jiaoyan Chen², Vincenzo Cutrona³, Vasilis Efthymiou⁴,
Oktie Hassanzadeh⁵, Madelon Hulsebos⁶, Ernesto Jiménez-Ruiz^{7,8}, Juan Sequeda⁹ and
Kavitha Srinivas⁵

¹Friedrich Schiller University Jena, Jena, Germany

²University of Manchester, UK

³SUPSI, Switzerland

⁴FORTH-ICS, Greece

⁵IBM Research, USA

⁶University of Amsterdam, The Netherlands

⁷City, University of London, UK

⁸SIRIUS, University of Oslo, Norway

⁹data.world, US

Abstract

SemTab 2022 was the fourth edition of the Semantic Web Challenge on Tabular Data to Knowledge Graph Matching, successfully collocated with the 21st International Semantic Web Conference (ISWC) and the 17th Ontology Matching (OM) Workshop. SemTab provides a common framework to conduct a systematic evaluation of state-of-the-art systems. In this paper, we give an overview of the 2022's edition of the challenge and summarize the results.

Keywords

Tabular data, Knowledge Graphs, Matching, SemTab Challenge, Semantic Table Interpretation

1. Motivation

Tabular data are the most frequent input to data analytics pipelines, thanks to their high storage and processing efficiency. Also, the tabular format allows users to represent the information in a compacted way, by exploiting the clear data structure defined by rows and columns. However, such clear structure does not imply a clear understanding of the semantic structure (*e.g.*, relationships between columns), as well as the meaning of the content (*e.g.*, if data are about a specific topic). The lack of understanding hinders data analytics processes, requiring additional effort to properly understand the data first. Gaining the semantic understanding is valuable for many applications, including data cleaning, data mining, data integration, data analysis and machine learning, and

SemTab@ISWC 2022, October 23–27, 2022, Hangzhou, China (Virtual)

*Corresponding author.

✉ nora.abdelmageed@uni-jena.de (N. Abdelmageed); jiaoyan.chen@cs.ox.ac.uk (J. Chen);
vincenzo.cutrona@supsi.ch (V. Cutrona); vefthym@ics.forth.gr (V. Efthymiou); hassanzadeh@us.ibm.com
(O. Hassanzadeh); m.hulsebos@uva.nl (M. Hulsebos); ernesto.jimenez-ruiz@city.ac.uk (E. Jiménez-Ruiz);
juan@data.world (J. Sequeda); kavitha.srinivas@ibm.com (K. Srinivas)

ORCID 0000-0002-1405-6860 (N. Abdelmageed)



© 2022 Copyright for this paper by its authors. Use permitted under Creative Commons License Attribution 4.0 International (CC BY 4.0).

CEUR Workshop Proceedings (CEUR-WS.org)

knowledge discovery. For example, the semantic understanding can help in assessing what kind of transformations are more appropriate for a dataset, or which datasets can be integrated to enable new analytics (*e.g.*, marketing analysis) [1].

In addition to their efficiency, the huge availability of tabular data on the Web makes Web tables a valuable source to consider for data miners (*e.g.*, open data CSV files). Adding semantic information to Web tables is useful for a wide range of applications, including web search, question answering, and knowledge base construction.

Tabular data to Knowledge Graph (KG) matching is the process of clarifying the semantic meaning of a table by mapping its elements (*i.e.*, cells, columns, rows) to semantic tags (*i.e.*, entities, classes, properties) from KGs (*e.g.*, Wikidata, DBpedia). The task difficulty increases when table metadata (*e.g.*, table captions, table description, or column names) being missing, incomplete or ambiguous.

The tabular data to KG matching process is typically broken down into the following tasks:

- cell to KG entity matching (CEA task),
- column to KG class matching (CTA task), and
- column pair to KG property matching (CPA task).

Over the last decade several approaches made advances in addressing one or several of above tasks, also constructing benchmark datasets ([2, 3, 4, 5]). The creation of **SemTab**¹ [6, 7, 8] aimed at putting this significant amount of work into a common framework, enabling the systematic evaluation of state-of-the-art systems. The ambition is to make **SemTab** becoming the reference challenge in the Semantic Web community, in the same way the OAEI² is for the Ontology Matching community.³

2. The Challenge

The **SemTab** 2022 challenge has been organised into two different tracks: the **Accuracy Track**, which is the standard track proposed in previous editions; and the **Datasets Track**, which focuses on applications in real-world settings where the output of matching systems can contribute. The datasets track was also open to the submission of novel benchmark datasets. **SemTab** 2022 also featured an *Artifacts Availability Badge* that was applicable to both tracks.

2.1. Accuracy Track

The Accuracy Track included 3 rounds, running from May 26 to October 15. Different target KGs were used across rounds (see Table 1):

- Wikidata (WD, W) [9]: <https://zenodo.org/record/6643443>
- Schema.org (SCH, S) [10]: https://gittables.github.io/downloads/schema_20210528.pkl
- DBpedia (DBP, D) [11]: <http://downloads.dbpedia.org/wiki-archive/> (version 2016-10 & 2022-03)

¹<http://www.cs.ox.ac.uk/isg/challenges/sem-tab/>

²<http://oaei.ontologymatching.org/>

³<http://ontologymatching.org/>

Table 1

Datasets used across SemTab 2022 rounds.

	Rounds			Tasks			Target KGs		
	R1	R2	R3	CTA	CPA	CEA	DBpedia	Wikidata	Schema.org
2T		✓		✓		✓	✓	✓	
HardTables	✓	✓		✓	✓	✓		✓	
BiodivTab			✓	✓		✓	✓		
GitTables			✓	✓			✓		✓

The different rounds of SemTab 2022 have been organised to evaluate participating systems on various datasets with variable difficulty. Unlike the previous editions of SemTab where the rounds were run with the support of Alcrowd, this year, we asked the participants to submit their solutions once per week via a Google Form. We created a submission form for each round and we evaluated the submitted results at the beginning of each week during each round.

2.1.1. Datasets

The different datasets used to run SemTab 2022 rounds are reported in Table 1, with some statistics available in Table 2. Unlike the previous editions where the ground truth was hidden from the participant, this year, we provided a partial ground truth data to the participants during the challenge itself. The teams could use these ground-truth labels for evaluating their methods locally. Thus, we report datasets’ statistics per split. All the datasets are available in Zenodo as follows:

- **Tough Tables (2T) [12]**: a dataset featuring high-quality manually-curated tables with non-obviously linkable cells, *i.e.*, where values are ambiguous names, typos, and misspelled entity names. These challenges are particularly relevant for the annotation of structured legacy sources to existing KGs.
Link: <https://doi.org/10.5281/zenodo.7419275>
- **HardTables (HT) [6]**: datasets with tables generated using an improved version of our data generator that creates realistic looking tables using SPARQL queries [6]. It is the largest dataset used in SemTab.
Link: <https://doi.org/10.5281/zenodo.7416036>
- **BiodivTab [13]**: a dataset with tables from real-world biodiversity research datasets. Original tables have been adapted for the SemTab challenge. This year featured DBpedia as a target KG instead of Wikidata
Link: <https://doi.org/10.5281/zenodo.7319654>
- **GitTables [14]**: a large-scale corpus of tables extracted from CSV files in GitHub. The main purpose of this dataset is to facilitate learning table representation models and applications in *e.g.*, data management. A subset of tables has been curated for benchmarking column type detection methods in SemTab. The GitTables set for this SemTab edition was larger than in 2021 to enable the training of potential data-driven methods.
Link: <https://zenodo.org/record/7091019>

⁴Target column train and test splits for CTA are distributed across all tables.

Table 2

Statistics of the datasets in each SemTab 2022 round. For target values: *W*=Wikidata; *D*=DBpedia; *S*=Schema.org, *C*=Class, *P*=Property.

Validation	HardTables		2T	BiodivTab	GitTables ⁴
	Round 1	Round 2	Round 2	Round 3	Round 3
# Tables	200	457	36	5	6,892
Avg. # Rows	5.74	5.54	683.3 _W 1,373 _D	119	61.1
Avg. # Cols	2.6	2.56	4.31 _W 4.42 _D	20.4	14.4
# CEA Targets	1,406 _W	1,983 _W	81,126 _W 177,453 _D	1,463 _D	NA
# CTA Targets	240 _W	398 _W	97 _W 111 _D	43 _D	6,228 _{D,P} 4,411 _{S,P} 1,001 _{S,C}
# CPA Targets	319 _W	348 _W	NA	NA	NA
Test	Round 1	Round 2	Round 2	Round 3	Round 3
Tables #	3,691	4,679	144	45	6,892
Avg. # Rows	5.69	5.57	1,180.69 _W 1,008.26 _D	275.73	61.1
Avg. # Cols	2.56	2.6	4.49 _W 4.47 _D	24.36	14.4
# CEA Targets	26,189 _W	22,009 _W	586,118 _W 48,6203 _D	31,942 _D	NA
# CTA Targets	4,511 _W	4,534 _W	443 _W 429 _D	526 _D	6,228 _{D,P} 4,411 _{S,P} 1,000 _{S,C}
# CPA Targets	5,745 _W	3,954 _W	NA	NA	NA

2.1.2. Participation

Table 3 shows the participation per round. Compared with previous editions, we had 10 participants (vs 11 in 2021) submitting to at least one round. Similar to last year, we identified 8 core participants, which completed in almost all of the 14 SemTab tasks, and submitted a system paper to the challenge: *KGCODE-Tab* [15], *DAGOBAB* [16], *s-elBat* [17], *TSOTSA* [18], *JenTab* [19], *Kepler-aSI* [20], *Mertens* [21], and *SemInt* [22].

2.1.3. Evaluation measures

As per the previous editions, systems have been evaluated on a single annotation for each provided target, for all the tasks; *i.e.*, in CEA, target cells are to be annotated with a single entity from the target KG; in CTA, target columns are to be annotated with a single type from the target KG (as fine-grained as possible).

Table 3

Participation in the SemTab 2022 challenge. W =Wikidata; D =DBpedia; S =Schema.org.

	Round 1	Round 2		Round 3	
	HardTables	2T	HardTables	BiodivTab	GitTables
CEA	6	$\begin{matrix} 5_W \\ 5_D \end{matrix}$	5	5	NA
CTA	8	$\begin{matrix} 6_W \\ 5_D \end{matrix}$	6	7	$\begin{matrix} 5_D \\ 4_S \end{matrix}$
CPA	6	NA	6	NA	NA
Total	9	6	6	7	5

The evaluation measures for CEA, CPA and CTA (DBpedia and Schema.org) are the standard Precision, Recall and F1-score, as defined in Equation 1:

$$P = \frac{|\text{Correct Annotations}|}{|\text{System Annotations}|}, R = \frac{|\text{Correct Annotations}|}{|\text{Target Annotations}|}, F1 = \frac{2 \times P \times R}{P + R} \quad (1)$$

where target annotations refer to the target cells for CEA, the target columns for CTA, and the target column pairs for CPA. We consider an annotation as *correct* when it is included within the ground truth set (a target cell usually has multiple annotations in the ground truth, because of redirect and same-as links in KGs).

Given the fine-grained type hierarchy in Wikidata, we adopted approximations of Precision and Recall in the CTA evaluation [7]. Approximations adapt their numerators to consider partially correct annotations, *i.e.*, annotations that are ancestors or descendants of the ground truth (GT) classes. The correctness score $cscore$ of a CTA annotation α considers the distance between the annotation and the GT classes in the type hierarchy, and it is defined as

$$cscore(\alpha) = \begin{cases} 0.8^{d(\alpha)}, & \text{if } \alpha \text{ is in GT, or an ancestor of the GT, with } d(\alpha) \leq 5 \\ 0.7^{d(\alpha)}, & \text{if } \alpha \text{ is a descendant of the GT, with } d(\alpha) \leq 3 \\ 0, & \text{otherwise;} \end{cases} \quad (2)$$

where $d(\alpha)$ is the shortest distance to one of the GT classes (as for CEA, also CTA GT columns may have multiple classes). For example, $d(\alpha) = 0$ if α is a class in the ground truth ($cscore(\alpha) = 1$), and $d(\alpha) = 2$ if α is a grandchild of a class in the ground truth ($cscore(\alpha) = 0.49$). Types in the higher level(s) of the KG type hierarchy are not considered in the GT (*e.g.*, `Q35120 [entity]` in Wikidata). Given the correctness score $cscore$, approximated Precision (AP), Recall (AR), and F1-score (AF1) for the CTA evaluation are as follows:

$$AP = \frac{\sum cscore(\alpha)}{|\text{System Annotations}|}, AR = \frac{\sum cscore(\alpha)}{|\text{Target Annotations}|}, AF1 = \frac{2 \times AP \times AR}{AP + AR} \quad (3)$$

2.1.4. Results

Table 4 contains the average F1-score achieved by the participating systems per round. The Tough Tables dataset still represents a challenge for almost all the systems, especially considering the fact

Table 4

Average F1-score consider for participating systems in each round. *W*=Wikidata; *D*=DBpedia; *S*=Schema.org.

	Round 1	Round 2		Round 3	
	HardTables	2T	HardTables	BiodivTab	GitTables
CEA	0.88 _W	0.84 _W 0.56 _D	0.69 _W	0.56 _D	NA
CTA	0.89 _W	0.35 _W 0.32 _D	0.77 _W	0.49 _D	0.33 _D 0.50 _S
CPA	0.80 _W	NA	0.76 _W	NA	NA

that the dataset is almost the same as in SemTab 2020.⁵ The BiodivTab and GitTables datasets brought additional complexity in Round 3, highlighting that real-world tables are challenging.

CEA task. Results for the CEA task are reported in Figure 1 for all the datasets. In Round 1, almost all the systems performed well on the HardTables dataset (automatically generated). Starting from Round 2, datasets have been selected to increase the number of tricky cases to solve. We used a new version of HardTables, mainly built of tables that we knew participants failed to annotate in previous rounds and editions, in addition to Tough Tables. As expected, performance on HardTable decreases, showing that systems still fail in annotating tricky cases (even if they have been already seen); by comparing results on Tough Tables with previous SemTab editions, more systems achieved an F1-score over 0.8, possibly because of the availability of the validation set. In this round, systems focused on Wikidata-based datasets; indeed, just 4 out of 8 systems submitted a solution for the DBpedia-based Tough Table dataset. In Round 3, the complexity brought by the (relatively small) tables in the BiodivTab dataset still represents a new problem to solve, showing a reduced performance by all systems. However, the scores are generally improved compared to the last year’s version of such a dataset [8]. This year, the top F1-score is 91% by KGCODE-Tab compared to 60% by JenTab last year for CEA task.

CTA task. In Figure 2, the results in the CTA task resemble the trend already seen from the CEA results. This is an indicator that most of the systems solve the CTA task based on annotations found in the CEA. Additional challenges have been included in Round 3 with the GitTables dataset, where we can see a critical performance drop for all the involved systems. The column type annotations from DBpedia seem harder than annotations from Schema.org. However, the performance of systems for both ontologies improved compared to last year [8]: the top F1-score for the CTA task for Schema.org is 66% compared to 19% last year, and for DBpedia 59% versus 0.04%. BiodivTab seems tough for s-elbat and SemInt, and moderate for JenTab. However, other systems achieved very good results. It is worth emphasising that, given the general picture provided by the results in CTA, more research is needed to make existing systems able to deal with real-world tables, where the cells may be missing a correspondence to the target KG.

CPA task. Results for the CPA task are plotted in Figure 3. Currently, only HardTables provides a GT for CPA. Results are overall positive for all the systems, with a slightly general decrease

⁵Minor changes to adapt Tough Tables to SemTab 2022: some Wikidata targets have been updated to the Wikidata version adopted in SemTab 2022; the original dataset has been split into validation and test sets.

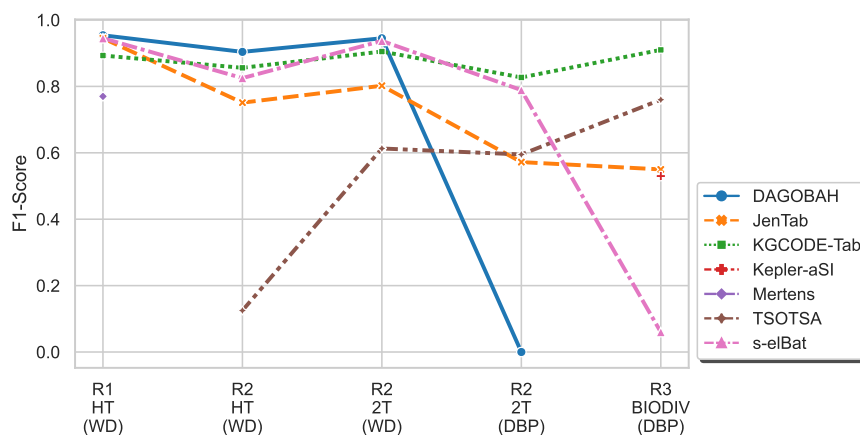


Figure 1: Results in the CEA task for the core participants.

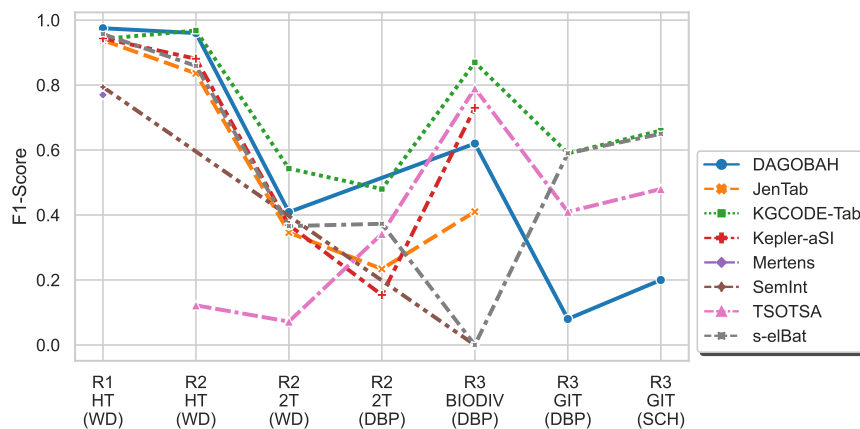


Figure 2: Results in the CTA task for the core participants.

from Round 1 to Round 2 for all participants. Again, this behavior can be explained by the fact that most systems use the CEA results to solve the CPA task, and the CEA scores for this dataset are high, overall.

2.2. Datasets Track

This new track aims at addressing applications in real-world settings that take advantage of the output of the matching systems. Challenging dataset proposals have also been accepted and may be included in future editions of SemTab.

2.2.1. Results

We opened the call for this track on 25 August 2022. We received four submissions, three of them were accepted and the fourth one was accepted as a poster (2 pages) since it is still at a

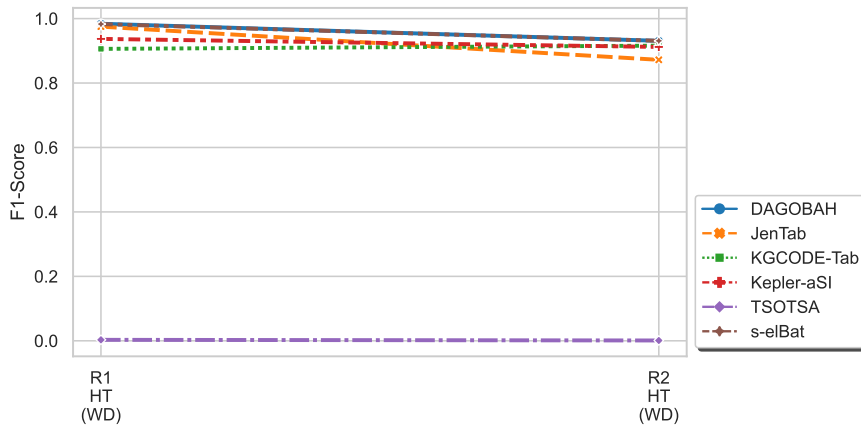


Figure 3: Results in the CPA task for the core participants.

preliminary stage. We give an overview of each of them as follows:

- **Wikary** [23] introduces a new gold standard for the Semantic Table Interpretation tasks. It proposes a promising expansion from binary to n-ary tuples to reflect the conditional properties of certain (object, predicate, subject) predicates. This will likely be helpful for matching tables with n-ary entities and lower KG coverage. Wikary contributes a large, diverse, multilingual and metadata-rich set of Wikipedia tables that are matched to n-ary statements and qualifiers from Wikidata. A subset of the dataset underwent a manual quality evaluation by asking annotators to annotate some tables.
- **MammoTab** [24] is a benchmark with a large number of tables extracted from Wikipedia, and annotations of Wikidata entities and semantic types (classes), and NIL (not matched with any KG entities). In comparison with state-of-the-art datasets, the annotations consider more key challenges including disambiguation, homonym, matching, NIL-mentions, literal and named-entity, missing context, etc. The annotations can be used for the CEA and CTA tasks of **SemTab**, but not for CPA. Such dataset is generated automatically, allowing it to be of a large scale (almost a million tables). The authors reported evaluation results on one of the top-performing **SemTab** participating systems, **MTab** [25], showing that this dataset is more challenging than previous **SemTab** benchmarks.
- **SOTAB** [26] is a benchmark dataset for the CTA and CPA tasks of **SemTab**, generated from the WDC Schema.org table corpus. The dataset consists of 88.5k tables, annotated with Schema.org as a target KG, exploiting the already existing annotations of the WDC table corpus. The subset of 88.5k tables was selected from the entire corpus, considering the language of the tables (targeting English tables mostly), columns with missing values, columns with heterogeneous formatting (e.g., date values expressed in different formats in the same table column), and corner cases (i.e., columns containing cell values that are tough to disambiguate). The dataset is split into train/validation/test sets, as usual, but more than that, it is also further divided into more test sets, targeting different challenges. For example, there is a test set just for addressing the missing values challenge, another one for corner cases, and another one for heterogeneous formats. In addition to those splittings,

the authors offer a small subset of the training data to check the effectiveness of systems that are trained with fewer examples.

- **FCTables** [27] is an interesting dataset that has the potential to be used in the **SemTab** challenge. It presents a construction for Food Composition Tables (FCT) benchmark from various data sources. The authors have the initiative of publishing such tables in CSV file format which is unlike the previous sources (INFOODS and LanguaL) both have tables in PDF. The benchmark covers various countries and languages. The authors highlighted potential applications for such a benchmark. Such benchmark is currently accepted as a poster paper (2 pages) since it is unclear how ready the dataset is in order to be included in the challenge. For example, the annotations for CTA, CEA and CPA are in progress.

2.3. Artifacts Availability Badge

In 2021, **SemTab** included a new track focusing on system usability. The main goal of this track was to mitigate a pain point in the community: the lack of publicly available, easy-to-use, and generic solution to address the needs of a variety of applications and settings. In 2022, the usability track was replaced by the *Artifact Availability Badge*, that applies to both tracks in 2022. This badge is given to submissions (regardless of the track) if all dependencies are verified to be accessible and sufficiently reusable. The goal of this badge is to motivate authors to publish and document their code and data, so that others can (re)use these artifacts and potentially reproduce the results.

2.3.1. Evaluation measures

The criteria used to assess submissions for eligibility of the Artifacts Availability Badge are:

- Publicly accessible data (if applicable).
- Publicly accessible source code.
- Clear documentation of the code and data.
- Open-source dependencies.

2.3.2. Results

Almost all core participants obtained good results in this track, by performing well on at least two criteria (open-source code and clear documentation). We report the evaluation details in Table 5. In general, tools requirements vary in complexity, but they are reasonable overall (*e.g.*, pre-processing required, like creating new indexes or embeddings). **JenTab** provides its pre-computed lookups and indices, thus, it has a checkmark under open-source data. Considering the other criteria, the evaluation panel concluded that most of the tools are pre-configured and can potentially be used out of the box: for example, **JenTab** has been packaged with Docker to ease deployment on local premises. In addition, **JenTab** is the only system released as open source under a permissive license (Apache 2.0).

For the Datasets Track, all submissions have open-source data and code and vary in the details of the provided documentation. **FCTables** is an exception, since the released code under its GitHub repository is for the table parsing component only, and not for their entire pipeline.

Table 5
Artifacts Availability Badge Evaluation details.

	Open-source data	Open-source code	Clear documentation (1-5)	Open-source dependencies
KGCODE-TAB				
s-elBat		✓	3	
DAGOBAB				
TSOTSA				
JenTab	✓	✓	5	✓
Kepler-aSI				
Mertens		✓	2.5	
SemInt				
Wikary	✓	✓	3.5	✓
MammoTab	✓	✓	2.5	
SOTAB	✓	✓	4.5	✓
FCTables	✓		2	

2.4. Awards

As in previous editions, the best systems in each track were awarded across different tracks:

- **Accuracy Track:** KGCODE-TAB and DAGOBAB (1st prize) were the highest performing systems in most of the tasks, showing appreciable improvements over previous years.
- **Dataset Track:** SOTAB won the first prize.
- **Artifacts Availability Badge:** SOTAB and Wikary got the badge from the dataset track, while JenTab was the only system to get this badge.

3. Lessons Learned and Future Work

Challenging HT Datasets. We have been using the same automated dataset generation process, with some variations that make it more challenging, since the first SemTab challenge. However, we applied a kind of filtration step to control the balance between the easy and hard cases. Unlike the previous editions of the challenge, this year the average F1-score is below 90% for the HT datasets (see Table 4). This proves the effectiveness of our filtration process and makes the dataset much harder to solve.

Ground-truth Quality. An accurate evaluation of systems is important for any benchmark or challenge and relies on the quality of the ground truth annotations. SemTab features diverse sets of tables that all are extracted and annotated differently. Some annotation procedures may yield inconsistent or erroneous annotations, which introduces noise in the development and evaluation of systems. Moreover, obtaining perfect ground truth for tables is hard as cell entities (CEA), column types (CTA), and column pair relations (CPA) could possibly correspond to multiple labels due to, for example, synonymy or hierarchy relations. In line with these challenges, a few participants indicated that some ground truth annotations in SemTab datasets might be questionable as well. This motivates an effort to improve the ground truth annotations for future editions, possibly in collaboration with the community, as suggested by a participant.

Artifacts Availability Badge. The introduction of the Artifacts Availability Badge extends the *Usability Track* from SemTab 2021. It encouraged the participating systems to provide publicly accessible resources. Our goal was exactly to emphasize this, despite the competitive nature that a challenge may have. We note that the current conditions for receiving this badge may have been too restrictive, so we are considering providing more badges of a narrower scope in the future (e.g., one badge for publicly accessible code, one for reproducible results, etc).

Dataset track. We believe that the call of the dataset track has grasped more attention from the community by introducing their own datasets. Compared to the last year, “Applications Track”, this year we received more contributions, four datasets, while last year we had only two submissions. Not all of them are ready to be used in the challenge, but they show a promising interest within the community. Such contributions from the community like Wikary and SOTAB help in extending the SemTab benchmark with new challenges that are hard to reproduce in synthetic datasets like HT. Thus, this new track has been an important addition to SemTab.

Visibility & Increased Impact. SemTab gained an increasing and broader attention from the community. This year, before the official start of the challenge, we presented SemTab in the Knowledge Graph Construction (KGC)⁶ Workshop, co-hosted with the ESWC conference. In addition, SemTab has contributed for four years to the Ontology Matching (OM)⁷ workshop, co-hosted with the ISWC conference. Such dissemination activities, in combination with the new Datasets Track, resulted in more contributions to SemTab 2022 overall. The proceedings in 2022 contains 12 papers (vs. 8 in 2021) with four of them belonging to the Datasets track. This shows diverse artifacts from SemTab 2022. We also plan to continue the challenge and propose a workshop with KG matching specificity. Either by organizing our own workshop or by joining one of the existing ones.

Acknowledgements

We thank the challenge participants, the ISWC & OM organisers, and our sponsor IBM Research that played a key role in the success of SemTab. We also thank Paul Groth and Çağatay Demiralp for their contribution to GitTables, and Sirko Schindler and Birgitta König-Ries for their contribution to BiodivTab. This work was also supported by the SIRIUS Centre for Scalable Data Access (Research Council of Norway), Samsung Research UK, the EPSRC projects UK FIRES and ConCur, and the HFRI project ResponsibleER (No 969). Finally, we like to acknowledge that the organization was greatly simplified by using the EasyChair conference management system and the CEUR-WS.org open-access publication service.

References

- [1] V. Cutrona, F. D. Paoli, A. Košmerlj, N. Nikolov, M. Palmonari, F. Perales, D. Roman, Semantically-Enabled Optimization of Digital Marketing Campaigns, in: International

⁶<https://kg-construct.github.io/workshop/2022/>

⁷<http://www.om2022.ontologymatching.org/>

- Semantic Web Conference (ISWC), Springer, 2019, pp. 345–362. URL: https://doi.org/10.1007/978-3-030-30796-7_22. doi:10.1007/978-3-030-30796-7_22.
- [2] G. Limaye, S. Sarawagi, S. Chakrabarti, Annotating and searching web tables using entities, types and relationships, *VLDB Endowment* 3 (2010) 1338–1347.
 - [3] D. Ritze, O. Lehmborg, C. Bizer, Matching HTML Tables to DBpedia, in: *Proceedings of the 5th International Conference on Web Intelligence, Mining and Semantics, WIMS, ACM, 2015*, pp. 10:1–10:6.
 - [4] O. Lehmborg, D. Ritze, R. Meusel, C. Bizer, A large public corpus of web tables containing time and context metadata, in: *WWW, 2016*.
 - [5] V. Efthymiou, O. Hassanzadeh, M. Rodriguez-Muro, V. Christophides, Matching Web Tables with Knowledge Base Entities: From Entity Lookups to Entity Embeddings, in: *ISWC, volume 10587, Springer, 2017*, pp. 260–277.
 - [6] E. Jimenez-Ruiz, O. Hassanzadeh, V. Efthymiou, J. Chen, K. Srinivas, SemTab 2019: Resources to Benchmark Tabular Data to Knowledge Graph Matching Systems, in: *The Semantic Web: ESWC, Springer International Publishing, 2020*.
 - [7] E. Jiménez-Ruiz, O. Hassanzadeh, V. Efthymiou, J. Chen, K. Srinivas, V. Cutrona, Results of SemTab 2020, in: *Proceedings of the Semantic Web Challenge on Tabular Data to Knowledge Graph Matching co-located with the 19th International Semantic Web Conference (ISWC 2020), 2020*, pp. 1–8.
 - [8] V. Cutrona, J. Chen, V. Efthymiou, O. Hassanzadeh, E. Jiménez-Ruiz, J. Sequeda, K. Srinivas, N. Abdelmageed, M. Hulsebos, D. Oliveira, C. Pesquita, Results of SemTab 2021, in: *Proceedings of the Semantic Web Challenge on Tabular Data to Knowledge Graph Matching co-located with the 20th International Semantic Web Conference (ISWC 2021), Virtual conference, October 27, 2021, 2021*, pp. 1–12. URL: <http://ceur-ws.org/Vol-3103/paper0.pdf>.
 - [9] D. Vrandečić, M. Krötzsch, Wikidata: a free collaborative knowledge base, *Commun. ACM* 57 (2014) 78–85.
 - [10] R. V. Guha, D. Brickley, S. Macbeth, Schema.Org: Evolution of Structured Data on the Web, *Commun. ACM* 59 (2016) 44–51. URL: <https://doi.org/10.1145/2844544>. doi:10.1145/2844544.
 - [11] S. Auer, C. Bizer, G. Kobilarov, J. Lehmann, R. Cyganiak, Z. Ives, DBpedia: A Nucleus for a Web of Open Data, in: *The Semantic Web, Springer Berlin Heidelberg, 2007*, pp. 722–735.
 - [12] V. Cutrona, F. Bianchi, E. Jiménez-Ruiz, M. Palmonari, Tough Tables: Carefully Evaluating Entity Linking for Tabular Data, in: *19th International Semantic Web Conference (ISWC), 2020*, pp. 328–343.
 - [13] N. Abdelmageed, S. Schindler, B. König-Ries, BiodivTab: Semantic Table Annotation Benchmark Construction, Analysis, and New Additions, in: *Proceedings of the 17th International Workshop on Ontology Matching co-located with the 21st International Semantic Web Conference (ISWC 2021), CEUR-WS.org, 2022*.
 - [14] M. Hulsebos, Ç. Demiralp, P. Groth, GitTables: A large-scale corpus of relational tables, *arXiv preprint arXiv:2106.07258 (2021)*.
 - [15] X. Li, S. Wang, W. Zhou, G. Zhang, C. Jiang, T. Hong, P. Wang, KGCODE-Tab Results for SemTab 2022, in: *Semantic Web Challenge on Tabular Data to Knowledge Graph Matching (SemTab), CEUR-WS.org, 2022*.

- [16] V.-P. Huynh, Y. Chabot, T. Labbé, J. Liu, R. Troncy., From Heuristics to Language Models: A Journey Through the Universe of Semantic Table Interpretation with DAGOBAN, in: Semantic Web Challenge on Tabular Data to Knowledge Graph Matching (SemTab), CEUR-WS.org, 2022.
- [17] M. Cremaschi, R. Avogadro, D. Chierigato, s-elBat: a Semantic Interpretation Approach for Messy taBle-s, in: Semantic Web Challenge on Tabular Data to Knowledge Graph Matching (SemTab), CEUR-WS.org, 2022.
- [18] A. Jiomekong, B. A. F. Tagne, Towards an Approach based on Knowledge Graph Refinement for Tabular Data to Knowledge Graph Matching, in: Semantic Web Challenge on Tabular Data to Knowledge Graph Matching (SemTab), CEUR-WS.org, 2022.
- [19] N. Abdelmageed, S. Schindler, JenTab: Do CTA solutions affect the entire scores?, in: Semantic Web Challenge on Tabular Data to Knowledge Graph Matching (SemTab), CEUR-WS.org, 2022.
- [20] W. Baazouzi, M. Kachroudi, S. Faiz, Yet Another Milestone for Kepler-aSI at SemTab 2022, in: Semantic Web Challenge on Tabular Data to Knowledge Graph Matching (SemTab), CEUR-WS.org, 2022.
- [21] L. Mertens, A low-resource approach to SemTab 2022, in: Semantic Web Challenge on Tabular Data to Knowledge Graph Matching (SemTab), CEUR-WS.org, 2022.
- [22] A. Sharma, S. Dalal, S. Jain, SemInt at SemTab 2022, in: Semantic Web Challenge on Tabular Data to Knowledge Graph Matching (SemTab), CEUR-WS.org, 2022.
- [23] I. Mazurek, B. Wiewel, B. Kruit, Wikary: A Dataset of N-ary Wikipedia Tables Matched to Qualified Wikidata Statements, in: Semantic Web Challenge on Tabular Data to Knowledge Graph Matching (SemTab), CEUR-WS.org, 2022.
- [24] M. Marzocchi, M. Cremaschi, R. Pozzi, R. Avogadro, M. Palmonari., MammoTab: a giant and comprehensive dataset for Semantic Table Interpretation, in: Semantic Web Challenge on Tabular Data to Knowledge Graph Matching (SemTab), CEUR-WS.org, 2022.
- [25] P. Nguyen, I. Yamada, N. Kertkeidkachorn, R. Ichise, H. Takeda, SemTab 2021: Tabular Data Annotation with MTab Tool, in: Semantic Web Challenge on Tabular Data to Knowledge Graph Matching (SemTab), CEUR-WS.org, 2021.
- [26] K. Korini, R. Peeters, C. Bizer, SOTAB: The WDC Schema.org Table Annotation Benchmark, in: Semantic Web Challenge on Tabular Data to Knowledge Graph Matching (SemTab), CEUR-WS.org, 2022.
- [27] A. Jiomekong, C. Etoga, B. Foko, M. Folefac, S. Kana, V. Tsague, M. Sow, G. Camara, A large scale corpus of food composition tables, in: Semantic Web Challenge on Tabular Data to Knowledge Graph Matching (SemTab), CEUR-WS.org, 2022.