From Tabular Data to Knowledge Graphs: A Survey of Semantic Table Interpretation Tasks and Methods

Jixiong Liu^{a,b}, Yoan Chabot^{a,*}, Raphaël Troncy^b, Viet-Phi Huynh^a, Thomas Labbé^a, Pierre Monnin^a

^aOrange, France ^bEURECOM, Sophia Antipolis, France

Abstract

Tabular data often refers to data that is organized in a table with rows and columns. We observe that this data format is widely used on the Web and within enterprise data repositories. Tables potentially contain rich semantic information that still needs to be interpreted. The process of extracting meaningful information out of tabular data with respect to a semantic artefact, such as an ontology or a knowledge graph, is often referred to as Semantic Table Interpretation (STI) or Semantic Table Annotation. In this survey paper, we aim to provide a comprehensive and up-to-date state-ofthe-artreview of the different tasks and methods that have been proposed so far to perform STI. First, we propose a new categorization that reflects the heterogeneity of table types that one can encounter, revealing different challenges that need to be addressed. Next, we define five major sub-tasks that STI deals with even if the literature has mostly focused on three sub-tasks so far. We review and group the many approaches that have been proposed into three macro families and we discuss their performance and limitations with respect to the various datasets and benchmarks proposed by the community. Finally, we detail what are the remaining scientific barriers to be able to truly automatically interpret any type of tables that can be found in the wild Web.

Keywords: Semantic Table Interpretation, Table annotation, Tabular data, Knowledge graph

1. Introduction

Data formats such as CSV/TSV, PARQUET, XML and JSON are commonly used to train machine learning algorithms. We focus on CSV, TSV and spreadsheet files and we argue that while this tabular data format is, compact, readable and simple to process, it does not selfexplain the meaning of the information even if headers are present. Hence, interpreting tabular data becomes a crucial task and it has attracted a lot of attention in recent years, with, in particular, the crystallization of research efforts around challenges such as the SemTab series [32, 53, 54]. The main idea to make tabular data intelligently processable by machines is to find correspondences between the elements composing the table with entities. concepts, or relations described in knowledge graphs (KG) which can be of general purposes such as DBpedia [17] and Wikidata [103], or enterprise specific. This problem is known as Semantic Table Interpretation (STI) or Semantic Table Annotation. KGs can be used to drive the semantic interpretation of tabular data while being themselves the artefacts that can be further enriched from the result of the interpretation process. In this latter case, tabular data becomes a means to either populate a nascent KG or improve the quality of an established one. Adding a semantic

 $Preprint\ submitted\ to\ Elsevier$

layer on top of tabular data, in order to make the latent meaning explicit and exploitable through a structured and shared format, is an invaluable step towards efficient and intelligent use of data. It opens up opportunities for new semantic-based services: leverage semantic annotation to better index datasets in search engines [13, 22], improve question/answering systems [77, 94, 113], enrich knowledge bases [84, 112, 119] or enhance dataset recommendation [117]. The emergence of specialized search engines for datasets such as Google Dataset Search [11, 75] is another prominent example.

Tabular data is challenging to interpret by machines because of the limited context available to resolve semantic ambiguities, the layout of tables that can be difficult to handle, and the incompleteness of KGs in general. Classical Natural Language Processing (NLP) tasks for unstructured text handle poorly such tables since they do not leverage the table structure and the underlying semantics [121]. For example, in the table depicted in Figure 7(b) (page 7), the mention "Rohr" is ambiguous as it can refer to a surname (Q16882196), a manufacturer (Q2391081), or a municipality in Germany (Q583512). However, this ambiguity can be resolved when taking into account the table structure and, in particular, the fact that the "Manufacturer" column only contains companies.

This work aims to comprehensively define the various sub-tasks that belong to STI and to review the many methods that have been proposed so far, along with their lim-

^{*}Yoan Chabot Email address: yoan.chabot@orange.com (Yoan Chabot)

itations and performances on well-established evaluation datasets of the STI community. To the best of our knowledge, such a survey is still missing in the community. Related surveys have recently been published, on Web tables retrieval and enrichment [118], on tables extraction, transformation and understanding [19], on deep learning with tabular data [8, 21] and on information extraction on the Web [62]. However, those surveys do not focus on the STI process. We, therefore, aim to complement these surveys by providing a new categorization reflecting the heterogeneity of tabular data that one can encounter and that yields new challenges.

The remainder of this paper is structured as follows. In Section 2, we describe the research methodology that has been used to make this survey. Next, we provide some preliminaries that are essential for understanding the context of STI and we propose a new fine-grained taxonomy of table types (Section 3). In Section 4, we define five subtasks that are relevant to STI, namely: cell-entity annotation, column-type annotation, columns-property annotation [53, 54], topic annotation, and row-to-instance annotation [84]. We also propose a macro-common pipeline that fulfils the tasks of STI from pre-processing of the input table to the final annotation results. In Section 5, we review the many approaches that have been proposed grouping them into three families (not mutually exclusive) respectively based on heuristics, feature engineering, and deep learning. In Section 6, we make an inventory of the gold-standard datasets commonly used to evaluate STI approaches and we analyse their strengths and weaknesses. We describe the current performances of STI systems on these datasets in Section 7. We elicit the open scientific challenges for the community in Section 8 before concluding the paper in Section 9.

2. Research Methodology

In this section, we describe our research methodology for collecting and to analysing scholar articles relevant to STI tasks, approaches and results.

We observe that most of the work was published between 2000 and 2021. Consequently, our work covers articles published within these two decades. To collect the relevant papers, we primarily make use of Google Scholar. We first use a combination of keywords from two lists: a list of terms reflecting the tasks we are interested in ("semantic annotation", "annotation", "semantic interpretation", "interpretation", "knowledge graph matching", "semantic matching", "semantic labelling", "type prediction", "entity linking", "entity typing", "knowledge graph mapping", "semantic mapping", "relation extraction") and another list of terms scoping the domain ("web table", "table", "tabular data", "structured data"). This list of terms was manually created when reading relevant papers such as the ones from the SemTab challenge¹. For example, a

valid combination is "web table semantic annotation". For each query, we retrieve the 10 most cited papers among the relevant results. We read the abstract of each retrieved paper and judge whether the paper is relevant for STI (e.g., presentation of a method for matching a table element with a given KG, introduction of a new dataset) or not. Then, we extend this initial corpus by retrieving, for each paper previously selected, their 10 most cited papers. The newly selected papers are then filtered according to the methodology used in the first phase. Besides Google Scholar, we use the same keyword-based methodology and co-citation network for finding relevant papers that have been made available on arXiv. This enables us to account for the very latest research work in this fast-moving field even if these papers have not yet been peer-reviewed or cited. Finally, we have also collected the papers from the SemTab challenge which is the most relevant competition for this domain.



Figure 1: Distribution of the publication years of the papers selected for this review.

Using this method, we generated 48 combinations of keywords to search on Google Scholar and we selected 38 distinct references in the first phase. The co-citation network has brought 58 additional papers in the second phase. We added 7 arXiv papers and 22 SemTab papers to this selection. In summary, while more than seven hundred papers have been collected and read, 114 papers have been assessed to be relevant in this two-phase selection process. Figure 1 depicts the time distribution of the publication of these articles. It shows the growing importance of STI in the literature over the last decade. We manually tagged each article into six categories (Figure 2). Note that some articles may address more than one research focus. For example, [60] provides both a dataset and a method for performing STI. The core of this survey is focused on STI systems but it also covers tangentially related topics such as the need for knowledge graphs and schemas to anchor the interpretation of tabular data or the importance of pre-processing tables, using NLP techniques for example.

¹https://www.cs.ox.ac.uk/isg/challenges/sem-tab/



Figure 2: Distribution of the topics of the papers selected for this survey.

3. Preliminaries

The first and main input of an STI system is the table itself. As there is an important heterogeneity of tables considering their layout, provenance, and usage, we propose a new fine-grained classification based on existing classifications with a deeper analysis of relational tables (Section 3.1). This classification of tables is intended to make it easier to define the scope of STI approaches proposed in the literature and to help identify the challenges related to STI tasks.

Tables do not always appear alone in real-life scenarios. Alongside the data contained in a table, metadata and the context in which a table appears are also valuable information for STI. For example, if a table has been published on a Web page describing the Bundesliga, it is probably more relevant to football than any other sport. Hence when extracting the table, it would be useful to collect both the table itself and its metadata. Section 3.2 provides a list of elements attached to the table carrying semantic information useful for interpretation.

Finally, STI uses KGs as sources of information and as references for producing annotations. Section 3.3 highlights the most commonly used KGs for table interpretation and discusses their specificities and what they imply for the STI tasks.

3.1. Tables

A table is a two-dimensional arrangement of data with n lines and m columns. This enables a compact visualization for reading. A cell is the basic element of a table where \mathcal{T}_{ij} ($0 \le i \le n-1, 0 \le j \le m-1$) indicates the cell from row i and column j of table \mathcal{T} . Tables are highly heterogeneous in terms of structure, content, and purpose. Therefore, before interpreting a table, it is important to identify its type so that potential specificities can be taken into account in the STI process.

We introduce a multi-level classification of tables based on several aspects. The first classification effort splits tables into two high-level categories: genuine and nongenuine [78, 106]. Genuine tables are firstly defined as two-dimensional structures with simple cells (i.e. short and without any complex structures) and a high level of coherence (syntactically and semantically) within rows and columns in [78]. Non-genuine tables are structures used to group contents for easy viewing [106]. One limitation of this dichotomy is that it does not consider tables with long and complex cell contents which are still semantically coherent. For example, we can observe cells containing a list of comma-separated entities (row "Celebration") or mixing text and entities (row "Significance") in the infobox depicted in Figure 6(b). According to the definition provided in [78] and used in [82], this table will not be considered a genuine table while, arguably, this table carries semantic information worth to be processed.

In more recent works, [30, 57] proposed two similar classes associated with a set of sub-classes: relational knowledge tables including vertical and horizontal listings, attribute/value table, matrix, etc. and layout tables including tables used for navigation and formatting purposes. This classification focuses mostly on relational knowledge and is therefore not comprehensive enough to cover all possibilities. For example, some tables do not have interrelation between table elements and are not for layout purposes either. This is the case of the table depicted in Figure 6(d) where there is no information about the common relationship between each cell.

In order to cope with these shortcomings, we propose a new classification of table types, shown in Figure 3, that rely on the existing work presented in [35, 57, 58, 78, 82, 108, 118] and consider overlapping dimensions. This classification also contributes to a better identification of relational tables in embracing their diversity.



Figure 3: Classification of table types with a finer-grained analysis of genuine tables along three dimensions: structure, inner-relationship and orientation.

We first consider that tables can be separated into two broad categories. **Layout tables** are used to format Web pages. Elements of these structures are not semantically consistent and are not linked by semantic relationships. They are used to visually organise the content of a page in order to maximise user comfort and site usability. A layout table used on Amazon to provide the order interface is given in Figure 4. Genuine tables represent in rows or columns human-understandable knowledge. In the literature, genuine tables are said to be short and without complex structure [78]. We relax this concept by considering that tables with a high level of coherence (syntactically and semantically) within rows and columns are genuine tables, without considering the table complexity. For example, in Figure 4, the semantic of the two genuine tables is the description attributes (e.g., price, color) of the product. Genuine tables contain relational knowledge that should be machine-interpretable, and thus, they constitute valid inputs for the STI process. On the contrary, layout tables are convenient for improved visual presentation but the semantic association between their cells is relatively sparse. Hence, they are not eligible for knowledge extraction and interpretation.



Figure 4: Illustration of genuine tables and layout tables on the Amazon website. a

$^{a} \rm https://www.amazon.com/Furinno-14035EX-Study-Table-Espressodp/B00NIYX9LC$

Previous works have also proposed to classify tables starting from different entry points and further segment genuine tables along different dimensions [35, 57, 82, 118]. However, the state of the art considers that these table types are mutually-exclusive which does not cover the heterogeneity and complexity of genuine tables. Consequently, we propose to categorize genuine tables using three nonmutually exclusive dimensions: structure, inner relationship, and orientation. Table types are then formed by a composition of these dimensions. For example, the table depicted in Figure 5(c) about railway lines is a concise table (structure dimension), a horizontal table (orientation dimension), and a composed-subject relational table (inner-relationship dimension). In the following sections, we further define each of these three dimensions.

3.1.1. Structure Dimension

[57] focuses on the layout structure of a table, which is mainly reflected in the table elements' composition. Accordingly, the subclass "Structure Dimension" of our classification is divided into the following four types of tables illustrated in Figure 5. Nested tables contain one or more tables in one or more of their cells. Figure 5(a) depicts a nested table as the main table contains a table about risk levels of hazardous materials in one of its cells. Split tables contain sub-tables with cells that are independent of those in the other sub-tables. [57] defines split tables as a sequential repetition of rows or columns. We enforce this definition by defining split tables as tables that can be split into sub-tables. To illustrate, in the table in Figure 5(b), the infobox of the city of Chicago is composed of several sub-tables. Each sub-table describes one concept of the subject of the original table, e.g., location, area, and population of the main Chicago entity. Concise tables contain merged cells in order to avoid repetitions of cells referring to the same content in rows and/or columns. In the table presented in Figure 5(c), the first cell of the column "lines" merges six individual cells with the same value "BART main lines". Multivalued tables contain multiple values in a single cell. For example, in Figure 5(d), the cells of the column "lines used" contain a list of route lines.

3.1.2. Inner-relationship Dimension

The inner-relationship dimension considers the topology of the semantic connection between the cells. [30] first gave a detailed classification of relational knowledge tables into listings, attribute or value tables, matrices, enumerations, and forms. [35] has extracted Web tables that are classified into listings, matrices and other tables to capture enumerations, calendars, etc. [58, 82] have extended this categorization with a fourth kind named "entity tables" while listings are considered as relational tables.

Accordingly, we propose the following types. Relational tables are structures in which each row (resp. column) provides information about a specific entity, and the corresponding columns (resp. rows) represent attributes that describe the entity. Hence, relational tables are oriented, either horizontally or vertically, depending on the arrangement of the entities and their attributes in the table. If the rows of a relational table contain the entities and the columns the attributes, then the table is horizontal. Otherwise, it is vertical. Relational tables may have a header, usually in the first row or the first few rows for horizontal tables. As an example, Figure 6(a) depicts a horizontal relational table where each row describes a tower with its attributes (e.g., height, year, country, and town). Entity tables, also known as attribute-value tables, are used to describe a unique entity. An entity table enumerates the attributes of the entity. Infoboxes from Wikipedia are examples of entity tables. For example, the distinct attributes and their values for the entity "Bastille Day" are shown in the table in Figure 6(b). It should be

	Fire dian	mond																							
Standard Representation	n	Tabular R	eprese	entation		Lines		Manufactu	irer	Class	Image	Number	Car numbers	Built	Notes										
	Risk levels o	of hazardo	us ma	terials in th	is facility			Rohr		A		59	1164-1276	1968-1975											
321	Health Risk	Flammal	bility	Reactivity	Special			Rohr		в		380	1501-1913	1971–1975											
\checkmark	Level 3	Level	2	Level 1	\]			Alstom		C1	I THE REAL PROPERTY	150	301-450	1987-1989	To be phased out by August 2023 and replaced by the "D" and "E" cars.										
Country State Counties	Multicol States Illinois Cook, DuPage circa 1780 March 4, 1837 Harch 4, 1837 Jean Baptiste Point du Sable Mayor-council Chicago City Council Lori Lightfoot (D) Anna Valencia (D) r Melissa Conyears-Ervin (T)			(a)		BART main lines		Morrison-Knuds		C2		80	2501-2580	1994–1996 ⁽⁶⁶⁾	- a										
Settled Incorporated (town) Incorporated (city) Founded by			circa 1780 August 12, 1833 March 4, 1837 Jean Baptiste Point du Sable Mayor-council		1-)				Bombardier		D		310	3001-3310	2012-	Order being filled/testing, entered service on January 19, 2018.								
Government • Type					Mayor-council		Mayor-council		Mayor-council		Mayor-council		(b)		(c)			Bombardier	ombardier			465	4001-4465 2012-		orosi osing inclutesung, entered service of January 13, 2010.
• Body • Mayor • City Clerk • City Treasurer						Oakland Airport O	Connector	DCC Doppelr	CC Doppelmayr Cable			4	1.3-4.3	2013	automated guideway transit trainsets										
Area ^[3] • City	234.21 sq mi			(d)	eBART		Stadler		GTW	N.C.	8	101-108	2014-2018	diesel multiple units										
1 and	(606.60 km ²)		_	(4	/																				
• Land	(588.98 km ²)		_		Route nan	ne	Firs	irst service		Lines used		Service times													
• Water	6.80 sq mi (17.62 3.0%	5.80 sq mi (17.62 km ²) 3.0%		Berryessa/N	orth San José–	Richmond line	September 11, 1972		R-Line, K-Line, Line, S-Line		Operates	Operates during all service hours.													
• Urban • Metro	2,122 sq mi (5,496 10,874 sq mi (28,160 km ²)	6 km ²)	2) Antioch–SFO/Millbrae line)/Millbrae line	May 21,		May 21, 1973 Li		Ne, K-Line, M W-Line, Y-Li RT	ne, Through-	Through-routed with the SFO-Millbrae line on Sundays. Uses DMU technology from Antioch to Pittsburg/													
Elevation ^[2] (mean)	597.18 ft (182.02	m)		Berryessa/N	orth San José–	Daly City line	Septem 1974	iber 16,	S-Line, A-Line, M- Line		No eveni	ng or Sund	day service.												
- near Blue Island	672 ft (205 m) 578 ft (176 m)			Richmond-M	Aillbrae line		April 19	9, 1976	R-Lin Line,	ie, K-Line, M W-Line	Terminate	es at Millb	rae on weekda	ys and at Daly	City on Saturdays; no evening or Sunday service.										
- at Lake Michigan				Dublin/Pleas	anton–Daly Ci	ty line	May 10	, 1997	L-Line	e, A-Line, M	Operates	during all	service hours.	Some Sunday	service terminates at Montgomery Street station.										
• City • Estimate (2019) ^[7]	2,695,598 2,693,976			SFO-Millbra	e line	February (previou 2004)		ry 11, 2019 usly 2003–	W-Lir	ne, Y-Line	Through-	routed wit	th the Antioch-	-SFO/Millbrae l	ine on Sundays.										
Density Urban Metro CSA	11,846.55/sq mi (4,573.98/km ²) 8,667,303 ^[5] 9,533,040 (3rd) ^[4] 9,901,711 (US: 3rd	46.55/sq mi '3.98/km ²) 7,303 ^[5] 3,040 (3rd) ^[4] 1,711 (US: 3rd) ^[4]		Coliseum-Oo	tional Airport line	nal Airport line November 22, 2014		Separate elevated automated guideway transit line (H-Line) not connected to other BART tracks		ine Operates er	Operates during all service hours.														

Figure 5: Examples of different structures of tables: (a) a nested table^{*a*}, (b) a split table^{*b*}, (c) a concise table (column "lines")^{*c*}, (d) a multivalued table (column "lines used")^{*d*}.

^ahttps://en.wikipedia.org/wiki/Table_(information)

^bhttps://en.wikipedia.org/wiki/Chicago

^chttps://en.wikipedia.org/wiki/Bay_Area_Rapid_Transit#Rollingstock

^dhttps://en.wikipedia.org/wiki/Bay_Area_Rapid_Transit#Infrastructure

noted that entity tables could also be seen as two-column vertical or two-row horizontal relational tables. Matrix tables present a two-dimensional arrangement of data that should be read simultaneously horizontally and vertically. A matrix associates pairs (row, column) with cell values through one unique property for the whole table. Generally, cells contain numeric or boolean values. Figure 6(c)shows a vowel confusion matrix that quantifies the understanding of vowels between people. It associates pairs (vowel produced, vowel perceived) with the number of persons having this perception of the produced vowel. For example, one person perceived an "a" when an "i" was produced. Bold numbers correspond to correct identifications. Other genuine tables contain semantic information but do not fit within the aforementioned types. Tables in this class include enumerations and calendars. To illustrate, Figure 6(d) is an enumeration table where each column is an independent enumeration of pronouns according to a pronoun type.

The literature considers relational tables as a leaf in the proposed table type taxonomies [58, 82, 108]. However, relational tables exhibit an important diversity, especially in the representations of entities. We propose to further classify them depending on the characteristics of their subjects. The **subject** of a row of a horizontal relational table (resp. column of a vertical relational table) is an entity that is described by the collections of cells in this row (resp. column). For example, in the table depicted in Figure 6(a), the entity "Tokyo Skytree" (Q57965) is the subject of the first row as it is described by the other entities of this row: "2011" (Q1994), "Japan" (Q17), and "Tokyo" (Q1490).

We introduce four subtypes of relational tables (Figure 7). **Single-cell-subject** tables associate each row of a horizontal table (resp. column for a vertical table) to a single subject. Labels of subjects are given in a single column (resp. row). To illustrate, in Figure 7(a), the column "department" contains the subjects. The other columns describe the subjects. **Composed-subject** tables require the combination of multiple cells to form the subject of each row (resp. column). For instance, in a table that describes persons with first and last names in different

	Bastille Day	Name	Pinnacle height		Year	Country	Tov		vn		marks	
1000		Tokyo Skytree		634 m (2,080	O ft)	2011	Japan	Tokyo	Tokyo			
1000		Kyiv TV Tower		385 m (1,263	53 ft) 197		Ukraine	Kyiv				
		Dragon Tower		336 m (1,102	2 ft) 2000		China Harbin					
1 mars	Emile -	Tokyo Tower		333 m (1,093	3 ft)	1958	Japan	Tokyo				
	A	WITI TV Tower	329.4 m (1,0	81 ft)	1962	United States	Shorewo	ood, V	Visconsii	ı		
to have one		St. Petersburg T	326 m (1,070) ft)	1962	Russia	Saint Peters		sburg			
Fireworks Also called	at the Eiffel Tower, Paris, 2017 French National Day		(a)			Perceive Produced	d i	е	а	ο	u
	(Fête nationale)						i	15		1		
	The Fourteenth of July (Quatorze juillet)	(b)			((c)	е	1		1		
Observed by	France		а			79	5					
Туре	National day						ο			4	15	3
Significance	Commemorates the Storming of the Bastille on 14, July 1789 [1][2]		(d)			u				2	2
	and the unity of the French	Demonstrative	elative			Indefinite		Interrogative			ve	
	people at the Fête de la	this	who / w	hom / whose	one /	one's	/ oneself		١	who / wh	om / \	whose
Celebrations	Hederation on 14 July 1790	these	what		some	ething	/ anything / noth	ing (thing	s) I	what		
Celebrations	concerts, balls	that which			some	omeone / anyone / no one (people)				which		
Date	14 July	those	that		some	ebody	/ anybody / nob	ody (peop	ole)			
Frequency	Annual	former / latter										

Figure 6: Examples of different inner-relationships of tables: (a) a horizontal relational table^{*a*}, (b) an entity table^{*b*}, (c) a matrix table^{*c*}, (d) an enumeration table^{*d*} which belongs to "other genuine tables".

^ahttps://en.wikipedia.org/wiki/Eiffel_Tower

^bhttps://en.wikipedia.org/wiki/Bastille_Day

 $^{c} https://en.wikipedia.org/wiki/Whistled_language\#Lack_of_comprehension$

^dhttps://en.wikipedia.org/wiki/Pronoun#English_pronouns

columns, it is necessary to merge these two columns to get the complete identifiers of entities. Similarly, in the table shown in Figure 7(b), one can identify subjects (particular train classes) by merging columns "Lines", "Manufacturer" and "Class". **Multi-subject** tables contain cells that refer to different subjects while being in the same row. In Figure 7(c), a row is composed of two subjects: "Artist/s" is the subject of the column "Nationality" while "Album" is the subject of columns "Release year", "Artist/s", "Worldwide sales", and "Ref(s)". **Hiddensubject** tables do not explicitly mention the subject of each row (resp. column). For example, in Figure 7(d), each row describes the result of a football match, but the mention of the match itself is not made explicit in the table.

3.1.3. Orientation Dimension

The orientation dimension considers the direction of the relationships inside a table [35, 57]. Indeed, knowing the direction of relationships within a table simplifies its interpretation, e.g., to read the attributes describing a subject.

In **Horizontal** tables, subjects are described horizontally, which means that each row describes a different subject. For example, in Figure 6(a), the subject "Dragon Tower" and its attribute "Harbin" are in the same row. In **Vertical** tables, subjects are described vertically, which means that each column describes a different subject. An example of a vertical table is depicted in Figure 6(b), where the attributes of the entity "Bastille Day" are in the same column. **Matrix** tables are defined as for the inner-relationship dimension. They cannot be interpreted row by row or column by column but rather cell by cell while simultaneously considering both horizontal and vertical headers. For example, the matrix in Figure 6(c) should be interpreted cell by cell while taking into account both horizontal and vertical headers to read the number of persons that have a specific perception of a produced vowel.

3.1.4. Table Types Statistics

We introduced in Figure 3 a multi-dimensional and fine-grained classification of table types. In this section, we aim to survey how frequent these types of tables are used in the wild. DWTC [35] has randomly selected 26,645 tables from the WDC Web Table Corpus [58] and has concluded that the resulting corpus was made of 96% layout tables and 4% genuine tables. Regarding the structure di-

(a)	D	epartment 🕈	Area (km ²) 🔹	Popu (201	1) ^[37] •	Municipaliti	es 🕈	(b)	Lines	Manufacturer	Class	Image	Number	Car numbers	Built
	Paris (75	Paris (75) 105.4 2 249 975		1 (Paris)						- / 103		1164-	1968-		
	Hauts-de-Seine (92) Seine-Saint-Denis (93)		176	1 581 628		36 (list)				Kohr	A		59	1276	1975
			236	15	29 928	40 (list)									
	Val-de-N	larne (94)	245	13	33 702	47 (list)				Rohr	в		380	1501-	1971-
	Petite Co	ouronne	657	44	45 258	123								1913	1975
	Paris + F	Petite Couronne	762.4	6 695 233		124									
(c)	Release 🖕	Album •	Artist/s	• 1	Nationality	Worldwide sales	Ref(s)	٠	BART main lines	Alstom	C1		150	301–450	1987– 1989
	2002	Come Away With Me	Norah Jones		United State	23.9	[3]			Morrison-	~		00	2501-	1994–
	2000	2000 The Marshall Mathers LP Eminem United		United State	\$ 23.29	[4]			Knudsen	C2		80	2580	1996 ^[67]	
	2002 The Eminem Show		Eminem	United States		s 22.95	[5]								
	2000	Hybrid Theory	Linkin Park	United States		s 20.8	[6]							3001-	
	2015	25	Adele	***	United Kingdo	m 20.41	101			Bombardier	D		310	3310	2012-
(d)	3	11 juin 1998 Italie		2	2 - 2 📕 Chili		Chili							4001-	
()	4	11 juin 1998	Cameroun	1	- 1 🚍	Autriche				Bombardier	E		465	4465	2012-
	19	17 juin 1998	Chili 🛯	1	-1 🚍	Autriche			Oakland Airport	DCC	Cable		Д	1 3-4 3	2013
	20	17 juin 1998	Italie 🛛	3	- 0 🔳	Cameroun			Connector	Doppelmayr	Liner	الدراجة	-	1.5 4.5	2010
	33 23 juin 1998 Italie 2 - 1 Autric		Autriche			-DADT	Charles	CTM		0	101 100	2014-			
	34	23 juin 1998	Chili 🖪	1	- 1 🔳	Cameroun			EDANI	Statier	GIW		0	101-108	2018

Figure 7: Examples of different relational tables: (a) a single-cell-subject relational table^a, (b) a composed-subject relational table^b, (c) a multi-subject relational table^c, (d) a hidden-subject relational table^d.

 ${}^{a} \rm https://en.wikipedia.org/wiki/France\#Major\%20 cities$

^bhttps://en.wikipedia.org/wiki/Bay_Area_Rapid_Transit

 $^{c} \rm https://en.wikipedia.org/wiki/List_of_best-selling_albums_of_the_21st_century$

 ${}^{d} https://fr.wikipedia.org/wiki/Coupe_du_monde_de_football_1998$

mension, [57] has extracted 342,795 Web tables (from various websites starting from Wikipedia, e-commerce, news and university Web sites) and has identified that 75.5%of the tables are for layout while the remaining tables are relational knowledge tables. [57] has also provided the distribution of nested tables, split tables, concise tables, and multivalued tables among the relational knowledge tables, which are respectively 3.7%, 2.6%, 12.9%, and 74.9%. Regarding the inner-relationship dimension, [58] applied the DWTC framework on the 233 million Web tables of the WDC Web Table Corpus to detect the type of each table w.r.t the inner-relationship dimension. Results show that relational tables, entity tables, and matrices respectively constitute 39%, 60%, and 1% of the corpus. Regarding the orientation, the distribution of horizontal tables and vertical tables is 54.9% and 45.1% for entity tables, and 94%and 6% for relational tables in the WDC Web Table Corpus. In [57], the authors show that 70% of the relational knowledge tables are horizontal.

Our proposed classification goes further into the details. However, identifying some table type such as hiddensubject tables remains an open scientific challenge. To date, we have not identified an approach to automatically classify tables with a level of granularity close to the classification proposed in this paper.

3.2. Metadata

Several STI works stress that one of the challenges to be addressed is the loss of context when annotating a table [118]. Indeed, tables do not constitute the unique source of information that can be used by STI processes since the context in which they appear may provide complementary or novel information. Such non-table information constitutes the metadata of tables and is defined as additional data that can be extracted from information sources to provide additional context for the interpretation. For example, metadata can describe the characteristics and content of the original data, and thus can be used to organize, retrieve, preserve, and manage extracted knowledge units. Depending on its structure, purpose, and provenance, metadata is split into descriptive, structural, and administrative metadata [80]. Such a definition was originally used in digital collection [115] and is applicable to table metadata as well. Each type of metadata in a table context can provide a deeper understanding of the table

Descriptive metadata is used to describe the target data by providing, e.g., its source, explanatory notes, or other contextual information. For example, the descriptive metadata of the table "Lattice towers taller than the Eiffel Tower" depicted in Figure 8(a) can include its provenance (i.e. the URL of the page) and its surrounding text. Indeed, texts surrounding tables are potential sources of

aller structures						Year
the Eiffel Tower was the w when the Chrysler Building vorld's tallest tower to the	vorld's tallest structu g in New York City w Tokyo Tower in 195	re whe as topp 8 but re	n completed in bed out. ^[107] The tains its status	1889, a distinction it retai e tower also lost its stand as the tallest freestanding	ned until 1929 ing as the a (non-auved)	
tructure in France.					3 (3-)/	Remarks
attice towers taller	than the Eiffel 1	ower				<pre><td< th=""></td<></pre>
Further information: List	st of tallest towers in	the wo	orld, Lattice tow	er, and Observation deck		2011 Japan
Name	Pinnacle height	Year	Country	Town	Remarks	Tokyo
Tokvo Skytree	634 m (2.080 ft)	2011	Japan	Tokyo		
Kyiv TV Tower	385 m (1,263 ft)	1973	Ukraine	Kyiy		385 m (1,263 ft) 1973
Dragon Tower	336 m (1,102 ft)	2000	China	Harbin		Vkraine Kyiv
Tokyo Tower	333 m (1,093 ft)	1958	Japan	Tokyo		<
WITI TV Tower	329.4 m (1,081 ft)	1962	United States	Shorewood, Wisconsin		Dragon Tower 336 m (1,102 ft)
St. Petersburg TV Tower	326 m (1,070 ft)	1962	Russia	Saint Petersburg		2000 2000 >>>>>>> <td< td=""></td<>
$\overline{(c)}$	1			1		> >
(0)						<pre>:</pre>
13:21, 8 April 2009 (diπ his 18:15, 27 March 2009 (diff I	t)(-188) <u>m</u> Ed Sil bist) (-141) Penta	mmons	(Reverted edits b Projects: removing	y 96.25.55.43 (talk) to last vi soam magnet)	ersion by ClueBot)	3338#160;m (1,0938#160;ft)
18:15, 27 March 2009 (diff I	hist) (-42) m Mor	ndrian O	LAP server (rem	oving advertising link)		Japan
16:14, 19 March 2009 (diff)	hist) (+347) User	talk:And	re Engels (→Rfa	nom)		lokyo
11:06, 19 March 2009 (diff)	hist) (+8) m User	talk:Car	lossuarez46 (→F	RedirectName listed at RfD)		<
11:05, 19 March 2009 (diff	hist) (+481) User	talk:Car	ossuarez46 (RFI	DNote)		>>>>>><
11:05, 19 March 2009 (diff	hist) (+241) Wikip	edia:Re	directs for discuss	sion/Log/2009 March 19 (No	minating Talk:Neftchala	>29.4dm (1,081 Tt) 1962
11:03, 19 March 2009 (diff	hist) (+8) Talk:Net	tchala ({{rfd}}})		- /	United States
11:00 19 March 2009 (diff I	hist) (+102) Neftc	hala (re	direct changed to	disambiguation page)		
11.00, 10 march 2000 (am)				/		(lt d) (lt -)

Figure 8: Metadata of the Web table "Lattice towers taller than the Eiffel Tower"^{*a*}: (a) descriptive metadata: its surrounding text, (b) structural metadata: , and tags indicate table cells, row ordering, and the presence of headers, (c) administrative metadata: the page history of the table^{*b*}.

^ahttps://en.wikipedia.org/wiki/Eiffel_Tower

^bhttps://en.wikipedia.org/w/index.php?title=Special:Contributions/Andre_Engels

contextual information, and thus valuable metadata, since they often explain a nomenclature or verbalize salient information. The different relationships between texts and tables, including titles and captions or even simple cooccurrences between a table and the surrounding texts, are useful indicators to guide and improve the annotation and knowledge extraction processes. However, this tabletext complementarity is little used in the STI domain so far. Descriptive metadata provides additional information that enhances the process of approaches such as [15, 34].

Structural metadata describes the structural schema of composite objects or relationships between objects. For example, the $\langle td \rangle$, $\langle tr \rangle$, and $\langle th \rangle$ tags in the Web table in Figure 8(b) allow to detect table cells, row ordering and the presence of headers. Such structural patterns could benefit STI approaches such as [97, 105].

Administrative metadata often captures information such as the process of creation or data acquisition for a table. For example, in Figure 8(c), the page history details how, when, by whom, and for which purpose data has been produced or altered, allowing to assess the quality and validity of the table. Additionally, the creation or modification date of a table can indicate the freshness of its information, and thus allows to assess the risk of extracting outdated information.

It should be noted that table metadata can appear in different forms since tables have different formats and structures. For example, in some approaches, table headers are available as metadata [24]. Furthermore, if a table element consists of a hyperlink (e.g., hyperlinks in infoboxes of Wikipedia), this mapping relationship also constitutes metadata. STI systems that leverage metadata as an input source are identified in Table 1.

3.3. Knowledge Graphs

Knowledge graphs are often associated with linked data technologies and projects since they focus on interrelations between concepts and entities [38]. In essence, KGs are semantic networks that formally describe things or entities of the real world and their relationships [45]. Each entity is identified by a globally unique URI [12]. Atomic elements of KGs are triples (subject, predicate, object). For example, the population of Paris can be represented by the triple (Paris, hasPopulation, 2m) where the predicate hasPopulation qualifies the relationship holding between the entity Paris and the value 2m.

KGs can be categorized into domain-specific (or vertical) KGs, encyclopedic KGs or common-sense KGs depending on their content. A domain-specific KG focuses on describing a particular field of interest. Such a KG is expected to present advantages in terms of accuracy and in-depth domain knowledge coverage. It can effectively support knowledge reasoning and knowledge retrieval for specific domain applications [111]. To illustrate, in the "Crop, Pest, and Diseases" field, domain-specific KGs play a substantial role in agriculture [99], greenhouse environment [114], and economic benefit analysis [100]. The Bio2RDF project which focuses on Life Sciences [10] is an example of a domain-specific KG that is largely used.

An encyclopedic KG is generally large, spanning multiple domains, and is often openly and collaboratively edited. This openness is reflected by the Linked Open Data cloud [16] which includes the following largest encyclopedic KGs. **DB**pedia [17] is one of the main hubs of the LOD cloud because of its numerous interlinks with other KGs. It was created by researchers from the University of Leipzig and the University of Mannheim in Germany by extracting multilingual structured data from Wikipedia (e.g., infoboxes). It is maintained up-to-date thanks to frequent extracts from Wikipedia. As of early 2016, DBpedia contained more than six million instances and 200 million facts. Moreover, the DBpedia project provides tools such as DBpedia Spotlight [63] that are convenient for mapping mentions contained in unstructured data with KG entities. Wikidata [103] is a project hosted by the Wikimedia Foundation which aims to fuel the infoboxes displayed on each Wikipedia page. Similarly to Wikipedia, it is collaboratively edited by thousands of volunteers. As of early 2021, Wikidata had facts about 90 million entities with labels expressed in more than 350 languages. Wikidata provides a separate page for each entity, has a unique digital identification mechanism, and a lineage system that allows to trace facts to their sources. Freebase [18] was developed by Metaweb since 2007 until Google acquires it in 2010. Its content comes also from collaborative editing and structural data automatically imported from Wikipedia and other websites. At the beginning of 2014, Freebase had 68 million entities and nearly one billion facts. Freebase ceased operations in May 2015, and most of its data was transferred to Wikidata. YAGO (Yet Another Great Ontology) [96] is a comprehensive knowledge base constructed by researchers from the Max Planck Institute (MPI). While the first versions of YAGO were made out of information extracted from the Wikipedia infoboxes and attached to a schema made of the WordNet synsets, the latest version of YAGO now contains entities extracted from Wikidata anchored to the Schema.org schema. In 2020, YAGO released its fourth version containing 67 million entities and 340 million facts.

KGs constitute essential assets to support the STI process. The column KG of Table 1 provides the specific KGs which have been used in each STI system reviewed in this paper. Indeed, understanding the content of a table comes down to identifying the entities mentioned in the table cells and the relationships between them. Therefore, mapping table content to KG entities can help identify latent relationships, and thus understand the table semantics. The key to map tables and target KGs is to examine the overlap of information between them. The wider the overlap is, the less difficult it is to find the mappings.

It should be noted that each KG may present its own (dis)advantages to support the STI process. For example,

Wikidata provides rich content and numerous aliases for each entity to cover a wide set of real-world synonyms. However, annotating a cell with such an encyclopedic KG based on a string-similarity mapping may lead to a significant number of candidates due to the presence of numerous homonyms. This would, in turn, make disambiguation more challenging. For example, over a hundred entities have labels or aliases that contain the word "France". The Wikidata data model is also complex as it provides qualifiers that may need to be specifically taken into account during the STI process. On the other hand, DBpedia provides a reduced number of types curated in the DBpedia Ontology, which makes the typing of table elements easier but potentially reduces as well the specificity of the annotations. Regarding vertical KGs, the lack of knowledge from other domains can lead to a reduced system generalization. Additionally, string matching may not be able to handle the specificities of sophisticated domain-specific relations, schema, or entities, increasing the interpretation complexity. For example, in the biology domain, genes and proteins often share the same labels. To guide the choice of the supporting KG, it is noteworthy that selecting KGs with the highest overlap with the dataset's content will maximize the system's performance. Besides, combining several KGs will maximize the coverage and granularity.

4. Annotation Tasks and Pipeline

The previous section introduced the type of tables to annotate and the KGs generally used for the annotation. In this section, we focus on the STI process. We first present five STI tasks in Section 4.1. Next, we describe a common pipeline to perform STI tasks in Section 4.2.

4.1. Annotation Tasks

An annotation task can be defined by the table elements required to be annotated and by the type of candidates (individuals, concepts, or properties of the KG). We propose to decompose STI into five main tasks: cell-entity annotation, column-type annotation, columns-property annotation [53], topic annotation, and row-to-instance [83] (illustrated in Figure 9).

Cell-Entity Annotation (CEA) is also known as Entity Linking. It aims to annotate a cell with a KG entity. For example, in Figure 9, a CEA task needs to match the cell mention "Suisse" with the entity Q165141 if Wikidata is the target KG. **Column-Type Annotation** (CTA) aims to map a column with a KG entity type. The difficulty of the CTA task lies in selecting an adequate type granularity in a potentially complex type hierarchy structure. An entity may have multiple types and types represented in complex hierarchical trees or even cyclic graphs (e.g., Wikidata type topology). The type selected for a given column must be representative of the individuals it contains and carry a maximum of information. If the selected type is too broad (e.g., the second column of the



Figure 9: Illustration of five STI tasks for a table describing the UEFA Euro 2008 group A results^a.

^ahttps://fr.wikipedia.org/wiki/Championnat_d%27Europe_de_football_2008#1er_tour_-phase_de_groupes

table in Figure 9 is annotated as a "geographic entity" (Q27096213) in Wikidata rather than "city of Switzerland" (Q1545591)), the annotation will carry little information. Conversely, a type that is too specific may not be representative for all values in a column, leading to an accuracy degradation in downstream tasks. In Figure 9, the label "city of Switzerland" (Q14770218) would no longer be compatible with the second column if other groups and cities that hold UEFA Euro 2008 games such as Vienna (Q1741) are included in the table. Columns-Property Annotation (CPA) aims to annotate a column pair of a relational table with a property. For example, the relationship between the last column and the numeric column circled in orange in Figure 9 should correspond to the predicate "number of points/goals/set scored" (P1351) in Wikidata. Topic annotation aims to annotate the entire table with a concept or an entity from the target KG. Figure 9 illustrates that the entire table is about the entity "UEFA Euro 2008" (Q241864) in Wikidata. Rowto-Instance annotates an entire row of a relational table with a KG entity. In this task, each row is treated as an entity, which is considered the subject of the row. Rowto-instance differs from the CEA task as it may be able to discover more entities leveraging the context of the row, especially in the case where the subject of the row is hidden (e.g., on hidden-subject tables). For example, in Figure 9, the fourth row is represented by ("Switzerland - Turkey, 11 Jun 2008" (Q12012827)) which can not be extracted by the CEA task.

4.2. Annotation Pipeline

The study of STI approaches in the literature allows identifying a recurrent pipeline of modules used by the vast majority of systems. This section introduces a macroscopic view of a four-stage STI pipeline: pre-processing (Section 4.2.1), candidate generation (Section 4.2.2), table elements processing (Section 4.2.3), and iterative disambiguation (Section 4.2.4). It should be noted that the order of these modules may vary from one approach to another. In addition, some approaches iterate between the annotation tasks and disambiguation stages to improve the accuracy of annotations [84, 121].

4.2.1. Pre-Processing

The quality of the STI output is greatly influenced by the quality of the input data. A pre-analysis of the data is, therefore, necessary and is often the first step of an efficient STI system. For example, knowing the orientation of a relational table can help to identify a column and its header information to disambiguate the cells.

The goal of the pre-processing module is to quickly and concisely summarize and analyse an input table to ease the annotation process by converting it into an interoperable format. The pre-processing analysis can be decomposed into the following two tasks. First, format normalization allows to transform the original data into a format acceptable to an STI system. Indeed, table sources are diverse, so does their representation in terms of formats (e.g., CSV, JSON, HTML), charsets being used (e.g., UTF-8, Unicode), languages (e.g., English, French), or content expressions (e.g., missing value). This task also aims to clean up invalid or erroneous data for better compatibility among different data sources, e.g., with syntactic corrections by fuzzy matching in [25, 73]. Second, informational analysis consists of extracting the potential information contained in the table as much as possible before the annotation. The information carried by a table includes the table types (e.g., relational, matrix) [58, 118], its orientation [23, 43], primitive types for its columns (resp. rows), its header positions, if any, and its key column position which carries the row's subject [23]. Such information helps the system to understand the scenario and to perform different operations in different situations. For example, the CTA task is only suitable for relational tables, and it is related to the table orientation: for horizontal tables (resp. vertical tables), it will assign a type to columns (resp. rows).

4.2.2. Candidate Generation

The annotation is usually selected within a candidate list that depends on the table elements and the corresponding task. For example, the first step of the CEA task is to generate a list of candidate entities while the first step of the CTA task is to generate a list of candidate types. Thus, STI systems automatically generate, manually add, or filter this candidate list in advance (e.g., an annotation system based on supervised learning should be trained with pre-set labels, those labels being the candidates of the annotation process).

String similarity-based lookup is a standard method to generate candidate entity sets in CEA and row-to-instance tasks. Specifically, the syntactic similarity is calculated between cell mentions and entity labels to select the most relevant entities as candidates. The choice of the matching algorithm depends on the application scenarios. For example, [51] uses Levenshtein-based distances whereas some public knowledge bases expose a public index allowing users to extract and generate candidates, e.g., DBpedia spotlight [63].

In other tasks like CTA (resp. CPA), candidates are the set of types (resp. predicates) available in the target KG. The generation of type candidates is sometimes equivalent to the retrieval of the CEA candidates' types [51, 65]. However, in some learning-based methods [79, 116], type candidates are manually selected for training.

Some mentions can correspond to a large number of entities, e.g. several thousands of entities whose labels contain the mention "France" can be retrieved from Wikidata. The number of candidates may affect the efficiency of the annotation system as calculating and evaluating all possible candidates may require a large amount of time. Hence, some STI systems prune the number of candidates [70, 98]. This filtering process can be applied to reduce the size of the final candidate set by tuning an acceptable entity-mention similarity threshold. However, one should be aware that an inadequate threshold risks to filter the correct entity out. Another way is to sort the importance of candidates by studying the attributes of the target ontology. For example, [25, 65] leverage the BM25 [86] weights from an Elasticsearch index, that can help to select the candidates with the highest usage rate.

4.2.3. Table Elements Processing

Processing different table elements is the core of an STI system. Each table element follows a particular rule according to the table type. Given an input table that is relational and horizontal, each row describes an entity and its attributes, while column elements share the same entity type. Given \mathcal{T}_{ij} , the target cell to be annotated from a horizontal table \mathcal{T} , we identify six table elements that can be leveraged to produce annotations.

 \mathcal{T}_{ij} indicates the target cell itself. Some studies use the string similarity as one of the components of candidate confidence [51, 65]. For example, in Figure 9, for searching CEA candidates for the cell "Suisse" in the first row, we should consider the entities containing the mention "Suisse" in their label. The correct annotation is Q165141 which has the English label "Swiss national football team" and the French label "équipe de Suisse de football".

 \mathcal{T}_{i*} indicates the row context of the target. Some studies leverage the matching degree between the attribute values of the target entity and the information provided in the table for the CEA task. Most of the annotation tasks consider a single-cell subject relational table scenario. Based on this assumption, it is reasonable to make this comparison for calculating the confidence of the candidate. For example, in Figure 9, knowing the date ("11 juin") and city ("Bâle") can help to annotate the third row with the right football game (Q12012827) based on its neighbouring nodes in the KG.

 \mathcal{T}_{*j} indicates the column context of the target. For a horizontal relational table, the information carried by cells in the same column of the table is somewhat similar (e.g., cells referring to the same concept or the same unit). For example, in Figure 9, the cells in the second column are cities. Having this information could help to choose the right candidate between the city of "Bâle" (Q78) and the family name "Bâle" (Q107983752).

 \mathcal{T}_{**} indicates intra-columns relationships from the target table. Intra-column relationships provide a global representation of a table. For example, in Figure 7(c), knowing the relation between the column "Album" and the column "Artist/s" could help to filter out people who did not publish any music album.

 \mathcal{T}_{0*} indicates the header of the target table. The table header often directly explains the contents of the column. Making full use of the information from the table header can help to find the column type or properties more efficiently. For example, in Figure 8(a), the headers "Year", "Country", and "Town" directly denote the concepts of the columns.

 \mathcal{T}_{out} indicates contextual elements encompassing the table. High-quality metadata can help the interpretation of the table. For example, the table's title can potentially determine the domain of the table content, or a hyperlink in a table cell can reveal the identity of the entity. In addition, the text surrounding the table is usually correlated with the content of the table. Leveraging this correlation is used by some STI approaches [34]. Another way to enrich the information used for the disambiguation of the target table relies on inter-table relationships [104].

From the elements mentioned above, the initial annotation step is based on row interpretation, column interpretation, entire table interpretation, or metadata interpretation. In row interpretation, each row in the table describes an entity's attributes. Column interpretation uses the entities in the same column of the relational table with high mutual similarity. This feature can help to constrain the range of candidates for a table cell. Entire table interpretation considers all cells from the table as the context for the disambiguation and is usually performed using deep learning models such as [34, 93].

4.2.4. Iterative Disambiguation

When an STI system jointly undertakes multiple tasks among the five tasks defined in Section 4.1, one task can provide additional useful information for solving the other tasks. For example, when knowing the type annotation of a column (CTA), candidate entities for its inner cells (CEA) that do not belong to the CTA type are less likely to be correct candidates [25]. We call this process iterative disambiguation. This iterative technique is frequently used in heuristic-based approaches (Section 5.1.2) in which a pipeline including specific ordered tasks and being executed once or in a loop is explicitly defined in two ways: (1) predefinition of a pipeline, e.g., [25] performing sequentially CEA, CPA, CTA, CEA disambiguation with CPA or (2) repeat a set of tasks multiple times and stop when it converges to a stable result [84, 121].

5. Semantic Table Interpretation Approaches

In this section, we review the notable approaches from the literature. Among the five tasks introduced in the previous section, the literature mostly focuses on CTA, CEA, and CPA. We propose to classify STI systems according to three representative paradigms of their intrinsic methodology: (1) heuristic methods (Section 5.1) which mostly rely on heuristic techniques such as entity matching, TF-IDF, majority voting, or simple probabilistic frameworks to predict a target; (2) feature-engineering based methods (Section 5.2) which require a feature engineering process to extract statistical and lexical features from the table that are then used to train Machine Learning models; (3) Deep Learning based methods (Section 5.3) that leverage a large number of tables and neural networks to learn deep and contextualised representations of elements of the table, requiring little feature engineering. More details on this classification are shown in Table 1 including representative algorithms, target tasks, table elements used, reference KG, and year of publication. Finally, we discuss these methods from different perspectives: the pros/cons of heuristic-based methods versus Machine Learning based methods, the importance of table elements and the KG structure for improving the accuracy, and the trade-off between efficiency and accuracy (Section 5.4).

5.1. Heuristic Approaches

The heuristic class gathers diverse approaches which are often considered as baseline STI approaches. The core of each system is algorithmically straightforward and does not require much effort in feature engineering or learning. Indeed, the STI tasks are carried out using heuristic techniques such as string similarity measures [65, 76, 105], majority voting [123], TF-IDF [76, 98] or probabilistic frameworks [70]. The context of the table, including the header, the title, and the neighbouring cells [51, 105] is also taken into account but not thoroughly. We further identify two subclasses of heuristic approaches: lookup based approaches (Section 5.1.1) and iterative approaches (Section 5.1.2).

5.1.1. Lookup Based Approaches

Approaches from this paradigm work with an initial candidate entity set determined by a lookup service. After generating candidates through lookup, these methods score the candidates using different metrics on table elements (e.g., cells, type of columns, etc.).

Venetis et al. [102] introduce a model for extracting the column type and the relationships between the key column and other columns. To increase knowledge coverage and avoid issues related to KG incompleteness, the authors present an isA database to carry out the CTA and a relation database to carry out the CPA. The isA database is built by using concept extraction techniques on 100 million English documents that contain the pattern $C[such \ as|including]e[and|, |.]$. They generate the relation database with the help of the TextRunner open extraction system [9]. The authors demonstrated that a hybrid model leveraging a Bayes rule and majority voting has the best performance. The Bayes rule measures the global relevance between cell values and column type labels in tables. The authors conclude that using a target knowledge base (YAGO) leads to higher precision. However, leveraging the isA database can significantly improve the coverage and allow to obtain more meaningful labels for complex or non-explicit table cells.

Wang et al. [105] focus on the table headers to better identify the concept associated with a given column. The approach uses a header detection module that leverages Probase [110] querying and rule-based filtering. In the absence of headers, a custom concept is employed with Probase queries by measuring the type occurrence of the column cells. The table interpretation is executed by studying the cells-header compatibility and entity-values compatibility. Through experiments on a search engine, the approach demonstrates that headers can help understand the columns of a table.

Deng et al. [33] focus on the production of top-k candidates for CTA. The authors first build a Directed Acyclic Graph mapping the entity labels from YAGO and Freebase within a type hierarchy tree. Then, they leverage a distributed system to make the process scalable and efficient without losing precision and accuracy for annotations. Specifically, a two-stage MapReduce system is built. (1) Multiple signatures for each cell mention and entity label are generated to support the cell-label fuzzy matching. For example, "Shar" is one of the signatures for the cell

Table 1: STI systems are classified into three families. We only consider the annotation tasks declared by the authors and when they have related evaluations. "R2I" indicates the task "Row-to-instance"; "TA" indicates the task "Topic annotation"; " \mathcal{T}_{i*} " indicates that, when labeling a cell, information from the same row is used; " \mathcal{T}_{i*} " indicates that the approach leverages information from the target column (CEA, CTA) or columns (CPA); " \mathcal{T}_{0*} " means that the approach has a special treatment on the headers of the table; " \mathcal{T}_{**} " indicates that the approach considers information from all tables elements, including inter-columns influence and training the model with the whole table; " \mathcal{T}_{out} " indicates that the approach not only uses the target table itself for annotation but also considers metadata, including other tables associated with the target table and the text near the original target table.

Approches			Annotation Tasks						Tab	le Ele	ments		V.O	D : 0	Published	
Cl	ass	Algorithm	CEA	CTA	CPA	R2I	TA	T_{i*}	\mathcal{T}_{*i}	T_{0*}	\mathcal{T}_{**}	Tout	KG	Data Source	Year	
		Venetis et al. [102]							V				Custom	Custom Web Tables	2011	
		Wang et al. [105]		V					V				Probase	Custom Wikipedia Tables	2012	
		Deng et al. [33]							V				FreeBase, YAGO	Custom Wikipedia Tables	2013	
		Sekhavat et al. [88]							V				DBpedia	Custom Web Tables	2014	
	Lookup based	TabEL [14]											YAGO	Limaye	2015	
		ADOG [76]											DBpedia	SemTab 2019	2019	
		Tabularisi [98]			\checkmark								DBpedia	T2D, VizNet	2019	
		C^{2} [56]		\checkmark					\checkmark	\checkmark	\checkmark		DBpedia, Wikidata	Limaye, ISWC2017, SemTab 2019, T2D, Semantification, Custom Data	2020	
Hanniatia		Magic [91]			\checkmark								DBpedia, Wikidata	SemTab 2021	2021	
Heuristic		Alobaid et al. [5]		V						V			DBpedia	SemTab 2021, T2D	2022	
		Zwicklbauer et al. [123]		v									DBpedia	Custom Wikipedia Tables	2013	
		T2K [84]			\checkmark		\checkmark						DBpedia	T2D	2015	
		TableMiner+ [121]											Freebase	Limaye, IMDB, MusicBrainz	2017	
		LOD4ALL [65]						V	V				DBpedia	SemTab 2019	2019	
	Thomations	CSV2KG [92]	\checkmark		\checkmark								DBpedia	SemTab 2019	2019	
	Iterative	MTab [70, 73, 74]											DBpedia, Wikidata	SemTab 2019-2021	2019	
		LinkingPark [25]	V	V				V	V				Wikidata	SemTab 2020	2019	
		DAGOBAH SL [49, 50, 51]	V	V	\checkmark			V	V				DBpedia,Wikidata	SemTab 2019-2022	2019	
		MantisTable [27, 26]						V					DBpedia, Wikidata	SemTab 2019-2021	2019	
		JenTab [1, 2, 3]	V	V	V			V	V				DBpedia, Wikidata	SemTab 2020-2021	2020	
		Limaye et al. [60]											YAGO	Limaye	2010	
E.		Mulwad et al. [67, 66]	V	V	V			V	V	V			Wikitology	Limaye	2010	
rea	ture	SemanticTyper [79]											DBpedia	Museum	2015	
engin	eering	DEL [64]											DBradia	City, Museum,	2010	
Da	sea	DSL [64]		\checkmark					\vee				DBpedia	Weather, Custom Soccer	2016	
		Neumaier et al. [69]											DBpedia	Government Data Portal	2016	
		NUMER [55]		v					V				DBpedia	NumDB	2018	
		Vasilis et al. [37]							V				Wikidata	Limaye, T2D, Wikipedia	2017	
	KG	Biswas et al. [15]					\checkmark						DBpedia	Custom Wikipedia inforbox	2018	
	modelling	DAGOBAH Embeddings [23]											DBpedia, Wikidata	SemTab 2019	2019	
	Ŭ	Radar Station [61]	V						V				Wikidata	Limaye, T2Dv2, SemTab 2020	2022	
		Sherlock [48]		\checkmark					V				DBpedia	T2D, VizNet	2019	
Deep		Sato [116]		v					V				DBpedia	VizNet	2019	
learning		ColNet [24]		v					V				DBpedia	Limaye, T2Dv2	2019	
based		Guo et al. [42]		V					V				DBpedia	T2Dv2	2020	
	TT. 1.1.	Zhang et al. [119]							V				DBpedia	T2Dv2	2020	
1	Table	TURL [34]	Ŵ	\checkmark	v			İŻ	, V				DBpedia	WikiGS, WikiTable, T2D	2020	
	modelling	TCN [104]		V	v			ĪV	V		v	√	-	Custom Web Tables, WikiTable [34]	2021	
		DUDUO [93]		v	v			1 V	v		v		-	WikiTable, VizNet	2021	
		Singh et al. [90]		· ·	V			1	V				DBpedia	T2Dv2	2021	
		Zhou et al [122]		\checkmark					V				DBpedia	Custom Wikipedia Tables	2021	

mention "Shark Night 3D". All candidate types according to the Directed Acyclic Graph are then aggregated with subordinative entities of each column. (2) The occurrence of each type is counted in order to select the top candidate type. To accelerate the computation, candidate types are aggregated into disjoint groups.

Sekhavat et al. [88] leverage NELL [95], a Web text corpus, and natural language patterns (PATTY [68]) to extract relations (CPA). The target KG is YAGO, in which only 23 relations are considered as possible annotation for a pair of columns. These are also relations extracted by PATTY from Wikipedia pages. Each relation r is represented by a set of textual patterns $p_1, ..., p_k$ provided by PATTY. Note that a pattern p_i can be associated with more than one relation. Semantic mentions in a pair of columns are linked to KG entities (YAGO) with exact matching. Two columns are connected via a relation rif each pair of entities (e_1, e_2) in the same row belongs to this relation. To determine whether (e_1, r, e_2) is a valid triple, textual contexts related to both e_1 , e_2 are extracted from NELL and are mapped to a list of patterns $p_1, ..., p_k$ in PATTY. The problem then comes down to computing the posterior probability of r given evidences $p_1, ..., p_k$: $Pr(r|p_1, ..., p_k)$. A Bayesian framework is used to compute this posterior.

TabEL [14] aims to provide an extensible framework. After a pre-processing, the approach first generates candidates for each cell using YAGO and ranks them according to their string similarity with the cell and their popularity. Every candidate participates in the calculation of the annotation. In the joint inference module, an undirected probabilistic graphical model is extracted to capture entity-context (elements from the same table) cooccurrence. These co-occurrence factors are updated with the connectivity between candidates, resulting in the final CEA annotations.

ADOG [76] considers scores combined with string similarities, frequencies of properties, and the normalized Elasticsearch score for each match from DBpedia for the CEA task. The system weights these scores with the IDF score of types. To be able to compute the Levenshtein distance and TF-IDF, ADOG uses ArangoDB² to load DBpedia

²https://www.arangodb.com/

and index its components. The frequency of classes and properties is then used to obtain the CTA and CPA results.

Tabularisi [98] adapts TF-IDF statistics to rank the CEA candidates in a given column. A candidate entity is represented by a binary feature vector in which each feature is an indicator (1 if present, else 0) of a property used to describe the entity (e.g., instanceOf). Different features have different expressiveness. They are thus weighted using TF-IDF. Specifically, the "Term Frequency" of a feature is the number of cells whose first candidate entity has that feature, and the "Document Frequency" is the total number of occurrences of that feature in all candidate entities in all cells. The score of a candidate entity is a weighted combination of its TF-IDF score, Levenshtein similarity and word similarity. The weights are either set or learned using a two-layer neural network. The CTA is performed by a top-down brute-force search in the KG class hierarchy tree. The system also sends the most frequent relation among columns with SPARQL queries of cell mention pairs.

 C^2 [56] aims to handle the CTA task with the help of nine datasets. C^2 classifies columns into string entity columns, number columns, and mix-type columns (e.g., emails, dates). The system relies on decision trees, which leverage the pattern of the cell mentions or numerical interval for splitting the branches, to annotate numeric columns and mix-type columns. For entity columns, voting from the concepts extracted from DBpedia or Wikidata for each entity cell is used. C^2 also considers the type co-occurrence within a table to adjust the annotation.

Magic [91] adopts the approach of generating comparison matrices (called INK embeddings) to speed up the computational efficiency. INK embeddings are representations of the attributes and values of an entity or the table context of a cell mention. The complete comparison matrix is generated by fusing multiple candidates. The system outputs CEA annotations by measuring the compatibility between the INK embeddings of the KG and the table. The INK embeddings of entities from the same column are collected to carry out CPA and CTA. For the annotation, the system focuses on the key column: they do the lookup (via public endpoints) for each cell in the key column, then use its neighbourhood to find the candidates for the neighbouring cells in the same row (they do not perform the lookup on the whole table due to limitations of public API usage). Misspellings might however be a challenge for Magic and it cannot detect synonyms of attributes. However, INK embeddings improve computational efficiency and provide a way to implement column wised similarity.

Alobaid et al. [5] handles CTA with a strong focus on the trade-off between type coverage and type specificity after generating type candidates via entities-wised queries on cell mentions. The type coverage is built on a weighted type hierarchy index inside a column, and the type specificity is associated with a distance to the root. The authors test different balance settings between these two factors using the T2D and SemTab datasets.

5.1.2. Iterative Approaches

Iterative approaches are usually built on top of a lookup system, with an additional multi-task disambiguation step for re-ranking candidate entities. The iterative disambiguation techniques, as described in Section 4.2.4, play a significant role in the improvement of the model performance.

Zwicklbauer et al. [123] pioneered iterative majority vote strategies. The idea is based on the majority voting of annotation candidates of the cells for the CTA task. The system generates these cell candidates using a search-based disambiguation method [124]. Since majority voting plays an essential role in the system, this approach is sensitive to the number of table rows. The extreme case is that the annotation precision can be less than 0.1 with single row tables.

T2K [84] annotates Web tables by mapping their columns to DBpedia properties, and their rows to DBpedia entities, associating the whole table with a DBpedia class. A key column's position for each table is firstly detected by a preprocessing step. T2K transforms row-to-instance and table topic tasks into CEA and CTA on the key column. The initial entity mapping is derived from a lexicographical comparison between the labels used in the table and those of the entities described in DBpedia with Jaccard, Levenshtein, and deviation similarities. An aggregation of these similarities is then used to choose the initial CTA annotation. The system adopts an iterative process between a CEA matching module and a CPA matching until the output is stable. The system achieves promising results on CEA and topic annotation.

TableMiner [120] and the following work TableMiner+ [121], first use a lexical expression to extract a primitive type for each column. Similarly to T2K, TableMiner sees row-to-instance and table topic as CEA and CTA on the subject column. The system identifies a subject column for each table considering evidence collected from the original Web pages. During the annotation phase, the authors iteratively label records corresponding to the subject column and their attribute columns using information from the HTML context of the tables to create a richer representation of cells and columns. TableMiner+ uses partial matching during the CTA annotation. The authors claim that partial matching is efficient and that eight rows are enough for supporting the annotation. A loop is used between CEA annotation and CTA annotation until the results remain stable. This system leverages the Freebase KG and was evaluated on the IMDB and MusicBrainz specific datasets, as well as on Limaye.

LOD4ALL [65] is initialized by building an RDF store database and a score DB database containing candidate types with scores reflecting the level of specificity generated by Okapi BM25 [86]. The system uses a similar approach as [123] for the candidate generation module, which uses a combination of Elasticsearchs score and the SimStrings score to select the top 100 candidates. The CTA leverages Okapi BM25 type scores and the type coverage on the table. The CEA and CPA are calculated after CTA type filtering. The system is targeting DBpedia as ontology and has participated in the SemTab 2019 challenge. A similar pipeline is used by **CSV2KG** [92]. However, CSV2KG considers applying a threshold on the normalized entropy of the two highest type counts from the annotated candidates parent types list to decide the level of granularity of the CTA.

MTab [70] employs four different lookup services. This approach analyses signals from server's lookup ranking, header context analyses, SpaCy type prediction³, Duckling type prediction⁴ and value similarities. Each signal is transformed into a normalized probability score. The system aggregates selected probabilities with learnable weights according to the associated task. In addition, the authors also used EmbNum [72] to help to produce annotations on numerical columns. Column types and column pair relations are computed based on entity scores. Entities, types and relations are iteratively calculated two times to disambiguate CEA, CTA, and CPA annotations with inter-tasks relatedness. In a more recent work, MTab4Wikidata [73, 74] adapts fuzzy matching and "two cells search" to enhance the support of misspelling and ambiguities in table content. The system won the first prize in both SemTab 2019 and SemTab 2020 challenges.

LinkingPark [25] leverages the Wikidata MediaWiki API to generate cell candidate entity lists. It also adapts a fine-grained Elasticsearch index to rank those candidates. This system firstly adopts a cascaded pipeline to generate candidate entities. Then the system disambiguates each cell through an iterative coarse-to-fine algorithm by considering the CPA annotation results. The system finally generates the CTA annotations from the disambiguated cell annotations. The authors also claim that Wikidata's type ontology is noisy which makes it difficult to assign types during CTA annotation.

DAGOBAH SL [49, 50, 51] calculates a score for each cell entity. This score combines the Levenstein similarity and the context similarity between the entity's neighbouring nodes and the row context of the target cell. The authors collect triples containing cell candidates between the two columns and calculate the sum of the weights from the corresponding cell candidates for each relation. The output of CPA is the relationship with the highest sum. The system then leverages the CPA results to perform CEA disambiguation. As for CTA, in addition to the majority vote based on CEA, DAGOBAH SL also leverages the distance to the root concept and the ranking of each Wikidata entity's class to select the most accurate type. In more recent works, DAGOBAH SL [50] enhances the system with CTA disambiguation. The entity context

made with multiple-hop neighbouring entities is also taken into account in calculating the scores. The system won the first prize in accuracy in the SemTab 2021 challenge. DAGOBAH SL [49] also uses language models to extract more meaning from the table headers and largely expands its aliases table using external sources for improving the lookup coverage. The system won again the first prize in accuracy in SemTab 2022 challenge.

MantisTable [27, 28] pipeline starts with classifying each column into three types: Named Entity column, Literal column, and Subject column. The candidate generation of this approach is based on SPARQL queries which extract all candidates containing the cell mentions. Then, the system handles the CEA using row-wise compatibility analysis and CPA using majority voting. For the CTA task, the authors list all candidate types in addition to their number of occurrences in the table (row coverage). After filtering with a threshold, the rest of the type candidates are transformed into a graph according to the ontology hierarchy. Type scores are then updated with the distance to the root. In the end, the highest score represents the most accurate and specific annotation. The recent version of MantisTable SE [26] optimizes the system by updating the scoring function, accessing the LamAPI⁵ API (instead of using a SPARQL endpoint) and adding a final disambiguation step. MantisTable also supports row-toinstance by applying CEA on a subject column detected in preprocessing.

JenTab [1, 2, 3] starts from analysing the row context and column context of a table. This system wraps each computational unit into independent modules so that they can be recalled easily and repeatedly. Those modules include row information processing, disambiguation using CTA output, etc. The system leverages different module combinations for supporting CEA, CTA and CPA tasks. The evaluation shows that JenTab has an excellent performance on synthetic datasets. The authors also investigate the implications of considering multi-hop links in type hierarchy relations. The result shows that considering two hops has a small probability of improvement, while multi hops lead to a significant decrease of the accuracy.

5.2. Feature Engineering based Approaches

This family of methods extracts statistical and lexical features (such as distribution of numerical values, occurrence of cell mentions, textual similarity, etc.) from the table rows and columns and uses them with machine learning models. Typical algorithms used for STI include SVM [67], Random Forest [64] and K-Nearest Neighbor [69] for example. A labelled dataset is required for the training. The amount and the quality of training data, and consequently the quality of input features, have a significant impact on the model performance, as discussed in [64]. In addition, we observe that ML methods target the CTA task more

³https://spacy.io/

⁴https://github.com/facebook/duckling

⁵https://bitbucket.org/disco-unimib/lamapi

than other tasks, as columns can provide more statistical features than other annotation targets.

Limaye et al. [60] introduce one of the first works on STI. The approach computes the TF-IDF cosine similarity between a cell mention and an entity label and the compatibility between the cell type and the column type to execute the CEA task. CTA task depends on TF-IDF cosine similarity between column header and each entity's type label. The CPA annotation depends on the compatibility between the relation and column pairs. All these features are weighted through a machine learning framework.

Mulwad et al. [67] leverage majority voting on the type of cells candidates for the CTA annotation. The PageRank algorithm weights each cell's type from the ontology during the CTA. For CEA annotation, the system collects entity features from entity PageRank, string similarities, entity index score and entity page length to generate a vector representation of each entity. An SVM classifier is then built upon these features to make predictions. This system supports both DBpedia and Wikitology.

SemanticTyper [79] processes each column of the table independently for CTA. The system first distinguishes between columns of numbers and columns of strings by voting on the types of cells in each column given a predefined threshold. Strings columns are trained upon cosine similarities on TF-IDF, considering a column as a document. To annotate numerical columns, the authors use a variety of distribution representation methods. In both cases, they adapt the training data that consist of a set of semantic labels associated with samples of data values. The prediction aims to find the most similar candidate by comparing the distance between the query column and each sample set in the training data corresponding to a distinct semantic label. The chosen distance metric depends on the column type. The training dataset was extracted from vertical domain datasets such as museums and cities, and the authors associate columns from these datasets with DBpedia classes. The main limitation of this work resides in not taking into account relationships between columns.

DSL [64] build their approach for CTA on datasets from four different domains: city, weather, museum, and soccer. Labels are partially manually added to these datasets DSL leverages features including string similarity and number distribution between chosen labelled datasets and the rest of the data during the prediction. The difference with SemanticTyper is that the distribution is also available for string columns in this approach. The system learns the weights between these features through two supervised learning algorithms: logistic regression and random forest. The evaluation shows that logistic regression achieves better results. The authors claim that the incorrect predictions come from the top decomposition point of the decision tree.

Neumaier et al. [69] focus on CTA labelling for numerical columns. Their work is not limited to predicting a

unique label but rather expands the scope of labelling to its surrounding information. For example, instead of labelling "height", this system will label it as "the height of an athlete playing basketball in the NBA". To do so, the authors constructed a background KG based on DBpedia. This background knowledge base is extracted as a hierarchical structure divided into multiple multi-level groups to provide context. Each node in the hierarchy represents a type or a predicate and provides statistical information (maximum, minimum, or distribution) of the relative number set as features. The authors use these features and KNN to make predictions. The authors also explore the system's performance at different hierarchy levels in the background KG built on DBpedia and Open Data. They pointed out that DBpedia has still limitations in terms of coverage and freshness compared with other open datasets. For example, the Austrian Open Data Portal has tables generated by weather stations every 15 minutes. However, DBpedia typically has numeric values only for "current" or "latest". Another limitation of this work is that the size of numerical columns and the popularity of numerical KG properties may influence the accuracy. Hence, NU-**MER** [55] proposes to link subject cells with KG entities first, and then only extract the linked properties for enabling the column-wise number distribution study.

5.3. Deep Learning based Approaches

Deep Learning has achieved many successes in various domains thanks to the availability of huge amounts of data and powerful computing resources. It has attracted more and more attention from the STI community over the past few years. We identify two main directions for the application of deep learning in the STI domain: KG modelling (Section 5.3.1) and table modelling (Section 5.3.2).

5.3.1. KG Modelling

This direction focuses on the entity level in which models learn embedding representations for entities of a table cell instead of the cell itself. Specifically, KG embedding techniques (e.g., TransE [20], TransH [107]) are used to encode the entities and their relationships into a vector space. STI models rely on the intuition that the entities in the same column should exhibit semantic similarities. Hence, they should be close to each other in the embedding space w.r.t. cosine similarity distance [37] or Euclidean distance [23].

Vasilis et al. [37] provide different methods. One of the proposed systems assumes that the correct CEA candidates in a column should be semantically close. From this assumption, a weighted correlation subgraph in which each node represents a CEA candidate is built. The edges are weighted by the cosine similarity between two related nodes. The best candidates are the ones whose accumulated weights over all incoming and outcoming edges are the highest. In addition, a hybrid system that combines the correlation subgraph method and an ontology matching system is also introduced and achieves a significant improvement in the end.

Biswas et al. [15] focus on topic annotation for Wikipedia infoboxes by leveraging metadata from the Wikipedia page. The annotation ignores the infobox content since the infobox information is often incomplete, incorrect and missing. Wikipedia page section headers and abstract are extracted as the source of information and are featured by Word2Vec embeddings for each word. Note that named entities present in the abstract are also transformed into RDF2Vec [81] vectors with DBpedia pre-trained embeddings. The global representation vector (called document embeddings) of an infobox is the concatenation of Word2vec and RDF2vec vectors. Two classifiers (Random Forest and CNN) are trained on top of the document embeddings on 150,000 tables with 30 preset types. The evaluation shows that CNN outperforms the Random Forest classifier.

DAGOBAH Embeddings [23] hypothesizes that all entities in the same column of the table should be close to each other in the embedding vector space. Consequently, the correct candidates are assumed to belong to a few clusters. The K-means clustering is performed using TransE's pre-trained embedding to cluster the candidate entities. The good clusters with high coverage are retained by a weighted voting strategy. Both CEA and CTA are selected from the chosen clusters. Experimental results prove that this approach has successfully improved the accuracy of the CTA task. However, the system is also misled by incorrect candidates in the selected clusters during the CEA task.

Radar Station [61] went a step further in proposing a hybrid system that aims to add a semantic disambiguation step after a previously identified CEA. Radar Station takes into account the entire column as context and uses graph embeddings to capture latent relationships between entities to improve their disambiguation. RadarStation has been evaluated on top of different heuristics-based systems (DAGOBAH SL, BBW, MTab) and have consistently demonstrated an accuracy improvement of around 3%. Furthermore, the system shows empirical evidences that among the various graph embeddings families, the ones relying on fine-tuned translation distance have superior performance compared to other models.

5.3.2. Table Modelling

This direction deals directly with the textual content of the table as well as intra-table and inter-table interactions. The contextualized representation of basic elements of the table (i.e. cell, column) is learned by using deep neural networks [24, 48] or language models like BERT [34, 93, 122].

Sherlock [48] learns to perform the CTA task using 1588 features extracted from a single column of a given relational table. The features are divided into four categories: character-wise statistics (e.g., frequency of the

character "c"), column statistics (e.g., mean, std of numerical values), word embedding, and paragraph embedding. Except for column statistics features, other features are compressed into a fixed-size embedding using a subnetwork. A two-fully connected layer network is trained on both the embedding features and column statistics features to predict a column type annotation among 75 types inherited from the T2Dv2 dataset. The evaluation shows good results on various column types, including Dates and Industry. However, it is less sensitive to the purely numerical values or values appearing in multiple classes. Facing the potential missing information in single-column annotation, Sato [116] extends Sherlock by considering the whole table context. The table topic embeddings with LDA features modelled by an additional subnetwork and the columnpairs-wise dependency modelled by a CRF layer are also studied.

ColNet [24] predicts the column type (CTA) using only intra-column contextual information. Specifically, all the cell mentions of a column are split and concatenated into a single word sequence. This word sequence is converted into an embedding vector using word embeddings models like word2vec. This embedding is later input into a CNN model. It is worth noting that ColNet predicts the type of each column independently, and thus ignores the inter-column contextual information. To generate a training dataset, ColNet makes use of a KG to collect candidate classes given the cells of the column. For each candidate class c, N different sets of h entities belonging to c in KG are retrieved and the associated word embeddings are stacked into a matrix \mathcal{M} . This forms N training samples $\{\mathcal{M}, c\}$ for class c. To alleviate the computational complexity, the candidate classes that appear rarely are not modelled and h rows of the column are randomly selected to predict a type annotation. The type prediction of ColNet is then refined by the matching entities of inner cells through majority voting. Similar work includes Guo et al. [42] where the authors also introduce the use of BiGRU-Attention rather than CNN and support multicolumn annotation using linear-CRF (linear-chain conditional random field) on the entire network table and an undirected graph model that directly models the conditional probability.

Zhang et al. [119] address the table-to-KG matching including CEA and CPA tasks. Additionally, in the pipeline, with the help of a table-to-KG matching step, their tool performs a novel entity discovery task. A classification model is based on syntactic similarity (e.g., edit distance, Jaccard distance) and semantic similarity (deep semantic matching method DRMM [40]) between a table mention and corresponding candidate entities to determine whether the mention is linkable. If yes, at most one entity is predicted for this mention. Another classifier is then built on the cell annotations (CEA), exploiting columnwised features such as naive features (e.g., length of the header), label similarity between header and properties of cell entities, value similarity between literal values (e.g., numerical, time) in the table and literal values of cell entities for the column property matching (CPA).

TURL [34] pioneers the application of pre-trained language models such as BERT in the STI domain. It provides a universal contextualized representation for each table element (i.e. caption, header, content cells) which can be fine-tuned and applied in various downstream tasks such as CEA, CTA, CPA, or table augmentation. The table augmentation task mainly involves enriching the semantics of the table by extending it with new columns (attributes). The model employs a Transformer-based encoder [101] to capture the information from table elements. To this goal, the input table is first serialized into a sequence of caption tokens, title tokens, header tokens, and row-by-row cells. A cell consists of its content (mention) and a candidate entity representing it in a KG. The sequence of tokens is then converted into embeddings using word embeddings for textual tokens and KG embeddings for entity tokens. To reduce the redundancy in the fully-connected attention learning and better draw the inter-column and intra-column, inter-row and intra-row, column-row interaction, the conventional attention layer is masked by a so-called visibility matrix which allows only a portion of table elements to participate in the modelling of a specific element. For example, cells in the same row or the same column can interact with each other. Apart from the BERT's Masked Language Model objective, TURL introduces an additional Masked Entity Recovery objective to reinforce the learning of factual knowledge embedded in the table and represented by KG entities. The model is trained on 570K relational Wikipedia tables.

Singh et al. [90] introduces a method based on BERT for relation extraction from tables (CPA). Only two-column tables are studied in this work. The table is tokenized and transformed into linearized rows and linearized column headers that are passed through a pre-trained BERT encoder to obtain two vector representations for table content and table header. A fully-connected layer takes as input these two vector representations and predicts scores over all candidate relations in the two-column table. The model is trained on synthetic tables generated automatically from a KG. However, synthetic tables lack metadata such as column headers and captions. According to the authors, such metadata is important for the relation extraction task. Hence, they propose a novel method which generates synthetic tables associated with metadata (context of table contents, meaningful column headers).

TCN [104] not only exploits intra-table contextual information but also inter-table contextual information for the two tasks CTA and CPA. According to the authors, the global context of a table can be complemented by discovering its implicit connections with other semantically related tables. Such inter-table connections can come from overlapping cell contents, consistent schemas or similar table topics between two tables. The embedding representations of a specific table cell and the table topic are jointly learned leveraging intra-table contexts (i.e. other cells in the same

column or the same row, the table topic) as well as intertable contexts (i.e. the cells sharing the same value, the columns with a similar header and the topic from other related tables). All of these contexts are fused into the embedding through the attention mechanism. As in DO-DUO model [93], the CTA and CPA tasks are trained in a supervised manner on two specific objective functions dedicated to column type prediction and column pair relation prediction, respectively. In addition, in the case where a large annotated training dataset is not available, TCN switches to transfer learning in which the ultimate cell embedding is fined-tuned with a BERT-like unsupervised pre-training. The evaluation shows that the intertable contextual information contributes positively to the model's performance. However, the utilisation of intertable context remains challenging since it requires prior knowledge about tables' schemas which are generally diverse.

DODUO [93] learns to jointly annotate the column type and column pair relation through multi-task learning. Similarly to other Transformer-based models, the key idea of DODUO is to incorporate table contexts (intra-column and inter-column contexts) into the prediction of a single column type or a single column pair relation using a tablewise attention mechanism. The model serializes the input table column by column into a sequence of tokens in which each token represents either a column header or a column cell. A special [CLS] token is appended at the beginning of each column to distinguish two different columns. This token is also considered as the embedding representation of the column itself. During the inference phase, two different output layers perform two different tasks: one takes a single column representation (i.e. the hidden vector of the [CLS] token) as input and produces a semantic type for this column accordingly, one takes a pair of column representations and predicts a relation between them.

Zhou et al. [122] focus on the column type detection (CTA) task. This work leverages the Star-Transformer model [41] to learn a vector representation for each column taking the inter-column interactions into account. The input embedding of a column is initialized as a concatenation of semantic features which are word embeddings averaged over cell contents and statistical features which are adopted from Sato's [116] Sherlock model. In the context of limited training tabular data and weak order-dependence of table columns, the Star-Transformer is preferable to the Transformer as it reduces the computational complexity by replacing the fully-connected attention with a sparser one.

5.4. Discussion

We have grouped the many STI approaches proposed so far into three families or paradigms. In this section, we further analyse these methods alongside different dimensions: the pros and cons of matching methods versus learning methods, the trend towards deep learning methods, the importance of table elements, the trade-off between efficiency and accuracy, and the influence of the target KG structure for improving the accuracy (but at the risk of adding noise).

5.4.1. Matching vs Learning

We observe that STI approaches rely on matching (a KG entity with a cell mention) and learning. Matching is key in heuristic-based approaches while feature engineering and deep learning based methods rely on representation learning of the input table. These can also be combined: the matching strategy is employed by learning-based models as a post-processing step where the annotations learned from neural networks are refined using mention-matching entities (DAGOBAH Embeddings [23] and Col-Net [24] are two examples).

From our observation, matching highly relies on the compatibility between the target table and the target KG. Consequently, it may be challenged by incompleteness in the table, knowledge shifting of KG and incompatibility between the table and the KG. Matching methods are less robust to noise than the learning methods. On the other hand, learning methods need large training datasets which are not always easy to collect or generate. Some learning approaches limit the number of target candidates to alleviate the lack of training data. [34, 48, 64, 79, 93, 104] predict the CTA within a predefined set of around one hundred types. ColNet [24] deals with the data shortfall using data augmentation. The model is trained on data generated automatically from a KG. However, it takes hours to do a single annotation. Another challenge for learning approaches is the size of target tables. To perform the CTA, [48, 64, 79] rely on the statistical features computed from the table (e.g., distribution of number, length of string for each cell, etc.). These features are not statistically stable if the number of samples (i.e. table cells) is low. Finally, we also discovered that learning approaches do not consider thoroughly the hierarchy of types possibly used in a KG impacting the type specificity returned by the CTA task [64, 79, 93].

5.4.2. Rise of Deep Learning

After 2017, deep learning techniques emerged in the STI field and have attracted research focus. Compared to feature engineering approaches, complex neuronal networks allow the system to process tabular features more efficiently as the feature engineering step is sometimes difficult and time-consuming to maintain. For example, Sherlock [48] is based on 1588 column-wised features. To mitigate this issue, an end-to-end learning framework is preferable and is more and more employed, for example, KG modelling with KG embedding techniques (e.g., Vasilis et al. [37]) and table modelling with BERT-like models (e.g., TCN [104]). However, we observe that table modelling approaches using language models always target class annotation (CTA) or relation annotation (CPA) tasks. The entity annotation task (CEA) still lacks dedicated work

and has a lot of room for improvement. At the time of conducting this survey, TURL may be the only one that handles the CEA task. Moreover, many systems [93, 104, 122] try to simplify the table representation to a collection of unordered lists for columns or rows, ignoring their index and other structural information. TURL [34] proposes a visibility matrix, as an attention matrix, to describe the connections between table elements (e.g., cells in the same columns, cells in the same rows, etc.). We argue that this design is only applicable to relational tables whereas the model is trained on a dataset containing all table types. This limitation still requires more effort to cover more complex scenarios.

5.4.3. Coverage of Table Elements

Annotation models use different table elements that are analysed in depth in Section 4.2.3. We observe that more and more table elements are considered in recent approaches. This phenomenon is characteristic of learningbased models. [48, 64, 69, 79] primarily concentrate on extracting features from a single column. With the rise of the deep learning and attention mechanism, [34, 93, 104, 116] start to pay more attention to other complex table elements. A typical example is TURL [34], in which the authors consider all of the six table elements discussed in Section 4.2.3. However, we can not compare the superiority of an approach only by the number of elements it uses. The search for the right number of elements taken into account to increase the accuracy without being subject to noise remains an open challenge.

5.4.4. Effectiveness vs Efficiency

Annotation systems usually deal with a trade-off between effectiveness and efficiency. TableMiner+ [121] introduces partial matching in which the CTA calculation relies on only eight table rows in order to improve the performance. This strategy indeed makes the systems faster but degrades the accuracy. For example, considering the annotation of a column containing ["Joe Biden", "Donald Trump", "Barack Obama", "Abe Shinzo"], applying partial matching on the first three column cells will output "American presidents" as the type of this column, while the correct answer is more likely to be "politicians" since "Abe Shinzo" is not an American president but a Japanese prime minister. Systems whose annotation pipeline includes a candidate generation step will heavily depend on the entity lookup service used. However, public lookup endpoints impose several limitations on their usage and it may take more time to obtain a candidate set with a desirable coverage of the target table. Furthermore, some systems (e.g., [70]) are not suitable for real-time applications due to the heavy computational requirements of their intrinsic algorithm. In addition, in scenarios where there is not enough data for training, learning-based models take advantage of transfer learning. While it helps to save time and resources, the system accuracy may be degraded if the fine-tuning is not carefully performed.

5.4.5. Public KGs vs Custom KGs

Many approaches to annotate tables rely on encyclopedic KGs such as Wikidata and DBpedia. Those KGs provide rich and high-quality information helping the annotation become more effective. However, more information also leads to more ambiguity, and KGs are usually incomplete. Knowledge base shifting is also a challenge for approaches based on public KGs. We observe that some approaches [34, 48, 64, 93, 104, 79] only treat the target KG as a dictionary of concepts but not a knowledge network, which means the relationships and hierarchy of concepts have not been used. The extreme case is that some attention-based models [93] directly abandon KGs and use only concept names. Knowing how to properly inject a KG structure into a statistical model is an open challenge. Some works build their custom KG to increase the coverage. For example, [102] built an isA database using Web documents which contains three times more types than Freebase.

6. Datasets and Benchmarks

Several datasets have been proposed to evaluate STI approaches. Some datasets attempt to establish gold standards in which table components (cells, rows, columns, or cell pairs) are associated with KG components (entity, class, or property), while others collect high-quality tables to support STI training. In this section, we detail the most popular datasets (Table 2) that are used by STI approaches.

Limaye [60] is one of the earliest gold standards used in the community. Limaye aims to annotate Web tables using the YAGO KG. The dataset is divided into four subsets according to the data source, the labelling method, and application scenarios. Three subsets are manually labelled while the fourth one is automatically labelled. The automatically labelled subset contains annotation errors [66] which were corrected by [14]'s work in 2015. Later on, [37] updated the disambiguation links to the DBpedia KG.

T2D [85] is taken from the Web Data Commons project.⁶ DBpedia is used as the target KG and extensive metadata such as the context of the table and whether the table has a header or not is provided. The addition of a small number of non-overlapping tables that do not have a mapping relationship with DBpedia makes this gold standard closer to real-world datasets. A second version of the gold standard adding negative examples has been published and named T2Dv2 [59]. T2D and its supplementary version T2Dv2 [59] have been largely used to evaluate approaches in [24, 48, 121], and together with Limaye, they became the de facto gold standard datasets to use for evaluating STI approaches. However, [39] pointed out that T2D has partial annotation errors and lacks fine-grained annotations since a large number of tables only point to the root

class owl:Thing. [39] has proposed a revised version of the dataset named T2D^{*}.

WDC [58] leverages the DWTC framework [36] to crawl tables from the Web and to distinguish between different types of Web tables according to the inner-relation dimension and the orientation dimension that have been defined in Section 3.1. The crawler has extracted a total of 10.24 billion tables from which 0.9% are relational tables, 1.4% entity tables, and 0.03% are matrix tables, amounting to 233 million tables available in the corpus. The rest are labelled as layout tables or other tables. WDC also provides a subset containing 90 million relational tables and a subset containing 50 million English relational tables.

TabEL [14] (also named WikiTables) has collected 1.6 million Wikipedia tables which contain the class attribute "wikitable" from the November 2019 XML English Wikipedia dump. During the extraction, the system collects the hyperlinks in table cells and metadata of the table, such as the table caption and the page title. Thus, the mappings between YAGO entities and table cells are easily derived. The types of tables in this dataset are unknown. **QuTE** [44] further enhanced the dataset by merging TabEL with 2.6 million tables from the TableL [52] dataset that were extracted from 1.5M Common Crawl Web pages using the DWTC framework [36]. The TableL dataset covers mostly five major topics: finance, environment, health, politics, and sports. These datasets have generally been used to train machine learning based STI approaches [14].

Zhang et al. [121] proposed datasets⁷ for evaluating entity linking in Web tables, as well as table header classification and relation annotation. The datasets contain 16,000+ annotated relational tables that can be used for many studies related to Web tables. In particular, it proposes the **IMDB** (movie) and **Musicbrainz** (music) datasets with cell entities and column headers annotated using Freebase topics.

The **SemTab** competition (Semantic Web Challenge on Tabular Data to Knowledge Graph Matching⁸) colocated with the International Semantic Web Conference in 2019, 2020 and 2021 provides the biggest datasets. This competition has attracted nearly 50 participant teams over the three years. The SemTab 2019 datasets [53] use DBpedia as the target KG, while SemTab 2020 [54] uses Wikidata and SemTab 2021 [32] uses both DBpedia and Wikidata in addition to Schema.org for the last round. The competition consists of four rounds in 2019 and 2020 and three rounds in 2021. The data sources vary depending on the rounds. SemTab 2019 Round 1 is a small number of high-quality tables extracted from T2Dv2. Tables from SemTab 2019 Round 2 are extracted from Wikipedia. Except GitTables [47] and BiodivTab [4], the rounds from SemTab contain a large number of artificially generated

⁶http://www.webdatacommons.org/webtables/

 $^{^{7} \}rm https://github.com/ziqizhang/data/tree/master/webtable\% 20 \rm entity\% 20 \rm linking$

⁸http://www.cs.ox.ac.uk/isg/challenges/sem-tab/

Table 2: Gold standard datasets for evaluating STI approaches. Table type refers to the classification introduced in Section 3.1 where R=relational table, SR=single-cell relational table, V=vertical, H=horizontal.

Gold st	andard	Table Type	#Tables	Avg. #Rows	Avg. #Col	#Entities	#Class	#Relations	Origin	KG	Year
Limovo ^a	Manual	SR-H	437	37	2	10,930	747	90	Web	Wikipedia, YAGO	2010
Limaye	Wiki_link	SR-H	6,085	20	2	131,807	-	-	Web	Wikipedia	2010
$T^{2}D^{b}$	$T2D^{c}$	SR-V	762	157	5	25,119	7983	-	Web	DBpedia	2015
12D	$T2Dv2^d$	SR-V	779	84	5	26,106	755	-	Web	DBpedia	2017
WDC^{e}		All	233,000,000	12	4	-	-	-	Web	-	2015
TabEL ^f		R	1,652,771	11	5	3,000,000	-	-	Wikipedia	YAGO	2015
Qu	ΓE^{g}	R	1,766,721	13	5	-	-	-	Wikipedia	-	2021
	R1	SR-H	64	142	5	8,418	120	116	Web	DBpedia	2019
SomTab 2010h	R2	SR-H	11,924	25	5	463,796	14,780	6,762	Web	DBpedia	2019
Semirab 2019	R3	SR-H	2,161	71	5	406,827	14,780	7,575	Synthetic	DBpedia	2019
	R4	SR-H	713	63	4	107,352	1,732	2,747	Synthetic	DBpedia	2019
	R1	SR-H	34,295	7	5	985,110	34,294	135,774	Synthetic	Wikidata	2020
	R2	SR-H	12,173	7	5	283,447	26,727	43,753	Synthetic	Wikidata	2020
SemTab 2020	R3	SR-H	62,614	6	4	768,325	97,586	166,633	Synthetic	Wikidata	2020
	R4	SR-H	22,207	21	4	1,662,164	32,462	56,476	Synthetic	Wikidata	2020
	ToughTables	SR-H	180	1,080	5	663,830	539	-	Synthetic	Wikidata	2020
	R1-DBP	SR-H	180	1,081	4	663,656	540	-	Synthetic	DBpedia	2021
	R1-WD	SR-H	180	1,081	4	667,244	540	-	Synthetic	Wikidata	2021
	R2-Bio	SR-H	110	2,449	6	1,381,325	657	547	Synthetic	Wikidata	2021
SemTab 2021	R2-Hard	SR-H	1,750	17	3	47,440	2,191	3,836	Synthetic	Wikidata	2021
	R3-Hard	SR-H	7,207	9	2	58,949	7,207	10,695	Synthetic	Wikidata	2021
	R3-Git	SR-H	1,101	59	17	-	123 / 60	-	Github	DBpedia, Schema.org	2021
	R3-BiodivTab	SR-H	50	260	24	31,468	614	-	Open data	Wikidata	2021

^ahttps://zenodo.org/record/3087000#.YbY5Lp7MJPY

 b The # class for the T2D dataset is the sum of the number of "table classes" and "column properties" in the original dataset. c http://webdatacommons.org/webtables/goldstandard.html

 ${}^{d}\rm http://webdatacommons.org/webtables/goldstandardV2.html$

^ehttp://www.webdatacommons.org/webtables/

fhttp://websail-fe.cs.northwestern.edu/TabEL/

 g https://www.mpi-inf.mpg.de/research/quantity-search/quantity-table-extraction

^hThe SemTab series are indexed at https://www.cs.ox.ac.uk/isg/challenges/sem-tab/

tables annotated for a target KG using SPARQL queries with the introduction of some noise such as misspellings.

7. Evaluation

The synthetic tables datasets can be used to test the scalability of a system given their size and the large number of lookup candidates that can be returned. Nevertheless, these datasets are not extremely challenging for the semantic interpretation approaches as they contain synthetically generated noise. This explains the very high accuracy of the top-performing approaches (F1 score up to 0,99 for some tasks [54]). There is also room for improvement in incorporating all real-world challenges in future editions of the challenge, for example, tables with cells containing multiple entities. To increase the difficulty, SemTab 2020 introduces during round 4 the socalled **Tough Tables** dataset [31]. Tough Tables is composed of specially designed tabular data simulating various difficulties: a large number of rows to evaluate the systems performance, non-Web tables and artificially added misspellings and ambiguities.

In SemTab 2021, the organizers bring new challenges with, in particular, the BiodivTab [4] dataset which contains 50 manually annotated tables from real-world biological datasets. BiodivTab also contains artificially generated noises, abbreviations and complex data formats. A subset of the GitTables dataset [47] has also been used in the last round of SemTab 2021. This dataset is a collection of CSV tables from GitHub which are annotated with DBpedia and Schema.org.

Traditionally, STI systems are evaluated using the information retrieval metrics: accuracy, recall, and F1 scores [23, 34, 53, 64, 121]. Among the various tasks, CTA is a special one since a given type and its parents can be all correct when determining the category of an entity or of a group of entities. For example, Paris can be equally typed as a capital or more generally as a city, which are not in conflict with each other. Considering only one of these types, such as capital, as the only correct answer would make it difficult to assess systems that can only predict the parent class. The SemTab 2019 challenge has defined the AH (Equation (1)) and AP (Equation (2)) scores for this purpose, where PerfectA indicates that the predictions made by an STI system exactly match the type declared in the ground truth and OkayA refers to annotations corresponding to one of the parent types. This evaluation method is also adopted by [24]. The SemTab Challenge evaluates systems using the Alcrowd evaluator⁹ and STILTool [29].

$$AH = \frac{\#\text{PerfectA} + 0.5 \times \#\text{OkayA} - \#\text{WrongA}}{\#\text{TargetColumns}}$$
(1)

$$AP = \frac{\#\text{PerfectA}}{\#\text{AnnotatedColumns}} \tag{2}$$

 $^{^{9}} https://github.com/sem-tab-challenge/aicrowd-evaluator$

SemTab 2020 and 2021 further enhance the CTA evaluation metrics by introducing the correctness score *cscore* (Equation (3)) for correctly positioning the annotated type in the hierarchy tree where $d(\alpha)$ indicates the distance between the annotation α and the ground truth.

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

Given the *cscore*, approximated Precision (AP), Recall (AR), and F1-score (AF1) for the CTA evaluation are then defined as follows:

$$AP = \frac{\sum cscore(\alpha)}{\#\text{Annotations}}, \ AR = \frac{\sum cscore(\alpha)}{\#\text{Targets}}$$
(4)

$$AF1 = \frac{2 \times AP \times AR}{AP + AR} \tag{5}$$

Table 3 gives the performance of the top three systems with the highest F1, AP, or AF1 scores for the CEA, CTA and CPA tasks on the datasets commonly used by the community. We observe a considerable diversity in the evaluation process across studies. This diversity is reflected in: i) the original version of the Limaye dataset had some errors that were corrected by the TabEL team. However, these corrections are not used by STI systems; ii) the TabEL and C^2 teams only report on accuracy scores, thus lacking an evaluation of the recall to compare with other systems; iii) ColNet has further enhanced the T2D dataset while other approaches do not consider these enhancements; iv) [56] provides aggregated results, and it is hard to get the performance details per dataset, and v) the performance of a system with the same dataset in different articles can largely vary. For example, the difference in accuracy between ColNet's T2D dataset between [24] and [56] is about 60%, probably because [56] has not used all of ColNet's settings. This diversity makes it difficult to compare the performance of different STI systems and stresses the importance of challenges such as SemTab in the STI community.

We adopted the following strategy to produce Table 3: i) we only consider datasets that have been used for evaluation more than three times. In addition, the datasets used solely for training purposes (e.g., TabEL [14] and Viznet [46]) have not been considered; ii) if the performance of a system on the same dataset has been reported multiple times, we only consider the accuracy from the original paper; iii) we prioritize comparing F1 scores except for CTAs in SemTab 2019, which take into account the AP score. If F1 scores are not provided, we use precision for the comparison; iv) for the SemTab challenge, we used the final results presented in the papers published by the participants rather than the results achieved during the time frame of the competition.

Based on our classification of STI approaches along three paradigms, we observe that heuristic systems appear in the top three systems for all datasets and all tasks. In particular, none of the feature engineering systems or deep learning systems reached the top three in the entity matching tasks (CEA and Row-to-instance). We believe that one of the main reason is that unlike CTA or CPA, which can extract features from columns or column pairs (all column cells can provide features such as entity embeddings, string length, or distribution), features that can be used to annotate individual cells or single rows are relatively rare. As a result, it limits the performance of such systems. Furthermore, the performance of a learning-based classifier is related to the number of candidates used. Annotating an entity means that a classifier should be trained to serve millions of candidates, which makes the task more difficult. From this point of view, rare feature engineering systems or deep learning systems position themselves for the tasks of CEA and row-to-instance. Ideal STI systems are therefore likely to be hybrid systems combining the best of heuristic-based and deep learning based methods.

We further group the selected datasets according to their provenance. These datasets are either synthetic tables automatically generated (e.g., SemTab datasets) or tables collected from the Web (e.g., Limaye and T2D). We observe that heuristic-based systems that follow an iterative approach such as MTab, DAGOBAH, LinkingPark and JenTab, achieve very strong performances on large synthetic datasets such as SemTab. These systems generally use the entire target KG in order to get a very good coverage. Leveraging inter-tasks process for iterative disambiguation further optimizes the performance of the system. As these datasets are synthetically generated from a KG, string matching based approaches have also more advantages since this matching setp is generally very reliable and will not face challenges with KG incompleteness issues. Systems that rely on statistical learning such as Col-Net or Guo et al. [42], on the other hand, perform well on smaller real-world datasets, like Limaye and T2D. It may be because smaller datasets provide a limited set of candidate entities/types/relationships. Short candidate lists make training more accessible, more efficient and more accurate. In addition, heuristic approaches are highly based on the closed world assumption, where KG incompleteness is always a big issue. Deep learning and feature engineering are more robust on this point since they are not highly dependent on the completeness of the KG. That may be one of the reasons explaining that learning-based methods such as ColNet or Guo et al. [42] have been able to become one of the top three systems for real-world datasets like Limaye and T2D.

8. Challenges and Future Directions

While recent works have made significant progress in the field of STI, existing approaches have several limitations: i) they mainly focus on single-cell subject within

Da	taset	(CEA / Row-to-ins	tance	CTA	^a / Topic annotation	n	CPA			
T.I.	b	$TabEL^{c}$	TabEAno ^d [71]	T2K ++	T2K ++	Guo et al	MantisTable	Mulwad et al.	T2K ++	TableMiner+	
Lin	naye-	0.894	0.88	0.87	0.88	0.852	0.84	0.89	0.80	0.76	
т	מפי	TabEAno	Zhang et al.	Kruit et al.	$ColNet^{e}$	Alobaid et al. [5]	MantisTable	T2K ++	Singh et al.	MantisTable	
12D		0.91	0.90	0.89	0.976	0.96	0.95	0.91	0.71	0.51	
	Do	MTab	CSV2KG	Tabularisi	MTab	CSV2KG	Tabularisi	CSV2KG	IDLab	Tabularisi	
	112	0.911	0.883	0.826	1.414	1.376	1.099	0.881	0.877	0.790	
SomTab 2010	R3	MTab	CSV2KG	ADOG	MTab	CSV2KG	Tabularisi	MTab	CSV2KG	Tabularisi	
Semilar 2019		0.970	0.962	0.912	1.956	1.864	1.702	0.844	0.841	0.827	
	R4	MTab	MantisTable	CSV2KG	MTab	CSV2KG	Tabularisi	MTab	CSV2KG	Tabularisi	
	104	0.983	0.973	0.907	2.012	1.846	1.716	0.832	0.830	0.823	
	D1	MTab	LinkingPark	MantisTable	JenTab	LinkingPark	MTab	MTab	LinkingPark	JenTab	
	ILI I	0.987	0.987	0.982	0.962	0.926	0.885	0.971	0.967	0.963	
	Do	MTab	DAGOBAH ^f	LinkingPark	LinkingPark	MTab	DAGOBAH	MTab	LinkingPark	DAGOBAH	
	112	0.995	0.993	0.993	0.984	0.984	0.983	0.997	0.993	0.992	
SomTab 2020	D2	MTab	LinkingPark	DAGOBAH	LinkingPark	MTab	DAGOBAH	MTab	DAGOBAH	bbw [89]	
Semilar 2020	165	0.991	0.986	0.985	0.978	0.976	0.974	0.995	0.993	0.989	
	D4	MTab	LinkingPark	DAGOBAH	MTab	bbw	DAGOBAH	MTab	bbw	DAGOBAH	
	104	0.993	0.985	0.984	0.981	0.98	0.972	0.997	0.995	0.995	
	от	MTab	bbw	DAGOBAH	DAGOBAH	MTab	LinkingPark	-	-	-	
	21	0.907	0.863	0.830	0.743	0.728	0.686	-	-	-	
	P1 (DPpodia)	DAGOBAH	GBMTab	JenTab	JenTab	DAGOBAH	Magic	-	-	-	
	iti (Dispetita)	0.945	0.692	0.607	0.46	0.422	0.159	-	-	-	
	D1 (WileData)	DAGOBAH	MTab	AMALGAM [6]	DAGOBAH	MTab	JenTab	-	-	-	
	Iti (WikiData)	0.923	0.907	0.658	0.832	0.728	0.697	-	-	-	
	D9 Hord	MTab	DAGOBAH	MantisTable	MTab	DAGOBAH	MantisTable	MTab	JenTab	DAGOBAH	
	nz-maru	0.985	0.975	0.968	0.977	0.976	0.955	0.997	0.996	0.996	
	D2 Bio	DAGOBAH	MTab	MantisTable	MTab	Magic	DAGOBAH	MTab	DAGOBAH	JenTab	
SemTab 2021	112-1010	0.970	0.964	0.93	0.956	0.916	0.916	0.947	0.899	0.899	
Sem1ab 2021	B3-Biodiv	JenTab	MTab	DAGOBAH	KEPLER-aSI [7]	DAGOBAH	MTab	-	-	-	
	10-Diodiv	0.602	0.522	0.496	0.593	0.391	0.123	-	-	-	
	B2 Hard	DAGOBAH	MTab	MantisTable	DAGOBAH	MTab	MantisTable	MTab	JenTab	DAGOBAH	
	10-maru	0.974	0.968	0.959	0.99	0.984	0.965	0.993	0.992	0.991	
	B3-Cit (DBn)	-	-	-	DAGOBAH	KEPLER-aSI	MantisTable	-	-	-	
	no on (DDp)	-	-	-	0.07	0.041	0.037	-	-	-	
	B3-Git (Sch)	-	-	-	MantisTable	DAGOBAH	-	-	-	-	
		-	-	-	0.205	0.183	-	-	-	-	

Table 3: Top-3 systems for each dataset and their corresponding F1 score unless otherwise stated in the footnote.

^aFor SemTab 2019, we consider the AH score, while for SemTab 2020 and 2021, we consider the AF1 score

^bWe consider only one of the Limaye subsets named Manual

 $^c\mathrm{TabEL}$ only reports about the accuracy of the system and not the F1 score

 $^d\mathrm{TabEAno}$ is a sub system of MTab

 e ColNet was evaluated on T2Dv2. Thus, we consider only the result from 237 PK columns, which is almost the 233 tables from T2D f This method is named DAGOBAH SL in [49, 50, 51]

relational or entity tables, and make strong assumptions about the coherence and the simplicity of their layout; ii) they are highly confident in both the completeness and the correctness of the target KG; iii) they only partially leverage the information of the table, in both substance and form. Based on these observations, we formulate some possible guidelines to sketch the future directions for the STI field.

8.1. Beyond Simple Table Type

From the literature, we observe that most of the works focus on single-cell subject tables, the simplest type of relational tables. A few approaches focus on entity tables such as the infoboxes from Wikipedia pages [15, 109] but the other types of tables defined in the classification we proposed in Section 3.1 are still hardly handled. Moreover, current approaches do not dig deeper into relational tables complexities such as hidden subjects or composed subjects, to name a few.

As a result, existing systems are far from being generalizable to any table type. To fill the gap and stimulate the search for new solutions, we believe it is important to broaden the spectrum of corpus complexities. To that end, we recommend creating new datasets with multiple table structures and complex contents to tackle the whole diversity of real-world data. We advise that content-level complexity should not be restricted to noise added to mentions, whether synthetically or manually, as these artefacts happened to be not so difficult to handle in the light of the SemTab experience, as well as not so close to real data tables. Introducing numerical mentions with heterogeneous units or lists within cells (multivalued tables), for instance, could be more challenging and therefore beneficial for the community. Last but not least, a ground truth shall be associated with these datasets to allow a fair comparison between the future approaches. This latter requirement suggests prioritizing quality over quantity for evaluation datasets to bootstrap new challenges quickly.

8.2. KG Incompleteness and Incorrectness

Existing approaches assume that the target KG is complete and error-free. As a consequence, an annotation can always be generated even if the correct result is not in the KG, whether it concerns an instance, a type or a relation. This situation can be harmful, especially as it may spread the error from one annotation to the whole column or even the whole table. Suppose for instance a table where a column contains the last names of writers (which can be very common), and another column related to books' titles (for the sake of the example, we assume the majority of these books have been adapted for the cinema). If the target KG covers extensively movies but only a few literature works (or is less accurate on books than movies), the annotation process might lead to type the second column as "film", which could lead to wrongly disambiguate the mentions in the first column (if some related actors have similar last names for example). As a result, this table will be interpreted as an actors-movies item instead of the correct writers-books target.

Some existing mechanisms, such as giving a confidence score for each candidate, can help filtering the incorrect annotations further. The T2Dv2 benchmark, for instance, added negative examples that can be leveraged to that end. However, rare studies focus on this challenge which is far from being trivial as it implies to have the capability to identify the KG coverage w.r.t. the tables to be processed, as well as to detect the possible errors. To improve both the completeness and the correctness, we believe that leveraging multiple KGs is the first step to make. Indeed, it could enhance the coverage and provide a basis for confidence scoring through popularity computation. However, evaluation procedures should be discussed and updated as judging from different sources might be challenging. Finally, we emphasize that future approaches should also consider to tackle domains where only nascent KGs exist, with the objective of using STI to augment these KGs in a virtuous iterative loop.

In Table 3, many annotation systems (e.g., MTab,-DAGOBAH SL) achieve very high performances on some synthetic datasets (e.g., SemTab 2020, SemTab 2021 R2-Hard). This can be explained since synthetic tables are automatically and accurately generated using a reference KG (see Section 6). Therefore, their content is almost fully represented in the KG and provides rich discriminative information for the disambiguation. In contrast, manually curated tables (R3-BiodivTable), complex tables (Tough-Table), or tables coming from diverse and specific domains (R3-GitTables) are still particularly challenging for annotation systems. In real situations, KGs are often incomplete. As a consequence, an existing entity may not be fully described in a KG (e.g., lack of literal attributes for a given entity), or an unpopular/heterogeneous domain may provide little information (e.g., R3-GitTables is made of tables from GitHub). Hence, the graph context provided by the KG can sometimes be insufficient for disambiguating tables in R3-BiodivTable, ToughTable, and R3-GitTables, which are much more ambiguous than synthetic tables. We argue that annotation systems should enrich and better handle both explicit and implicit contextual information by exploiting knowledge graph reasoning or table representation learning (e.g., Transformer) to improve the performance on these kinds of table datasets.

8.3. Table Context

We observe that many approaches leverage only partially the elements of the table (see Table 1), even if more recent ones tend to extend their view. As we discussed in Section 5.4.3, we believe that leveraging as many elements as possible should increase the accuracy by adding more contextual information. In that sense, language models generated from transformers could be better used. Indeed, one could consider a table as a way of structuring the language: in the simplest case, one table row can be seen as a sentence describing a subject with some attributes. The same applies to the corresponding sub-graph in the target KG. Thus, sentence representation could be used to compute similarities. Nonetheless, the specificity of tabular data as well as KGs should be taken into account, which implies adapting attention mechanisms to this very structure. The visibility matrix used in [34] is an attempt to do so in relational tables, but it should be extended to other types of tabular data.

We also notice that most approaches treat tables independently. However, some tables are related to each other since they can be generated with the same template, be part of a coherent corpus of tables or related to keys such as SQL database tables. Inter-table relations as studied by [104] can constitute an interesting complementary approach with appropriate target tables. We believe that STI systems could significantly take advantage of combining table elements with inter-table connections, which can be considered as another context layer added to capture richer prior information about the data to be processed.

9. Conclusion

The last few years have seen significant growth in the field of STI, with the development of many new approaches and the proposal of ever more complete datasets, notably under the impulse of initiatives such as the SemTab challenge. In this survey, we have provided a comprehensive and up-to-date overview of the STI field. Our work also gathers and proposes a set of definitions to structure and unify the field. It first defines the inputs (the different types of tables, their context and their metadata) of the STI process as well as the KGs used both as a support and as a repository to be enriched via the STI process. Then, we propose a generic pipeline for STI and a description of the different tasks performed by existing approaches. These are classified into three families, heuristics, feature engineering and deep learning, with an emphasis on the strengths and weaknesses of each.

We have also summarized what are the performances of the various approaches on the datasets that have been proposed in the community. We have highlighted the best performing systems for each dataset. Finally, we have listed several challenges to address for improving STI systems. We observe that recent works have tried to develop modern web-based user interfaces on top of STI systems such as [87] which will undoubtedly empower end-users when adopting these systems.

References

- Abdelmageed, N., Schindler, S.: JenTab: Matching Tabular Data to Knowledge Graphs. In: Semantic Web Challenge on Tabular Data to Knowledge Graph Matching (SemTab). pp. 40–49 (2020)
- [2] Abdelmageed, N., Schindler, S.: JenTab: A Toolkit for Semantic Table Annotations. In: 2^{nd} International Workshop on Knowledge Graph Construction (2021)
- [3] Abdelmageed, N., Schindler, S.: JenTab Meets SemTab 2021s New Challenges. In: Semantic Web Challenge on Tabular Data to Knowledge Graph Matching (SemTab) (2021)
- [4] Abdelmageed, N., Schindler, S., König-Ries, B.: BiodivTab: A Tabular Benchmark based on Biodiversity Research Data. In: Semantic Web Challenge on Tabular Data to Knowledge Graph Matching (SemTab) (2021)
- [5] Alobaid, A., Corcho, O.: Balancing coverage and specificity for semantic labelling of subject columns. Knowledge-Based Systems p. 108092 (2022)
- [6] Azzi, R., Diallo, G.: AMALGAM: making tabular dataset explicit with knowledge graph. In: Semantic Web Challenge on Tabular Data to Knowledge Graph Matching (SemTab). pp. 9–16 (2020)
- [7] Baazouzi, W., Kachroudi, M., Faiz, S.: KEPLER-asi at SemTab 2021. In: Semantic Web Challenge on Tabular Data to Knowledge Graph Matching (SemTab) (2021)
- [8] Badaro, G., Saeed, M., Papotti, P.: Transformers for Tabular Data Representation: A survey of models and applications. Tech. rep., EURECOM (2021), https://www.eurecom. fr/publication/6721
- [9] Banko, M., Etzioni, O.: The tradeoffs between open and traditional relation extraction. In: ACL-08: HLT. pp. 28–36 (2008)
- [10] Belleau, F., Nolin, M.A., Tourigny, N., Rigault, P., Morissette, J.: Bio2RDF: towards a mashup to build bioinformatics knowledge systems. Journal of biomedical informatics 41(5), 706-716 (2008)
- [11] Benjelloun, O., Chen, S., Noy, N.: Google Dataset Search by the Numbers. In: International Semantic Web Conference (ISWC), In-Use Track (2020)
- Berners-Lee, T.: Linked Data Design Issues. http://www.w3. org/DesignIssues/LinkedData.html (2006)
- [13] Bhagavatula, C.S., Noraset, T., Downey, D.: Methods for exploring and mining tables on wikipedia. In: ACM SIGKDD workshop on interactive data exploration and analytics. pp. 18–26 (2013)
- [14] Bhagavatula, C.S., Noraset, T., Downey, D.: Tabel: Entity linking in web tables. In: 14th International Semantic Web Conference. pp. 425–441. Springer (2015)
- [15] Biswas, R., Türker, R., Moghaddam, F.B., Koutraki, M., Sack, H.: Wikipedia Infobox Type Prediction Using Embeddings. In: DL4KGS@ ESWC. pp. 46–55 (2018)
- [16] Bizer, C., Heath, T., Berners-Lee, T.: Linked data: The story so far. In: International Journal on Semantic Web and Information Systems, pp. 205–227. IGI global (2009)
- [17] Bizer, C., Lehmann, J., Kobilarov, G., Auer, S., Becker, C., Cyganiak, R., Hellmann, S.: DBpedia-A crystallization point for the Web of Data. Journal of web semantics 7(3), 154–165 (2009)
- [18] Bollacker, K., Evans, C., Paritosh, P., Sturge, T., Taylor, J.: Freebase: a collaboratively created graph database for structuring human knowledge. In: ACM SIGMOD International Conference on Management of Data (2008)
- [19] Bonfitto, S., Casiraghi, E., Mesiti, M.: Table understanding approaches for extracting knowledge from heterogeneous tables. Wiley Interdisciplinary Reviews: Data Mining and Knowledge Discovery p. e1407 (2021)

- [20] Bordes, A., Usunier, N., Garcia-Duran, A., Weston, J., Yakhnenko, O.: Translating embeddings for modeling multirelational data. Advances in neural information processing systems 26 (2013)
- [21] Borisov, V., Leemann, T., Seßler, K., Haug, J., Pawelczyk, M., Kasneci, G.: Deep neural networks and tabular data: A survey. arXiv:2110.01889 (2021)
- [22] Cafarella, M.J., Halevy, A., Wang, D.Z., Wu, E., Zhang, Y.: Webtables: exploring the power of tables on the web. In: VLDB Endowment. pp. 538–549. VLDB Endowment (2008)
- [23] Chabot, Y., Labbe, T., Liu, J., Troncy, R.: DAGOBAH: an end-to-end context-free tabular data semantic annotation system. In: Semantic Web Challenge on Tabular Data to Knowledge Graph Matching. pp. 41–48 (2019)
- [24] Chen, J., Jimenez-Ruiz, E., Horrocks, I., Sutton, C.: ColNet: Embedding the Semantics of Web Tables for Column Type Prediction. In: 33^{rd} AAAI International Conference on Artificial Intelligence (2018)
- [25] Chen, S., Karaoglu, A., Negreanu, C., Ma, T., Yao, J.G., Williams, J., Gordon, A., Lin, C.Y.: Linkingpark: An integrated approach for semantic table interpretation. In: Semantic Web Challenge on Tabular Data to Knowledge Graph Matching (SemTab) (2020)
- [26] Cremaschi, M., Avogadro, R., Barazzetti, A., Chieregato, D.: MantisTable SE: an Efficient Approach for the Semantic Table Interpretation. In: Semantic Web Challenge on Tabular Data to Knowledge Graph Matching (SemTab) (2020)
- [27] Cremaschi, M., Avogadro, R., Chieregato, D.: MantisTable: an Automatic Approach for the Semantic Table Interpretation. In: Semantic Web Challenge on Tabular Data to Knowledge Graph Matching (SemTab). pp. 15–24 (2019)
- [28] Cremaschi, M., De Paoli, F., Rula, A., Spahiu, B.: A fully automated approach to a complete semantic table interpretation. Future Generation Computer Systems 112, 478–500 (2020)
- [29] Cremaschi, M., Siano, A., Avogadro, R., Jimenez-Ruiz, E., Maurino, A.: STILTool: a semantic table interpretation evaluation tool. In: European Semantic Web Conference. pp. 61–66. Springer (2020)
- [30] Crestan, E., Pantel, P.: Web-scale table census and classification. In: 4th ACM International Conference on Web Search and Data Mining (WSDM). pp. 545–554 (2011)
- [31] Cutrona, V., Bianchi, F., Jiménez-Ruiz, E., Palmonari, M.: Tough tables: Carefully evaluating entity linking for tabular data. In: 19th International Semantic Web Conference (ISWC). pp. 328–343. Springer (2020)
- [32] Cutrona, V., Chen, J., Efthymiou, V., Hassanzadeh, O., Jiménez-Ruiz, E., Sequeda, J., Srinivas, K., Abdelmageed, N., Hulsebos, M., Oliveira10, D., et al.: Results of SemTab 2021. In: CEUR Workshop Proceedings (2021)
- [33] Deng, D., Jiang, Y., Li, G., Li, J., Yu, C.: Scalable column concept determination for web tables using large knowledge bases. In: PVLDB. pp. 1606–1617. VLDB Endowment (2013)
- [34] Deng, X., Sun, H., Lees, A., Wu, Y., Yu, C.: TURL: Table Understanding through Representation Learning. arXiv:2006.14806 (2020)
- [35] Eberius, J., Braunschweig, K., Hentsch, M., Thiele, M., Ahmadov, A., Lehner, W.: Building the dresden web table corpus: A classification approach. In: IEEE 2nd International Symposium on Big Data Computing (BDC). pp. 41–50. IEEE (2015)
- [36] Eberius, J., Thiele, M., Braunschweig, K., Lehner, W.: Top-k Entity Augmentation Using Consistent Set Covering. In: SS-DBM (2015)
- [37] Efthymiou, V., Hassanzadeh, O., Rodriguez-Muro, M., Christophides, V.: Matching web tables with knowledge base entities: from entity lookups to entity embeddings. In: 16th International Semantic Web Conference (ISWC). pp. 260–277. Springer (2017)
- [38] Ehrlinger, L., Wöß, W.: Towards a Definition of Knowledge Graphs. SEMANTICS 48(1-4), 2 (2016)
- [39] Ermilov, I., Ngomo, A.C.N.: TAIPAN: automatic property mapping for tabular data. In: European Knowledge Acquisi-

tion Workshop. pp. 163–179. Springer (2016)

- [40] Guo, J., Fan, Y., Ai, Q., Croft, W.B.: A deep relevance matching model for ad-hoc retrieval. In: the 25th ACM international on conference on information and knowledge management. pp. 55–64 (2016)
- [41] Guo, Q., Qiu, X., Liu, P., Shao, Y., Xue, X., Zhang, Z.: Startransformer. arXiv:1902.09113 (2019)
- [42] Guo, T., Shen, D., Nie, T., Kou, Y.: Web table column type detection using deep learning and probability graph model. In: International Conference on Web Information Systems and Applications. pp. 401–414. Springer (2020)
- [43] Habibi, M., Starlinger, J., Leser, U.: DeepTable: a permutation invariant neural network for table orientation classification. Data Mining and Knowledge Discovery pp. 1–21 (2020)
- [44] Ho, V.T., Pal, K., Weikum, G.: QuTE: Answering Quantity Queries from Web Tables. In: International Conference on Management of Data. pp. 2740–2744 (2021)
- [45] Hogan, A., Blomqvist, E., Cochez, M., d'Amato, C., de Melo, G., Gutierrez, C., Gayo, J.E.L., Kirrane, S., Neumaier, S., Polleres, A., et al.: Knowledge graphs. arXiv:2003.02320 (2020)
- [46] Hu, K., Gaikwad, S., Hulsebos, M., Bakker, M.A., Zgraggen, E., Hidalgo, C., Kraska, T., Li, G., Satyanarayan, A., Demiralp, Ç.: Viznet: Towards a large-scale visualization learning and benchmarking repository. In: CHI Conference on Human Factors in Computing Systems. pp. 1–12 (2019)
- [47] Hulsebos, M., Demiralp, c., Groth, P.: GitTables: A Large-Scale Corpus of Relational Tables. arXiv:2106.07258 (2021)
- [48] Hulsebos, M., Hu, K., Bakker, M., Zgraggen, E., Satyanarayan, A., Kraska, T., Demiralp, Ç., Hidalgo, C.: Sherlock: A deep learning approach to semantic data type detection. In: 25th ACM International Conference on Knowledge Discovery & Data Mining (SIGKDD). pp. 1500–1508 (2019)
- [49] Huynh, V.P., Chabot, Y., Labbé, T., Liu, J., Troncy, R.: From Heuristics to Language Models: A Journey Through the Universe of Semantic Table Interpretation with DAGOBAH. In: Semantic Web Challenge on Tabular Data to Knowledge Graph Matching (SemTab) (2022)
- [50] Huynh, V.P., Liu, J., Chabot, Y., Deuzé, F., Labbé, T., Monnin, P., Troncy, R.: DAGOBAH: Table and Graph Contexts for Efficient Semantic Annotation of Tabular Data. In: Semantic Web Challenge on Tabular Data to Knowledge Graph Matching (SemTab) (2021)
- [51] Huynh, V.P., Liu, J., Chabot, Y., Labbé, T., Monnin, P., Troncy, R.: DAGOBAH: Enhanced Scoring Algorithms for Scalable Annotations of Tabular Data. In: Semantic Web Challenge on Tabular Data to Knowledge Graph Matching (SemTab) (2020)
- [52] Ibrahim, Y., Riedewald, M., Weikum, G., Zeinalipour-Yazti, D.: Bridging quantities in tables and text. In: 2019 IEEE 35th International Conference on Data Engineering (ICDE). pp. 1010–1021. IEEE (2019)
- [53] Jiménez-Ruiz, E., Hassanzadeh, O., Efthymiou, V., Chen, J., Srinivas, K.: SemTab 2019: Resources to Benchmark Tabular Data to Knowledge Graph Matching Systems. In: European Semantic Web Conference (ESWC). pp. 514–530. Springer (2020)
- [54] Jiménez-Ruiz, E., Hassanzadeh, O., Efthymiou, V., Chen, J., Srinivas, K., Cutrona, V.: Results of SemTab 2020. In: CEUR Workshop Proceedings. vol. 2775, pp. 1–8 (2020)
- [55] Kacprzak, E., Giménez-García, J.M., Piscopo, A., Koesten, L., Ibáñez, L.D., Tennison, J., Simperl, E.: Making sense of numerical data-semantic labelling of web tables. In: European Knowledge Acquisition Workshop. pp. 163–178. Springer (2018)
- [56] Khurana, U., Galhotra, S.: Semantic annotation for tabular data (2019)
- [57] Lautert, L.R., Scheidt, M.M., Dorneles, C.F.: Web table taxonomy and formalization. ACM SIGMOD Record 42(3), 28–33 (2013)
- [58] Lehmberg, O., Ritze, D., Meusel, R., Bizer, C.: A large public

corpus of web tables containing time and context metadata. In: 25^{th} International Conference Companion on World Wide Web. pp. 75–76 (2016)

- [59] Lehmberg, O., Ritze, D., Meusel, R., Bizer, C.: A Large Public Corpus of Web Tables containing Time and Context Metadata.
 In: 25th International Conference Companion on World Wide Web (WWW Companion). pp. 75–76 (2016)
- [60] Limaye, G., Sarawagi, S., Chakrabarti, S.: Annotating and searching web tables using entities, types and relationships. Proceedings of the VLDB Endowment 3(1-2), 1338–1347 (2010)
- [61] Liu, J., Huynh, V.P., Chabot, Y., Troncy, R.: Radar Station: Using KG Embeddings for Semantic Table Interpretation and Entity Disambiguation. In: 21st International Semantic Web Conference (ISWC) (2022)
- [62] Martinez-Rodriguez, J.L., Hogan, A., Lopez-Arevalo, I.: Information extraction meets the semantic web: a survey. Semantic Web Journal 11(2), 255–235 (2020)
- [63] Mendes, P.N., Jakob, M., García-Silva, A., Bizer, C.: DBpedia spotlight: shedding light on the web of documents. In: 7th International Conference on Semantic Systems. pp. 1–8 (2011)
- [64] Minh, P., Suresh, A., Craig, A.K., Pedro, S.: Semantic Labeling: A Domain-Independent Approach. In: 15th International Semantic Web Conference (ISWC). pp. 446–462 (2016)
- [65] Morikawa, H.: Semantic Table Interpretation using LOD4ALL. In: Semantic Web Challenge on Tabular Data to Knowledge Graph Matching (SemTab). pp. 49–56 (2019)
- [66] Mulwad, V., Finin, T., Joshi, A.: Semantic message passing for generating linked data from tables. In: 12th International Semantic Web Conference (ISWC). pp. 363–378. Springer (2013)
- [67] Mulwad, V., Finin, T., Syed, Z., Joshi, A., et al.: Using linked data to interpret tables. In: 1st International Workshop on Consuming Linked Data (COLD) (2010)
- [68] Nakashole, N., Weikum, G., Suchanek, F.: PATTY: A taxonomy of relational patterns with semantic types. In: Empirical Methods in Natural Language Processing and Computational Natural Language Learning. pp. 1135–1145 (2012)
- [69] Neumaier, S., Umbrich, J., Parreira, J.X., Polleres, A.: Multilevel semantic labelling of numerical values. In: 15th International Semantic Web Conference (ISWC). pp. 428–445 (2016)
- [70] Nguyen, P., Kertkeidkachorn, N., Ichise, R., Takeda, H.: MTab: Matching Tabular Data to Knowledge Graph using Probability Models. In: Semantic Web Challenge on Tabular Data to Knowledge Graph Matching (SemTab) (2019)
- [71] Nguyen, P., Kertkeidkachorn, N., Ichise, R., Takeda, H.: TabEAno: table to knowledge graph entity annotation. arXiv:2010.01829 (2020)
- [72] Nguyen, P., Nguyen, K., Ichise, R., Takeda, H.: Embnum: Semantic labeling for numerical values with deep metric learning. In: Joint International Semantic Technology Conference. pp. 119–135. Springer (2018)
- [73] Nguyen, P., Yamada, I., Kertkeidkachorn, N., Ichise, R., Takeda, H.: Mtab4wikidata at semtab 2020: Tabular data annotation with wikidata. In: Semantic Web Challenge on Tabular Data to Knowledge Graph Matching (SemTab) (2020)
- [74] Nguyen, P., Yamada, I., Kertkeidkachorn, N., Ichise, R., Takeda, H.: SemTab 2021: Tabular Data Annotation with MTab Tool. In: Semantic Web Challenge on Tabular Data to Knowledge Graph Matching (SemTab) (2021)
- [75] Noy, N., Burgess, M., Brickley, D.: Google Dataset Search: Building a search engine for datasets in an open Web ecosystem. In: 28th The Web Conference (WWW) (2019)
- [76] Oliveira, D., d'Aquin, M.: Adog-annotating data with ontologies and graphs. In: Semantic Web Challenge on Tabular Data to Knowledge Graph Matching (SemTab) (2019)
- [77] Pasupat, P., Liang, P.: Compositional semantic parsing on semi-structured tables. arXiv:1508.00305 (2015)
- [78] Penn, G., Hu, J., Luo, H., McDonald, R.: Flexible web document analysis for delivery to narrow-bandwidth devices. In: 6th International Conference on Document Analysis and Recognition. pp. 1074–1078. IEEE (2001)

- [79] Ramnandan, S.K., Mittal, A., Knoblock, C.A., Szekely, P.: Assigning semantic labels to data sources. In: European Semantic Web Conference (ESWC). pp. 403–417. Springer (2015)
- [80] Riley, J.: Understanding metadata. http://www.niso.org/ publications/press/UnderstandingMetadata.pdf (2017)
- [81] Ristoski, P., Paulheim, H.: Rdf2vec: Rdf graph embeddings for data mining. In: International Semantic Web Conference. pp. 498–514. Springer (2016)
- [82] Ritze, D.: Web-scale web table to knowledge base matching. Ph.D. thesis, University of Mannheim (2017)
- [83] Ritze, D., Bizer, C.: Matching Web Tables To DBpedia A Feature Utility Study. In: International Conference on Extending Database Technology (EDBT). pp. 210–221 (2017)
- [84] Ritze, D., Lehmberg, O., Bizer, C.: Matching html tables to dbpedia. In: 5th International Conference on Web Intelligence, Mining and Semantics. pp. 1–6 (2015)
- [85] Ritze, D., Lehmberg, O., Bizer, C.: Matching HTML Tables to DBpedia. In: 5th International Conference on Web Intelligence, Mining and Semantics (WIMS). pp. 1–6 (2015)
- [86] Robertson, S., Zaragoza, H.: The probabilistic relevance framework: BM25 and beyond. Now Publishers Inc (2009)
- [87] Sarthou-Camy, C., Jourdain, G., Chabot, Y., Monnin, P., Deuzé, Huynh, V.P., Liu, J., Labbé, T., Troncy, R.: DAGOBAH UI: A New Hope For Semantic Table Interpretation. In: 19th European Semantic Web Conference (ESWC), Poster and Demo Track (2022)
- [88] Sekhavat, Y.A., Di Paolo, F., Barbosa, D., Merialdo, P.: Knowledge base augmentation using tabular data. In: LDOW (2014)
- [89] Shigapov, R., Zumstein, P., Kamlah, J., Oberländer, L., Mechnich, J., Schumm, I.: bbw: Matching CSV to Wikidata via Meta-lookup. In: Semantic Web Challenge on Tabular Data to Knowledge Graph Matching (SemTab). vol. 2775, pp. 17– 26. RWTH (2020)
- [90] Singh, G., Singh, S., Wong, J., Saffari, A.: Relation Extraction from Tables using Artificially Generated Metadata. arXiv:2108.10750 (2021)
- [91] Steenwinckel, B., De Turck, F., Ongeane, F.: MAGIC: Mining an Augmented Graph using INK, starting from a CSV. In: Semantic Web Challenge on Tabular Data to Knowledge Graph Matching (SemTab) (2021)
- [92] Steenwinckel, B., Vandewiele, G., De Turck, F., Ongenae, F.: Csv2kg: Transforming tabular data into semantic knowledge.
 In: Semantic Web Challenge on Tabular Data to Knowledge Graph Matching (SemTab) (2019)
- [93] Suhara, Y., Li, J., Li, Y., Zhang, D., Demiralp, Ç., Chen, C., Tan, W.C.: Annotating Columns with Pre-trained Language Models. arXiv:2104.01785 (2021)
- [94] Sun, H., Ma, H., He, X., Yih, W.t., Su, Y., Yan, X.: Table cell search for question answering. In: 25th International Conference on World Wide Web. pp. 771–782 (2016)
- [95] Talukdar, P.P., Wijaya, D., Mitchell, T.: Acquiring temporal constraints between relations. In: the 21st ACM international conference on Information and knowledge management. pp. 992–1001 (2012)
- [96] Tanon, T.P., Weikum, G., Suchanek, F.: YAGO 4: A Reasonable Knowledge Base. In: European Semantic Web Conference (ESWC). pp. 583–596. Springer (2020)
- [97] Tao, C., Embley, D.W.: Automatic hidden-web table interpretation, conceptualization, and semantic annotation. Data & Knowledge Engineering 68(7), 683–703 (2009)
- [98] Thawani, A., Hu, M., Hu, E., Zafar, H., Divvala, N.T., Singh, A., Qasemi, E., Szekely, P.A., Pujara, J.: Entity Linking to Knowledge Graphs to Infer Column Types and Properties. In: Semantic Web Challenge on Tabular Data to Knowledge Graph Matching (SemTab). vol. 2019, pp. 25–32 (2019)
- [99] Thornton, P.K., Stroud, A., Hatibu, N., Legg, C., Ly, S., Twomlow, S., Molapong, K., Notenbaert, A., Kruska, R., von Kaufmann, R.: Site selection to test an integrated approach to agricultural research for development: combining expert knowledge and participatory Geographic Information System

methods. International Journal of Agricultural Sustainability 4(1), 39–60 (2006)

- [100] Van Eeden, W., De Villiers, J.P., Berndt, R., Nel, W.A., Blasch, E.: Micro-Doppler radar classification of humans and animals in an operational environment. Expert Systems with Applications 102, 1–11 (2018)
- [101] Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A.N., Kaiser, L., Polosukhin, I.: Attention is all you need. arXiv:1706.03762 (2017)
- [102] Venetis, P., Halevy, A.Y., Madhavan, J., Pasca, M., Shen, W., Wu, F., Miao, G.: Recovering semantics of tables on the web. PVLDB 4(9), 528–538 (2011)
- [103] Vrandečić, D., Krötzsch, M.: Wikidata: a free collaborative knowledge base. Communications of the ACM 57(10), 78–85 (2014)
- [104] Wang, D., Shiralkar, P., Lockard, C., Huang, B., Dong, X.L., Jiang, M.: Tcn: Table convolutional network for web table interpretation. arXiv:2102.09460 (2021)
- [105] Wang, J., Wang, H., Wang, Z., Zhu, K.Q.: Understanding tables on the web. In: International Conference on Conceptual Modeling. pp. 141–155. Springer (2012)
- [106] Wang, Y., Hu, J.: Detecting tables in HTML documents. In: 5th IAPR International Workshop on Document Analysis Systems. pp. 249–260. Springer (2002)
- [107] Wang, Z., Zhang, J., Feng, J., Chen, Z.: Knowledge graph embedding by translating on hyperplanes. In: AAAI Conference on Artificial Intelligence (2014)
- [108] Wang, Z., Dong, H., Jia, R., Li, J., Fu, Z., Han, S., Zhang, D.: TUTA: Tree-based Transformers for Generally Structured Table Pre-training. In: 27th ACM SIGKDD Conference on Knowledge Discovery & Data Mining. pp. 1780–1790 (2021)
- [109] Wu, F., Weld, D.S.: Automatically refining the wikipedia infobox ontology. In: 17th international conference on World Wide Web. pp. 635–644 (2008)
- [110] Wu, W., Li, H., Wang, H., Zhu, K.Q.: Inferring a universal probabilistic taxonomy from the web. Tech. rep., Technical report, Microsoft Research (2010)
- [111] Xiaoxue, L., Xuesong, B., Longhe, W., Bingyuan, R., Shuhan, L., Lin, L.: Review and trend analysis of knowledge graphs for crop pest and diseases. IEEE Access 7, 62251–62264 (2019)
- [112] Yakout, M., Ganjam, K., Chakrabarti, K., Chaudhuri, S.: Infogather: entity augmentation and attribute discovery by holistic matching with web tables. In: ACM SIGMOD International Conference on Management of Data. pp. 97–108 (2012)
- [113] Yin, P., Neubig, G., Yih, W.t., Riedel, S.: TaBERT: Pretraining for Joint Understanding of Textual and Tabular Data (2020)
- [114] Yusof, M.M., Rosli, N.F., Othman, M., Mohamed, R., Abdullah, M.H.A.: M-DCocoa: M-agriculture expert system for diagnosing cocoa plant diseases. In: International Conference on Soft Computing and Data Mining. pp. 363–371. Springer (2018)
- [115] Zhang, A., Gourley, D.: Creating digital collections: a practical guide. Elsevier (2008)
- [116] Zhang, D., Suhara, Y., Li, J., Hulsebos, M., Demiralp, Ç., Tan, W.C.: Sato: Contextual Semantic Type Detection in Tables (2019)
- [117] Zhang, S., Balog, K.: Recommending related tables. arXiv:1907.03595 (2019)
- [118] Zhang, S., Balog, K.: Web table extraction, retrieval, and augmentation: A survey. ACM Transactions on Intelligent Systems and Technology (TIST) 11(2), 1–35 (2020)
- [119] Zhang, S., Meij, E., Balog, K., Reinanda, R.: Novel entity discovery from web tables. In: The Web Conference. pp. 1298– 1308 (2020)
- [120] Zhang, Z.: Towards efficient and effective semantic table interpretation. In: 3th International Semantic Web Conference. pp. 487–502. Springer (2014)
- [121] Zhang, Z.: Effective and efficient semantic table interpretation using tableminer+. Semantic Web 8(6), 921–957 (2017)
- [122] Zhou, Y., Singh, S., Christodoulopoulos, C.: Tabular Data

Concept Type Detection Using Star-Transformers. In: Proceedings of the 30th ACM International Conference on Information & Knowledge Management. pp. 3677–3681 (2021)
[123] Zwicklbauer, S., Einsiedler, C., Granitzer, M., Seifert, C.: To-

- [123] Zwicklbauer, S., Einsiedler, C., Granitzer, M., Seifert, C.: Towards Disambiguating Web Tables. In: International Semantic Web Conference (Posters & Demos). pp. 205–208 (2013)
- [124] Zwicklbauer, S., Seifert, C., Granitzer, M.: Do We Need Entity-Centric Knowledge Bases for Entity Disambiguation? In: 13th International Conference on Knowledge Management and Knowledge Technologies. pp. 1–8 (2013)