Global Journals $end{transformula} \mathbb{A}T_{\mathbf{E}} X$ JournalKaleidoscope

Artificial Intelligence formulated this projection for compatibility purposes from the original article published at Global Journals. However, this technology is currently in beta. *Therefore, kindly ignore odd layouts, missed formulae, text, tables, or figures.*

CrossRef DOI of original article:

1	Entity Matching for Digital World: A Modern Approach using
2	Artificial Intelligence and Machine Learning
3	K. Victor Rajan ¹ and Edward Lambert ²
4	¹ Atlantic International University, USA
5	Received: 1 January 1970 Accepted: 1 January 1970 Published: 1 January 1970

7 Abstract

Entity matching is the field of research solving the problem of identifying similar records 8 which refer to the same real-world entity. In today?s digital world, business organizations deal 9 with large amount of data like customers, vendors, manufacturers, etc. Entities are spread 10 across various data sources and failure to correlate two records as one entity can lead to 11 confusion. Relationships and patterns would be missed. Aggregations and calculations won?t 12 make any sense. It is a significant data integration effort that often arises when data originate 13 from different sources. In such scenarios, we understand the situation by linking records and 14 then track entities from a person to a product, etc. There is appreciable value in integrating 15 the data silos across various industries. 16

17

18 Index terms— entity matching, entity resolution, record linkage, de-duplication, machine learning

¹⁹ 1 Introduction

ur world is moving towards digitized business. This opens up numerous avenues to increase revenue through 20 digital marketing, sales forecast, etc. Huge amount of historical data is available to analyze customer behavior, 21 buying patterns and make predictions for future. However, it also comes with challenges along the way. A 22 substantial amount of the value to be harvested from digitization depends on successful integration of large 23 24 volume of data from different sources. Unfortunately, many of the existing data sources do not share a common 25 frame of reference. For example, let us say, a marketing team wants to use statistics from retail stores, e-commerce sites etc., to find out potential buyers for a product. Sadly, these two systems do not refer to customers in the 26 same way -i.e., there are no common identifiers or names across the two systems. Duplicate emails or messages 27 may be sent to same customer again and again unless customer records are tagged uniquely. Recommendations 28 to a customer and an effective marketing scheme cannot be performed based on distinct data silos. A group of 29 similar problems has been studied for a long time in a variety of fields under different names like entity resolution, 30 de-duplication etc. Entity matching is the field of research dedicated to solving the problem of matching which 31 records refer to the same real-world entity. Organizations often struggle with a plethora of customer data captured 32 multiple times in different sources by various people in their own ways. Despite having been studied for decades, 33 entity matching remains a challenging problem in practice. In general, there are several factors that make it 34 35 difficult to solve: 36 Poor Data Quality: Real-world data is seldom completely structured, cleansed, and homogeneous. Data

originating from manual insertion may contain alternative spellings, typos, or fail to comply with the schema (e.g., mixing of first and last name).

Dependency on Human Knowledge: Same data may be represented in different formats by various users like abbreviations, suffixes, prefixes, etc. To perform matching, our solution must interact with human experts and make use of their knowledge. Human interaction in itself is a complex domain.

For example, let's look at a customer table from which analyst is trying to identify distinct customers. Without manual inspection and good understanding of geographical locations, it is difficult to guess whether record 2 is

duplicate of 1 or 3. Somewhat ironically, as often pointed out, entity matching suffers from the problem of 44 being referenced by different names, some referring to the exact same problem, while others are slight variations, 45 generalizations, or specializations. In addition, the names are also not used completely consistently. Deduplication 46 47 or duplicate detection is the problem of identifying records in the same data source that refer to the same entity and can be seen as the special case 1 = 2. Given such representation variations, an unprecedented number of 48 permutations and combinations, the entity matching would be a herculean job when we handle large volume of 49 data. Artificial intelligence and machine learning has become an essential part of multiple research fields in recent 50 years, most notably in natural language processing and computer vision, which are concerned with unstructured 51 data. Its most prominent advantage over systematic approaches is its ability to learn features instead of relying 52

53 on step-by-step calculations.

$_{54}$ 2 a) Problem Definition

Researchers have already realized the potential advantage of machine learning for entity matching. In this paper,
 we aim to propose a machine learning model for entity matching.

Let E be a data source containing entities. E has the attributes (??1,??2, ...,????), and we denote entities as e = (e1, e2, ..., e??) ? E. A data source is a set of records, and a record is a tuple having a specific schema of attributes. An attribute is defined by the intended semantics of its values. So, entities e?? = e?? if and only if attributes ???? of e?? are intended to carry the same information as attributes a?? of e??, and the specific syntactics of the attribute values are irrelevant. Attributes can also have metadata (like a name) associated with them, but this does not affect the equality between them.

The goal of entity matching is to find the largest possible binary relation ?? ? $E \times E$ such that ?? and ?? refer to the same entity for all (??, ??) ? ??. In other words, we would like to find all record pairs across data source that refer to the same entity. We define an entity to be something of unique existence. Attribute values are often assumed to be strings, but that is not always the case. The records are assumed to operate with the same taxonomic granularity. In this research, we will stick to the definition of deduplication (or duplicate detection) as the problem of identifying which records in the same data source refer to the same entity.

The remainder of this paper is organized as follows. We discuss related work in section 2. In Section 3, we formally formulate the problem and propose our methodology. Section 4 describes how our approach is used to detect similarity in a real-world data set and the results of our experiment are explained. Finally, the paper is concluded in Section 5.

73 **3 II.**

74 4 Related Work

Entity resolution, record linkage, deduplication and entity matching are frequently used for more or less the same 75 problem as we mentioned earlier. It is a technique to identify data records in a single data source or across 76 multiple data sources that refer to the same real-world entity and to correlate the records together. In entity 77 matching, the strings that are nearly identical, but not exactly the same, are matched without explicitly having a 78 unique identifier. Entity matching is crucial as it matches non-identical records despite all the data inconsistencies 79 without the constant need for formulating rules. By combining databases using fuzzy matching, we can refine 80 the data and analyze the information. Comparing big data records having nonstandard and inconsistent data 81 from diverse sources that do not provide any unique identifier is a complex problem. In this section, we present 82 an overview of the previous work done by researchers in entity matching. 83

⁸⁴ 5 Researchers use two major techniques as shown below:

Rule-Based: Rule-based systems perform matching based on a set of manually crafted rules. To match any two records of the same entity, various string-based comparison rules are defined. Each record then would run with every other record on all these rules to decide if the two are identical.

Automatic: These systems rely on machine learning algorithms to learn from data. Computers first learn from data provided for training so that they can later make predictions on unknown input data items.

Usually, a rule-based system uses a set of human-crafted rules to help identify subjectivity. As the number of 90 records increases, the number of comparisons increases exponentially in rule-based systems. With large volume 91 of records, rule-based data matching becomes computationally challenging and unscalable. Automatic methods, 92 93 contrary to rule-based systems, do not rely on manually crafted rules but on machine learning algorithms. There 94 has been an uptick in interest on machine learning as a solution for entity matching in recent years. We note 95 that this process is machine-oriented and does not highlight any iterative human interactions or feedback loops. 96 First, there are several books that provide an overview. Christen [15] is a dedicated and comprehensive source on entity matching. Anhai Doan et al. [2] and Talburt [10] introduce entity matching in the context of data quality 97 and integration. Quite early on, statisticians dominated the field of entity matching. Probabilistic methods 98 were first developed by Newcombe et al. [15]. A solid theoretical framework was presented by Fellegi and 99 Sunter [9]. Blocking, which is surveyed by Papadakis et al. [8,9], is considered an important subtask of entity 100 matching. This is meant to tackle the quadratic complexity of potential matches. Christophides et al. [24] 101

specifically review entity matching techniques in the context of big data. Significant research has gone into active 102 learning approaches by Arvind [3], Jungo [11] and Kun [12]. Interestingly, Jungo et al. [11] use a deep neural 103 network in their active learning approach. Such human-in-the-loop factors are often crucial for entity matching 104 in practice as analyzed by Anhai et al. [2]. Many state-of-the-art models for natural language processing are 105 based on deep learning networks. Central to all these approaches is how text is transformed to a numerical 106 format suitable for a neural network. This is mainly done through embeddings, which are translations from 107 text units to a vector spacetraditionally available in a lookup table. The text units will usually be characters 108 or words. An embeddings lookup table may be seen as parameters to the network and can be learned together 109 with the rest of the network endto-end. That way the network is able to learn good, distributed character or 110 word representations for the problem at hand. The words used in a data set are often not unique to that data 111 set, but rather just typical words from some language. Therefore, one may often get a head start by using 112 pretrained word embeddings like word2vec, GloVe or fastText, which have been trained on enormous general 113 corpora. One particular influential recent trend is the ability to leverage huge pretrained models that have 114 been trained unsupervised for language modeling on massive text corpora similar to what the computer vision 115 community has done for image recognition. They produce contextualized word embeddings that consider the 116 surrounding words. These contextual embeddings can be used as a much more powerful variant of the classical 117 118 word embeddings, but as popularized by BERT. However, with neural networks, the actual line between the 119 initial feature extraction part and the rest is an artificial one and not necessarily indicative of how the networks 120 actually learn and work. But they do reflect design decisions to a certain degree and help us compare them in that regard. Often these approaches use pre-built word embeddings for a specific set of values. Our research 121 focuses on entity matching based on attributes where the number of attributes may vary from one use case to 122 another. Also, we try to address the problem of multiple domains, i.e., the machine learning model must be 123 suitable for entities from various categories like customers, products, vendors, etc. In this paper, we present a 124 machine learning model which will perform attribute-based matching of entities. The type, number of attributes 125 may vary over the time, but our approach does not require re-design. Merely a re-training of the model on the 126 new data set will suffice. The model is robust enough to handle slight variations in ordinality and type of the 127 attributes. 128

129 6 III.

130 7 Methodology

Most neural network-based methods perform entity matching by producing so-called knowledge graph embeddings, embeddings of entries which incorporate information about their relationship with other entries. The embeddings work mainly at word level or character level. Embeddings offer neural networks an initial mapping from the actual input to a suitable numeric representation. When we surveyed the earlier methods, we found that researchers focus on explicit levels of representation of entities into single word or text. However, we try to address two problems mainly,

? How to perform matching of entities containing attributes of different data types, say string, boolean, and 137 categorical? ? Will the machine learning algorithm continue to work even if the number of attributes change over 138 the time? Let's say there are few entities in a data set as shown in Table 1. It has two duplicates. Following is a 139 generalized notation. The entities e1 and e2 are same, though they might vary slightly in their attribute values 140 but have similar meanings. Our aim is to design an approach which will combine the attribute level similarity and 141 artificial intelligence to classify entities as unique or duplicate. We propose a two-step methodology where the 142 first step involves calculating attribute level similarity scores and the second step is classification using supervised 143 learning. Feature extraction involves use of a distance function for every pair of attributes. It transforms every 144 pair of entities into numerical vector. For any give pair of attributes (??????, ??????), the distance function 145 146

147 If the two attributes are exactly same, then the distance metric is zero. If they are completely unrelated, 148 then the distance is 1. Partial match will result in value between 0 and 1. We call it as similarity score of the 149 attributes.

A sample set of vectors for a set of three entities will be as shown below. The extracted values correspond to two class labels duplicate (D) and unique (U). If we extract feature vectors of a data set and plot the points in a 3-dimensional space, then we will see two clusters as shown below. Our approach takes every pair of entities and produces a numerical vector. This is in turn fed to a machine learning algorithm for classification. We use supervised learning algorithm for classification. The ML model learns from the training data set and makes accurate predictions on the incoming test data.

¹⁵⁶ 8 a) Feature Extraction using Similarity Score

The first step in ML modeling is data preprocessing, which is usually a crucial step in many data analytics tasks. Typical transformations involve lowercasing all letters, removing excess punctuation, normalizing values, and tokenizing. There are two other major steps in our process. Second one being the feature vector construction using similarity score and For example, consider a similarity function Levenshtein Distance. The Levenshtein distance between 'new yrk' and 'new york' is one since it needs at least one edit (insertion, deletion, or substitution) to transform from 'new yrk' to 'new york'. It is advisable to normalize the similarity scores between 0 and 1 for improved accuracy of the machine learning algorithm.

¹⁶⁴ 9 b) Classification using Supervised Learning

The matching phase aims to develop the prediction model, which takes a candidate pair as input and predicts whether they are matching or nonmatching. Figure 2 illustrates that the model predicts an output label Duplicate (D) or Unique (U). This is a binary classification problem. Data scientists need to decide which algorithm is most suitable for their classification task. Based on our study and experiments, we found three classification algorithms suitable for this task.

¹⁷⁰ 10 i. K-Nearest Neighbors (KNN)

K-Nearest Neighbors (KNN) is an algorithm that learns all available cases from data set and classifies new data item by a majority vote of its K neighbors. A case assigned to the data is majority of its K nearest neighbors measured by a distance (metric) function. The metric functions include Euclidean, Manhattan, Minkowski, and Hamming distances. KNN can be used for both regression and classification problems. However, it is widely used in classification problems in the industry.

ii. XG Boost XG Boost stands for Extreme Gradient Boosting. It is a scalable, distributed gradient-boosted
 decision tree (GBDT) machine learning library. It provides parallel tree boosting and is the leading machine
 learning library for regression, and classification problems.

¹⁷⁹ 11 iii. Support Vector Machines (SVM)

Support Vector Machine is a supervised algorithm in which the learning algorithm analyzes data and recognizes patterns. We plot the data as points in an n-dimensional space. The value of each feature is then tied to a particular co-ordinate, making it easy to classify the data.

183 And finally, we need to tune hyper-parameters in order to get the best model performance.

184 IV.

185 12 Experiments and Results

Automatic entity matching makes the life of commercial organizations easier. A company that maintains thousands of customer records cannot afford to employ many people to verify manually and identify duplicates.

Artificial Intelligence based entity matching is an efficient and cost-effective analytics tool for operational efficiency. We used open-source data sets for our experiments. While several open-source datasets are available,

we picked up few commercial data sets for analysis. In this section, we describe the evaluation tasks, the data

191 sets used, and the experimental results of our approach.

Evaluation Tasks: 1. We evaluate our approach on real-world data set. 2. We evaluate our approach on popular benchmarks.

Our goal is to provide real-life solution using our approach. We aim to evaluate the quality of entity matching. The empirical result is compared with realtime data to harness the accuracy. The results show promising output.

196 13 a) Data Set

We conducted extensive experiments on realworld benchmark entity datasets to evaluate the performance of approach. Following are few opensource data sets available for evaluating entity matching algorithms. Many commercial organizations are nowadays struggling with customer de-duplication. Automatic deduplication has significance in various sectors like Banking and Finance, Insurance, Telecom, Retail, etc. Hence our results mainly focus on the evaluation metrics accuracy on the customer data set.

²⁰² 14 b) Popular Metrics

In this section, we first describe a set of metrics commonly used for evaluating the performance of our classification
 model. Then we present a quantitative analysis of the performance using popular benchmarks.

Accuracy and Error Rate: These are primary metrics to evaluate the quality of a classification model. Let TP, FP, TN, FN denote true positive, false positive, true negative, and false negative, respectively. The classification Accuracy and Error Rate are defined in Equation 1.

where ?? is the total number of samples. Obviously, we have Error Rate = 1 -Accuracy.

210 Precision, Recall, and F1 Score: These are also primary metrics and are more often used than accuracy or

- error rate for imbalanced test sets. Precision and recall for binary classification are defined in Equation ??. The
- F1 score is the harmonic mean of the precision and recall, as in Equation ??. F1 score reaches its best value at
- 213 ~ 1 (perfect precision and recall) and worst at 0.

^{208 ()1}

²¹⁴ 15 c) Empirical Results

We aim to use our entity matching systemin real world applications like retail, e-commerce etc. We analyzed the presence of duplicate customer data and results showed more that 80% accuracy in read-world data sets. Following is a set of predictions made by our system from Fodors-Zagats dataset. From the above table, we observe that customers, vendors can easily get their ambiguities resolved using automatic entity matching system. AIbased entity matching is an alternative to traditional manual or other text analysis-based tools, and it is costeffective solution for decision-makers.

221 V.

222 16 Conclusion and Future Work

The proposed method accomplished superior performance in terms of time and cost. The overall benefits of AI-based entity matching include: Sorting Data at Scale: Manually screening thousands of customer records, or product details is complex and time-consuming. AI-based entity matching helps businesses process large amount of data in an efficient and cost-effective way.

227 17 Real-Time Analysis:

The automatic entity matching can help organizations quickly identify duplicates on realtime basis and act swiftly before duplicate marketing or promotional offers are sent out. Though many deep learning models are being developed nowadays for entity matching, we propose a supervised learning model for few major reasons.

Explainability and Ease of Debugging: For many applications, it is crucial to trust the data source, and try 231 to understand why something does not work is key. Unfortunately, deep learning models are notoriously hard to 232 interpret. As steps in the entity matching process increasingly coalesce into a large neural network, we get fewer 233 checkpoints along the way in the process that can easily be inspected. We can't see the output from each step 234 in the same way anymore. Therefore, figuring out why two records where matched or not matched is usually 235 nontrivial while inspecting deep learning models. There are a few techniques that are already used, such as 236 looking at alignment scores, but we are still far away from a comprehensive way of debugging neural networks 237 for entity matching. Our model addresses the challenges of explainability, running time in interactive settings, 238 and the large need for training examples. Explainability of our supervised learning algorithm helps researchers to 239 240 improve accuracy through inspection, comparison of algorithms and meet the real-world demands. We also see 241 a lot of opportunities in trying to develop more open datasets, standardized benchmarks, and publicly available pretrained models for entity matching.



 $\mathbf{1}$

Figure 1: Figure 1:



Figure 2: Figure 2:

Accuracy =
$$\frac{(TP + TN)}{N}$$
, Error rate = $\frac{(FP + FN)}{N}$

Figure 3: (2)

242

 $^{^{1}}$ © 2023 Global Journals

Pre	$cision = \frac{TP}{TP + FP}$, Recall =	TP TP + FN	, F1-se	$\operatorname{core} = \frac{2 * P}{Pre}$	Prec * Rec ec + Rec
		Figure 4:	Figure 3	:		
No. 1 2 3	Name Alexander Great Alexander G Alexander Graham	Address 2/13, Philip Fra 2/13, Philip Stru 10, Middle Stree	nce Stree eet, Paris et, New Y	t, Paris, fork	Email alex.gr© n/a alex.gr©)gmail.com)yahoo.com
		Figure 5	Table 1	:		
Entity	y Attribute1 Attribu	te2 Attribute3			Label	
el		all	a12	a13	Duplicate	
e2 e3		a21 a31	a22 a32	a23 a33	(e1 = e2) Unique	
		Figure 6	Table 2	:		
Entity e1,e2	y Pair	$\frac{\text{Score1}}{??(??}$		Score2	Score3	Label
Note: 1. ?(?? 22 ?(?? 13	$\begin{array}{l}1 \ , \ ?? \ 21 \) = 0.8 \ ??(?? \\2 \ , \ ?? \ 32 \) = 0.6 \ ??(?? \\3 \ . \ ?? \ 33 \) = 1 \ U \end{array}$	$\begin{array}{l} 12 \;,\; \ref{eq: 12 } : 22 \;) = 0.6 \; \ref{eq: 12 } : 23 \;,\; \ref{eq: 12 } : 23 \;) = 0 \; U \; e \end{array}$?(?? 13 , ? 1,e3 ??(??	? 23) = 11 11, ?? 31	$D \ e2, e3 \ ??(?? \ 21) = 0.6 \ ??(?? \ 12)$, ?? 31) = , ?? 32) =

Figure 7: Table 3 :

 $\mathbf{4}$

		popular
		algorithms.
No.	Data Type	Similarity Function
1		Exact Match
2		Levenshtein Distance
3	Single Word String	Jaro Distance

[Note: last step is machine learning. One might also view second as feature extraction, since records are transformed to a feature space. The success of this entity matching systems depends upon careful selection of right algorithms. Attributes are often assumed to be strings, but that is not the case always. Attributes of an entity may be of any data type like string, numeric, categorical, boolean etc. One single function will not be able to calculate similarity score for various attributes to attribute. It is useful to compare various functions available for similarity score and pick the right choice. To this end, we present a high-level overview of few]

Figure 8: Table 4 :

5					
No.Dataset		Description	Training esting No.		
			Size	Size	of At- tributes
1	Fodors- Zagats	Customer records with name, address, city, phone, type, and category code.	757	189	6
2	iTunes- Amazon	Records of songs with song name, artist name, album name, genre, etc.	430	109	8
3	DBLP- ACM	Publication dataset with paper title, author, venue etc.	9890	2473	4
4	DBLP- Scholar	Publication dataset with title, authors, venue, and year.	22965	5742	4
5	Amazon- Google	Software product dataset with attributes product title, manufacturer, and price.	9167	2293	3
6	Walmart- Amazon	Electronic product dataset with attributes product name, category, brand, model number, etc.	8193	2049	5
7	Abt-Buy	Product dataset with attributes product name, price, and description.	7659	1916	3

Figure 9: Table 5 :

6

CLASSIFICATION ALGORITHM XGBoost

Figure 10: Table 6 :

KNN

SVM

$\mathbf{7}$

No.	Name	Address	City	Phone	Label
1	restaurant ritz-carlton	181 Peachtree st.	Atlanta	404/659	
	atlanta			-0400	
2	ritz-carlton restaurant	181 Peachtree st.	Atlanta	404/659	D
				-0400	
3	posterior	545 post st.	San Francisco $415/$	776 -7825	
4	postrio	545 post street.	San Francisco $415/$	776 -7825	D
5	tavern on the green	in central park at 67 th st	New York	212/873	
				-3200	
6	tavern on the green	central park west	New York	212/873	D
				-3200	
7	carey's	1021 cobb pkwy . se	marietta	770-422-	U
				8042	

Figure 11: Table 7 :

- [Christen (2012)] 'A Survey of Indexing Techniques for Scalable Record Linkage and Deduplication'. P Christen
 IEEE Trans. Knowl. Data Eng 2012. Sept. 2012. 24 p. .
- [Ivan et al. (1969)] 'A Theory for Record Linkage'. P Ivan , Alan B Fellegi , Sunter . Jr. Am. Stat. Assoc 1969.
 Dec. 1969. 64 p. .
- [Kun Qian et al. ()] 'Active Learning for Large-Scale Entity Resolution'. Lucian Kun Qian , Prithviraj Popa ,
 Sen . Proceedings of the 2017 ACM on Conference on Information and Knowledge Management (CIKM '17),
 (the 2017 ACM on Conference on Information and Knowledge Management (CIKM '17)) 2017.
- [Naumann and Herschel ()] An Introduction to Duplicate Detection, Felix Naumann, Melanie Herschel. 2010.
 Morgan and Claypool Publishers.
- [Newcombe et al. (1959)] 'Automatic linkage of vital records'. H B Newcombe , J M Kennedy , S J Axford , A
 P James . Science 1959. Oct. 1959. 130 p. .
- 254 [Papadakis and Svirsky (2016)] 'Comparative Analysis of Approximate Blocking Techniques for Entity Resolu-
- tion'. George Papadakis , Jonathan Svirsky . Proceedings VLDB Endowment May 2016. 9 p. . (Avigdor Gal,
 and Themis Palpanas)
- [Christen ()] Data Matching: Concepts and Techniques for Record Linkage, Entity Resolution, and Duplicate
 Detection, Peter Christen . 2012. Springer Science & Business Media.
- [Thomas et al. ()] Data Quality and Record Linkage Techniques, N Thomas , Herzog , J Fritz , William E
 Scheuren , Winkler . 2007. Springer Science & Business Media.
- [Li et al. (2020)] Deep Entity Matching with Pre-Trained Language Models, Yuliang Li , Jinfeng Li , Yoshihiko
 Suhara , Anhai Doan , Wang-Chiew Tan . arXiv:cs.DB/2004.00584. 2020. April 2020.
- [Kooli et al. ()] Deep Learning Based Approach for Entity Resolution in Databases. In Intelligent Information and
 Database Systems, Nihel Kooli, Robin Allesiardo, Erwan Pigneul. 2018. Springer International Publishing.
 p. .
- [Muhammad et al. (2018)] 'Distributed representations of tuples for entity resolution'. E Muhammad , T
 Saravanan , Shafiq Joty , Mourad Ouzzani , Nan Tang . Proceedings VLDB Endowment 2018. July 2018. 11
 (11) p. .
- [A K Elmagarmid et al. (2007)] 'Duplicate Record Detection: A Survey'. P G A K Elmagarmid , V S Ipeirotis ,
 Verykios . *IEEE Trans. Knowl. Data Eng* 2007. Jan. 2007. 19 p. .
- [Christophides et al. (2019)] End-to-End Entity Resolution for Big Data: A Survey, Vassilis Christophides ,
 Vasilis Efthymiou , Themis Palpanas , George Papadakis , Kostas Stefanidis . arXiv:cs. 2019. May 2019.
- [Fu et al. ()] 'End-to-end Multi-perspective Matching for Entity Resolution'. Cheng Fu , Xianpei Han , Le Sun , Bo Chen , Wei Zhang , Suhui Wu , Hao Kong . *Proceedings of the 28th International Joint Conference* on Artificial Intelligence (IJCAI'19), (the 28th International Joint Conference on Artificial Intelligence
- 276 (IJCAI'19)Macao, China) 2019. AAAI Press. p. .
- [Brunner and Stockinger (2020)] 'Entity matching with transformer architectures-a step forward in data integration'. Ursin Brunner , Kurt Stockinger . International Conference on Extending Database Technology, (Copenhagen) 2020. 30 March-2 April 2020.
- [John R Talburt ()] Entity Resolution and Information Quality, John R Talburt . 2011. Elsevier.
- [Deepak Gottapu et al. (2016)] 'Entity Resolution Using Convolutional Neural Network'. Ram Deepak Gottapu
 , Cihan Dagli , Bharami Ali . *Procedia Comput. Sci* 2016. Jan. 2016. 95 p. .
- [Getoor and Machanavajjhala (2012)] 'Entity Resolution: Theory, Practice & Open Challenges'. Lise Getoor ,
 Ashwin Machanavajjhala . *Proceedings VLDB Endowment* 2012. Aug. 2012. 5 p. .
- [Köpcke et al. (2010)] 'Evaluation of Entity Resolution Approaches on Real-world Match Problems'. Hanna
 Köpcke , Andreas Thor , Erhard Rahm . Proceedings VLDB Endowment 2010. Sept. 2010. 3 p. .
- 287 [Doan et al. ()] 'Human-in-the-loop challenges for entity matching: A midterm report'. Anhai Doan , Adel
- Ardalan, Jeffrey Ballard, Sanjib Das, Yash Govind, Pradap Konda, Han Li, Sidharth Mudgal, Erik Paulson, G C Paul Suganthan, Haojun Zhang. Proceedings of the 2Nd Workshop on Human-In-the-Loop Data
- Analytics (HILDA'17). dl.acm.org, (the 2Nd Workshop on Human-In-the-Loop Data Analytics (HILDA'17).
 dl.acm.orgNew York, NY, USA) 2017. 12 p. 6.
- [Kasai et al. ()] 'Low-resource Deep Entity Resolution with Transfer and Active Learning'. Jungo Kasai , Sairam
 Kun Qian , Yunyao Gurajada , Lucian Li , Popa . Proceedings of the 57th Annual Meeting of the
 Association for Computational Linguistics, (the 57th Annual Meeting of the Association for Computational Linguistics, the 57th Annual Meeting of the Association for Computational Linguistics.
- [Arasu et al. ()] 'On Active Learning of Record Matching Packages'. Arvind Arasu, Michaela Götz, Raghav
 Kaushik. Proceedings of the 2010 ACM SIGMOD International Conference on Management of Data
 (SIGMOD '10), (the 2010 ACM SIGMOD International Conference on Management of Data (SIGMOD
 (10) New York, NY, USA) 2010, ACM, p.
- ²⁹⁹ '10)New York, NY, USA) 2010. ACM. p. .

17 REAL-TIME ANALYSIS:

Ioannou et al. (2013)] 'On Generating Benchmark Data for Entity Matching'. Ekaterini Ioannou , Nataliya
 Rassadko , Yannis Velegrakis . Jr. Data Semant 2013. March 2013. 2 p. .

302 [Papadakis and Skoutas (2019)] George Papadakis , Dimitrios Skoutas . arXiv:cs.DB/1905.06167. Emmanouil

- Thanos, and Themis Palpanas. A Survey of Blocking and Filtering Techniques for Entity Resolution, May 2019.
- 305 [Arasu et al. ()] 'Transformation-based Framework for Record Matching'. A Arasu , R Chaudhuri , Kaushik .
- 306 IEEE 24th International Conference on Data Engineering.ieeexplore.ieee. org, 2008. 2008. p. .