# Unsupervised Matching of Data and Text

Naser Ahmadi
EURECOM, France
naser.ahmadi@eurecom.fr

Hansjörg Sand
KPMG, Germany
hsand@kpmg.com

Paolo Papotti
EURECOM, France
paolo.papotti@eurecom.fr

*Abstract*—Entity resolution is a widely studied problem with several proposals to match records across relations. Matching textual content is a widespread task in many applications, such as question answering and search. While recent methods achieve promising results for these two tasks, there is no clear solution for the more general problem of matching textual content and structured data. We introduce a framework that supports this new task in an unsupervised setting for any pair of corpora, being relational tables or text documents. Our method builds a fine-grained graph over the content of the corpora and derives word embeddings to represent the objects to match in a low dimensional space. The learned representation enables effective and efficient matching at different granularity, from relational tuples to text sentences and paragraphs. Our flexible framework can exploit pre-trained resources, but it does not depends on their existence and achieves better quality performance in matching content when the vocabulary is domain specific. We also introduce optimizations in the graph creation process with an "expand and compress" approach that first identifies new valid relationships across elements, to improve matching, and then prunes nodes and edges, to reduce the graph size. Experiments on real use cases and public datasets show that our framework produces embeddings that outperform word embeddings and fine-tuned language models both in results' quality and in execution times.

## I. Introduction

In data integration, matching records referring to the same real world object is an important task, usually referred to as *entity resolution* (**ER**) [1], [2], [3]. In other communities, such as in Natural Language Processing (NLP), *text matching* (**TM**) is also a widespread task in many applications, such as question answering [4] and information retrieval [5]. However, in many scenarios the borders between the two tasks are not clearly defined. Several datasets have long textual cell values, such as product descriptions. Text documents have structural properties and content organized in hierarchies. Finally, some applications *match textual content to structured data*, such as relational tuples [6], [7].

The best **ER** results are obtained by methods exploiting deep learning techniques [1], [8], [2], [9], but they rely on the presence of multiple attributes (or fields) in the schema, which is missing in text. Transformer-based approaches have enabled important improvements in **TM** [10], [11]. However, transformers are designed to capture the (hidden) relationships in the language. Methods designed to learn table relationships need a large corpus in the pre-training [12], [13], [14] and their pre-trained models do not achieve top quality performance on unseen tables for our task, as we show experimentally.

Consider the following examples.

**Example 1: Text and relational data.** A corpus of product reviews is gathered from the Web. A company must link tuples in a relation to these reviews for a promotional campaign. However, the product reviews have no identifier (Figure 1).

**Example 2: Structured texts.** An enterprise manual about auditing processes is hard to navigate for the final users. To support search, the paragraphs in the manual must be matched to a large taxonomy of concepts (Figure 2).



Fig. 1: Text and data: paragraph $p1$ matches tuple $t2$.



Fig. 2: Structured texts: $1^{st}$ paragraph matches $4^{th}$ node.

This problem is difficult because matching is based on overlapping (or similar) content in the input objects, but this signal can be missing or ambiguous in this new setting.

**Missing matches.** In Example 1, movie *Pulp Fiction* is reported as *Drama* in the table but *comedy* is mentioned in the review. This problem can be partially tackled with the pre-trained embedding of *Tarantino*, which models him as a director for both comedies and dramas [15]. But in Example 2, modeling the connection between *PDCA* and its full spelling is crucial to match the paragraph to the right taxonomy node. While pre-trained embeddings can be used to identify synonyms for common words and popular entities, they fail for *domain specific* terms. **Challenge 1**: as specific vocabularies are not well modeled by pre-trained resources, we need to learn embeddings *across the heterogeneous corpora* at hand to discover similarities in their content.

**Ambiguous matches.** In Example 1, an actor named *Willis* appears in different paragraphs and tuples, but only one tuple
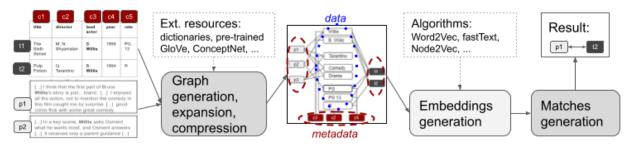
Fig. 3: The proposed unsupervised solution: (1) jointly model text and structured data documents in a graph, (ii) produce embeddings for data and metadata nodes (representing texts, taxonomy nodes, tuples), (iii) match metadata nodes.

in the relation is the correct match in this case. Similarly, the term *audit* appears in most paragraphs and taxonomy nodes in the second example. This suggests the need of a weighting mechanism for combining the matching tokens across two candidate objects. **Challenge 2**: there is the need to learn how to combine matching signals, but the lack of training data rules out solutions based on fine tuning of pre-trained models.

Previous methods lack the flexibility to cover such challenges. Existing approaches handle traditional **ER** and **TM**, but they fail in terms of accuracy in the new use cases above. What is missing is a unified representation that is at the same time modeling the *relationships in the structured content* (for learning a good representation and identify similarities) and the *importance of every matching word* when comparing heterogeneous objects. The last point highlights the need of learning a comparable representation for sets and sequences of tokens, such as tuples and text documents.

To overcome these issues, we propose a framework for learning representations of data and text that (i) is tailored at the domain at hand with a joint modeling of heterogeneous corpora and (ii) exploits structured information whether available to improve the quality of the generated embeddings and of the matching process.

Figure 3 shows our solution. First, it represents text documents and tables as nodes and edges in an undirected graph. This graph contains two main types of nodes. *Data* nodes represent tokens (words) in the corpora, either in text paragraphs or in table cells. *Metadata* nodes represent IDs for tuples, attributes and paragraphs. Graph edges represent the relationship between data and metadata, e.g., a tuple/attribute/paragraph contains the token in a data node. As our goal is to match metadata nodes, we aim at creating more paths between related nodes and at removing spurious connections. The first goal is achieved in an expansion step that exploits external resources, such as ConceptNet [16]. The second goal is obtained by pruning edges and nodes with graph compression techniques designed for our matching task. Next, we generate an embedding for every graph node. We rely on existing solutions for this step and the algorithm at hand can be replaced as the community makes progress in this task. Finally, we use the embeddings for the metadata nodes in an unsupervised algorithm to identify the matching ones, such as the paragraph and the tuple in the first example.

The framework enables users to improve the solution according to the requirements and the resources at hand. If relevant external resources exist, such as word dictionaries, they can be plugged in graph construction to merge data nodes. Knowledge graphs and ontologies are plugged in the expansion step to find more relationships across metadata nodes. New embedding generation algorithms can be plugged in the second step to improve the quality of the embeddings.

Our work extend recent solutions designed for the **ER** task (data to data only) to this new problem [1], [17], [2] by enabling matching for schema-less text paragraphs and taxonomies. The module for graph generation is our first contribution as it creates a rich representation which is then reflected in the *domain specific* embeddings for text and structured data metadata nodes, thus tackling the first challenge. Our second contribution is the expansion and compression approach that, together with matching generation module, exploits the benefits of embeddings metadata nodes in an *unsupervised solution*, thus tackling the second challenge. While our work extends [1], this paper introduces the novel problem of matching text paragraphs to tuples and nodes in taxonomies, it extends the original graph to data sources without a schema, and it handles long free texts with a novel expansion and compression approach. Our proposal outperforms state-of-the-art methods by increasing quality performance up to 45% in absolute terms while taking a fraction of their time in matching. Finally, while we focus on unsupervised applications, any downstream classifier can be trained using the embeddings from our solution.

In the following, Section II describes the proposed framework and how we generate and refine the graph at the core of the proposal. Section III introduces algorithms to expand the original graph with external resources and compress it to keep its size manageable. Section IV introduces the methods for producing embeddings and matching metadata nodes. Section V evaluates our work with datasets from real applications. Section VI discusses related work. Section VII concludes the paper with open challenges and the next steps we plan to take.

## II. A GRAPH FOR HETEROGENEOUS CORPORA

In this section, we describe the algorithm for the generation of a graph across heterogeneous corpora. We discuss the case with text documents and tables as corpora, but the same algorithm applies for the case with documents (tables) only.

The input of the graph creation are two corpora. A *corpus*, based on the task, is a table, some structured text, or simple text. The *document* to match is a tuple, for tables, while the granularity of the text is user-defined and can span from a single sentence to a paragraph. These corpora are matched in three possible combinations, or *tasks*: text to structured text matching, text to data matching, or text to text matching. The purpose of each task is to find the top-$k$ closest documents in the second corpus for all the documents in the first corpus. External resources, such as pre-trained embeddings, are not needed to run the pipeline, but they can be naturally exploited as we discuss in Section II-C.

We jointly represent the document corpora and tables in a graph, from which we then generate embeddings. We first perform some pre-processing steps for every corpus. This includes stop-words removal and stemming on the tokens coming from the texts and the cell values of the tables. We call *terms* these processed values and a term can be composed of one or multiple tokens. For example, "The Sixth Sense" is a term composed of three tokens.

We define two types of nodes. *Data nodes* represent the terms after pre-processing. If a term is contained in multiple documents across the corpora, it still appears as a single node in the graph. *Metadata nodes* represent a group of tokens, such as a sentence, a tuple, or an attribute. Undirected and unweighted edges connect metadata with the respective data nodes, i.e., a tuple node is connected to its tokens.

The graph creation is presented in Algorithm 1. The algorithm takes as input two corpora of pre-processed documents. It creates a metadata node for each document in the first corpus (lines 3-4). For example, Figure 4 shows metadata nodes *t1* and *t2* for the tuples. If the document is a table, it also creates a metadata node for every attribute (lines 5-10), such as nodes *c2*, *c3*, *c4*, and *c5*. These columns are modeled as nodes as they add 2-hop paths across values from the active domain of an attribute. If the document is a structured text, its nodes are modeled as metadata nodes and edges are added to represent relations (lines 12-15). For each term associated to the document (metadata) node, the algorithm then creates term nodes (lines 18-20) and connects each data node to its respective metadata node (lines 21-24). For example, edges are created to connect *t1* to *Shyamalan*, *Willis*, *B._Willis*, *PG*, and *Thriller*. Next, metadata nodes for the documents in the second corpus are created (lines 27-28), thus adding *p1* and *p2* in Figure 4. Term nodes for these documents are connected to their metadata nodes (lines 29-34). For example, metadata node *p1* is connected to data nodes *Willis* and *Comedy*.

Algorithm 1 creates different metadata nodes based on the input documents. In this example, it outputs metadata nodes to represent tuple, columns, and text. For the other two tasks (text to text and text to structured text), only text metadata nodes are produced.

### A. Connecting metadata nodes

The graph creation algorithm never connects metadata nodes from different corpora, as we assume that these connections

---

**Algorithm 1:** Graph Creation

1 **Input.** Two sets of documents;
2 G = An un-directed graph;
3 **foreach** *document doc_i in the first set* **do**
4     G.addNode($doc\_i$) ;
5     **if** $get\_type(first\ set)$=*table* **then**
6         **foreach** *column col_j in the document* **do**
7             **if** $\neg$ *G.hasNode(col_j)* **then**
8                 G.addNode($col\_j$)
9             **end**
10         **end**
11     **end**
12     **if** $get\_type(first\ set)$=*structured* **then**
13         parent ← get parent of document i ;
14         **if** *G.hasNode(parent)* **then**
15             G.addEdge($doc\_i$,parent)
16         **end**
17     **end**
18     terms ← get list of terms in document i ;
19     **foreach** *term tm_k in terms* **do**
20         G.addNode($tm\_k$)
21         G.addEdge($doc\_i$,$tm\_k$)
22         **if** $get\_type(set)$=*table* **then**
23             G.addEdge($col\_j$,$tm\_k$)
24         **end**
25     **end**
26 **end**
27 **foreach** *document doc_i in the second set* **do**
28     G.addNode($doc\_i$);
29     terms ← get list of terms in the document;
30     **foreach** *term tm_j in terms* **do**
31         **if** *G.hasNode(tm_j)* **then**
32             G.addEdge($doc\_i$,$tm\_j$)
33         **end**
34     **end**
35 **end**
36 **Output.** Graph G

---

are hard to infer and are indeed the results of the downstream task in our system. However, those relations can be provided by the users and consumed by our algorithm. On the other hand, metadata text nodes from the same structured document can be connected. For example, for the taxonomy in Figure 2, the corresponding graph connects metadata text nodes for *Audit programme* and *ISO 19001*. This edge represents the hierarchical relation between the concepts.

### B. Filtering nodes

The graph can become extremely large with real text corpora. This is a problem in terms of performance both for the execution time and for the quality, as it may lead to model several terms that do not contribute to the final matching tasks. To address this problem, we filter out irrelevant terms in the graph creation. Algorithm 1 does not create data nodes for
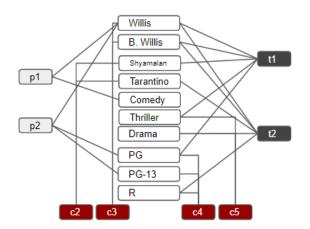
Fig. 4: Graph with a sample of the nodes for Example 1.

all terms in both corpora. It starts by creating data nodes for documents in the corpus with the smaller number of distinct tokens and filters out from the second corpus the terms that are not already in the graph.

As our goal is to model the connections across the two corpora, we compromise the loss of some words (and possibly relationships) in the second corpus to focus the learning in the next step on the terms that create bridges between metadata nodes. To limit the loss of possibly relevant words and to create a more compact graph, we present next some techniques to merge token nodes across corpora.

### C. Merging nodes

Intuitively, (correctly) merging data nodes increases the connectivity between related metadata nodes across corpora and ultimately improves the matching tasks. For example, merging data nodes *Bruce Willis* and *B. Willis* in Figure 1 decreases the distance in many paths connecting metadata nodes *P1* and *t1*. Merging data nodes is easier that solving the metadata matching task and there are several resources available for this operation. For this matter, we use different techniques to merge data nodes:

- *Stemming* merges different forms of a word. For example, in Figure 2, stemming merges *planning* from the first paragraph with node *Plan* from *Plan Do Check Act Steps*.
- *Bucketing* merges data nodes with numeric values by using equal width binning and the *Freedman–Diaconis* rule [18] to compute buckets' width.
- For merging synonyms, acronyms, and typos we use external resources, such as WordNet [19].

For the last case, we merge two data nodes if the cosine similarity between their embedding vectors is higher than a threshold $\gamma$. For calculating $\gamma$, we use a list of 17K synonym terms from *WordNet* and define $\gamma$ as the average cosine similarity between their vectors in the pre-trained model that we use for merging. Specifically, for Wikipedia2Vec [20] we identify and set $\gamma = 0.57$. This approach is widely used in tasks such as entity linking [21] and retrieval [22].

### D. Tokens and terms

One important aspect in graph creation is handling multi-tokens data nodes. There are a lot of meaningful multi-tokens words in document (e.g., movie names) and information is lost if they are split over different single-word data nodes in the graph, e.g., *The Sixth Sense* split over data nodes *Sixth* and *Sense* with *The* filtered out as stop word. A possible solution for tables is to represent the whole cell value as a single data node (*The_Sixth_Sense*). But this granularity has also drawbacks as it may lead to graphs that miss important connections across corpora. For example, if *B. Willis* in the review and node *Willis* are not merged, a strong connection between the correct metadata nodes is lost.

We use a combination of the two approaches that solves both problems. For each text in a corpus, we generate possible $n$-gram tokens for $n = 1, \ldots, n$. For example, for $n = 3$, the graph represents *The Sixth Sense* using five data nodes: *Six*, *Sense*, *The_Six*, *Six_Sense*, and *The_Six_Sense*. This increases the chance of connecting terms of the second corpus with nodes from the first corpus. We identified the default value of $n$ for every scenario by profiling a file of titles of *Wikipedia* articles. About $99\%$ of the titles have at most three tokens. Experiments in Section V-F1 show that increasing $n$ up to three improves the quality performance but we observe diminishing return with higher values. Other results also show that by increasing $n$ in character tokenization, the lexicon size grows rapidly and precision diminishes for most languages [23].

### III. GRAPH EXPANSION AND COMPRESSION

By generating the graph with Algorithm 1, the paths between metadata nodes represent the relationships which are present in the documents and tables. This network of connections leads to embeddings that ultimately guide the metadata matching process. However, the data relationships are not all the existing relations between two real objects represented in the graph. Real entities and concepts are connected by more relationships that are missing from the corpus at hand.

For example, an actor and a director may have worked together in a movie that is not in the movie table or in any of the reviews. Such external relationships can be very valuable if represented in our graph as they lead to better embeddings and ultimately enable better matching. However, while expanding the graph is a valid solution to include external information, we should be careful in trying to remove useless or even misleading new nodes and edges in order to keep the graph as little as possible in terms of size. We present solutions to address these two tasks in this section.

### A. Expanding the graph with external information

A natural approach to expand the generated graph is to employ external resources that model information with nodes and edges, such as ontologies and knowledge bases. By exploiting existing resources, we add information to the graph.

For example, in Figure 4, *p1* is the review related to tuple *t2*. Even though there are seven paths between these two metadata, only one of them has three or less nodes: p1 → Willis →
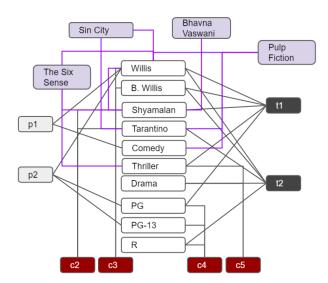
Fig. 5: Expanded graph for Example 1.

t2. By expanding this graph with new nodes and edges, we can add new meaningful paths between these nodes, improve their embeddings, and increase their chance of being matched. Consider as a resource of external information the knowledge graph DBpedia [24]. Among the relations for entity *Tarantino* in DBpedia, there is the following triple: *style(Tarantino, Comedy)*. Adding this new edge to the graph creates nine new path between *p1* and *t2* including one with less than three nodes: p1 → Comedy → Tarantino → t2.

Different external resources can be exploited to expand a graph. In graphs with named entities, we are interested in finding more information about them, such as data about their spouse, country, university, workplace, etc. This information can be extracted from existing entity-centric knowledge bases (KBs) such as *DBpedia* and *Wikidata*. For example, the graph presented in Figure 4 contains information about movies and their casts, those are entities for which we can use a KB for expansion. By expanding this graph with DBpedia, we enrich it with new edges such as *starringOf(Willis, Pulp Fiction)* and *spouse(Shyamalan, Bhavna Vaswani)* as shown in Figure 5.

Textual corpora do not contain only named entities, but also concepts, generic nouns and verbs. For example, by expanding the word *management* in Figure 2, we connect it with the relevant words in the correct text paragraph, such as *planning*. In these cases, other external resources can be exploited, such as *ConceptNet* [25] and *Wordnet*.

Algorithm 2 shows how we exploit any external resource to fetch all connections for every data node in the graph. We also remove any sink node, i.e., nodes that are not connected to more than one other node. E.g., in Figure 5 node *Bhavna Vaswani* is only connected to *Shyamalan* and can be removed.

The expansion technique introduces new graph paths that affect the matching process between metadata nodes. For example, in the original graph in Figure 4, there is one path with less than five nodes between metadata nodes *p1* and *t2*, but the their shortest path is of size two after expansion.

---

**Algorithm 2:** Graph Expansion Algorithm

**Input:** Un-directed graph G

1    External resource E

    // Expanding by fetching connections from E
2  **foreach** *node in G* **do**
3     **if** *node is not a metadata node* **then**
4        relations ← all connections of node in *E*;
5        **foreach** *(node,m) in relations* **do**
6           **if** ¬ *G.hasNode(m)* **then**
7              G.addNode(m)
8           **end**
9           G.addEdge(node,m)
10       **end**
11    **end**
12 **end**
    // Cleaning the graph
13 **foreach** *node in G* **do**
14    **if** *degree(node) == 1* **then**
15       G.removeNode(node)
16    **end**
17 **end**
18 **Output.** Expanded graph G

---

### B. Pruning nodes and edges for compression

Expansion introduces relevant connections between related nodes, but it also increases its size by adding new nodes and edges that are not helpful for our tasks. For example, there are more than 800 relations for entity *Quentin Tarantino* in DBpedia but only a few of them increase the chance of matching *p1* and *t2* (e.g., *directorOf(Quentin Tarantino, Pulp Fiction)*, *redirectsOf(Quentin Tarantino, Samuel Jackson)*). As new nodes and relationships increase the execution time for random walks and embedding generation, we should avoid keeping nodes and edges that do not contribute to the connections among metadata nodes.

We therefore introduce a graph compression techniques to reduce the size of our graph after the expansion phase. Compression for static graphs has been studied for its benefits in terms of reduction of data volume and storage, which in turn enables speedup of algorithms and queries [26]. Noise elimination has also been reported as an important effect of the compression, with the removal of erroneous nodes and labels [27]. Compression methods can be based on node sampling [26], [28], [29], edge sampling [30], [31], or exploration based sampling [32], [33], [34]. Most methods are configurable w.r.t. the desired compression ratio, i.e., the desired size of the output graph compared to the input graph.

As these methods are very general, they are not application specific and cannot make use of the node types in our graph. Our experiments in Section V show that these techniques can help in reducing the size of the graph by filtering nodes, but do not preserve good performance in the matching task.

Our key observation is that the goal of our graph and our embeddings it to match metadata nodes. A crucial component

in determining the distance between the embeddings for two metadata nodes in their distance in the graph. We therefore start the design of our compression algorithm from the idea that it should preserve the shortest path across all the metadata nodes in the two corpora. This is a quadratic number of paths w.r.t. the number of metadata nodes, but this is unavoidable as we do not know at compression time what are the metadata to match.

Inspired by an existing graph compression technique that exploits shortest paths, namely *SSP* [33], we introduce an algorithm tailored at our graph and matching application. The original *SSP* is an exploration based sampling method which takes a sampling size as input. It randomly picks a pair of nodes in each iteration, computes their shortest path, and adds nodes and edges of the shortest path to the output graph. In our setting, the idea is to use metadata nodes in distinct corpora for the selection of node pairs. This guarantees that metadata nodes are connected and keeps in the graph the data nodes that are modeling their relationship concisely.

---

**Algorithm 3:** Graph Compression MSP

**Input:** Un-directed graph G
1        Compression ratio $\beta$
2 CG = empty un-directed graph
3 i = 0
4 L = $\beta$ * size(G.nodes())
5 **while** $i < L$ **do**
     // Select two random metadata nodes
6     first $\leftarrow$ a random node from the first corpus
7     second $\leftarrow$ a random node from the second corpus
8     shortest_paths $\leftarrow$ find all shortest paths between **first** and **second** in G
     // Add nodes and edges of the paths to CG
9     **foreach** *path in shortest_paths* **do**
10        | CG.add(path)
11     **end**
12     i += 1
13 **end**
14 **Output.** Compressed graph CG

---

Algorithm 3 shows our graph compression based on the idea of using Metadata Shortest Path (*MSP*). It takes an un-directed graph and a compression ratio $\beta$ as input and returns a compressed graph. We define the number of iterations of the Algorithm by multiplying $\beta$ and the number of nodes in the graph. We also make sure that all metadata nodes, even if not sampled at lines 6-7, are connected to the graph with at least one shortest path.

There is however an orthogonal challenge in such an aggressive compression. Keeping only shortest paths among metadata nodes may lead to a graph with similar embeddings for all nodes. Assume that all shortest paths are of length three and all of them pivot on a single, very popular data node. This extreme situation leads to all metadata nodes ending up with the same embedding, therefore making the matching process

impossible. This observation highlights that compression is not suitable in all cases. We discuss in the experiments how we can recognize setting that are less likely to benefit from the compression.

## IV. MATCHING TEXT AND STRUCTURED DATA

In this section we discuss how to generate and compare node embeddings for the unsupervised matching of objects.

### A. Embeddings Generation

We generate embeddings for the graph nodes. Multiple methods can be employed to generate embeddings directly from the graph [35], [36]. A less resource intensive solution is to use word embedding models on walks over the graph [37]. In our default setting, we use the second approach as we found the results with different methods comparable in quality, but the latter is faster and less demanding in terms of resources.

---

**Algorithm 4:** Embedding Generation

**Input:** Graph G
1        Number of RWs $n$
2        Length of RWs $l$
3 docs = []
4 **for** *i in range(0,n)* **do**
5     **for** *node in G* **do**
6        rw = []
7        **while** *len(rw) < l* **do**
8           next $\leftarrow$ a random neighbor of node
9           rw.append(next.label())
10        **end**
11        docs.append(' '.join(rw))
12     **end**
13 **end**
14 M $\leftarrow$ Word embedding model
15 M.train(docs)
16 **Output.** Embeddings for data and metadata nodes

---

Algorithm 4 shows the process of generating embeddings. The output of this algorithm are word embeddings that map words in the random walks (data and metadata nodes) to vectors of real numbers. In a word embedding model, words with more co-occurrences in sentences show smaller distances in the corresponding vectors. In our model, related metadata nodes have a higher chance of appearing together in a random walk, thus their vectors show smaller distances. In our default method, a random walk starts from every graph node and at each step it randomly chooses the next node among the current node's neighbors. A sentence is derived with the concatenation of nodes traversed by the walk. The union of the sentences is then processed with *Word2Vec* or similar methods. Multiple parameters of the random walks affect the quality of the embeddings, including whether random walk should be limited to some nodes or not, the length of a random walk, and how many random walks should be generated for every graph node. In Section V-F1, we report the impact of these parameters on the quality of the proposed model.

We stress the importance of the unified graph and of the walks in our context. While any data structure can be serialized as a sequence of sentences (e.g., row by row), this is not effective in practice for learning embeddings for two reasons. First, the resulting sentences do not follow meaningful patterns as in real languages; the relationships that are nicely captured in real text are missing. This is especially true for transformer based solutions based on attention. Second, a simple serialization misses the structural dependencies in a relation; existing features in the data are not exploited. Given text documents and a relation, our graph and walks enable the joint representation for the given corpora.

### B. Matching Metadata Nodes

In the final step, the input of the matching module is a metadata node for a document in the first corpus and the metadata nodes representing all documents in the second corpus. The embedding vectors are used to match such nodes and ultimately the documents they represent (text paragraphs, text sentences, or tuples). The distance between two nodes' vectors is used for matching them. Given the embeddings, we use cosine similarity to identify the top-$k$ neighbours in the second corpus for the metadata node from the first corpus.

Differently from the graph creation, in the metadata matching we found more effective to start the process from the larger corpus and this is our default configuration. However, this decision can be changed according to the specific task. For example, in text matching we do it claim by claim (from the smaller corpus) as this is the natural setting for the application.

## V. EXPERIMENTS

We first introduce our execution setting and the baseline methods. We then report the results in three matching tasks: i) text to data, ii) text to structured text, iii) text to text. We do not report performance for the data to data task since it has already been studied in previous work [1]. Finally, we report execution times and discuss the impact of the different parameters and optimizations in our solution. Code and datasets are available at https://github.com/naserahmadi/TDmatch.

**Execution setting.** Experiments have been conducted on a laptop with CPU Intel i5-7300U, 4x2.6GHz cores and 8GB RAM. For fine tuning the baselines, we used a *Google Colaboratory* instance with CPU Intel 2x2.20GHz, 13GB RAM, and NVIDIA Tesla T4 GPU (16GB memory). Algorithms are written in *Python* with the *Numba* compiler.

**Baselines.** We compare our approach to several baselines; we report only the best performing baselines for every task. Some baselines rely on *training* over the given documents and data (as our approach), while others use *pre-trained* resources. We further distinguish unsupervised and supervised solutions.

For unsupervised methods based on *training*, we test *Word2Vec* (W2VEC) for word embedding and *Doc2Vec* (D2VEC) for document embedding. We obtain embeddings from the documents at hand and then use such embeddings to identify matches. We generate embeddings for longer texts

with the mean of the vectors of their tokens [38]. We use vectors of size 300, Skip-Gram for *Word2Vec* and DBOW for *Doc2Vec*. For unsupervised approaches using *pre-trained* embeddings, we report on *SentenceBERT* (S-BE). Unsupervised methods match objects with the algorithm in Section IV-B.

We denote supervised methods using pre-trained models with * for clarity and always report results for 5-fold cross validation. The first approach is based on fine-tuning for a multi-label classification task on *BERT* large (L-BE*). We also report the results for two supervised state-of-the-art entity matching methods for the text to data task: *Ditto* [2] and *DeepMatcher* [9]. These methods take two tables as input and compute the matching probability for tuples from different tables. We represent text documents as tuples of a table with one attribute. We report the results for text-to-data matching also using *TAPAS* [13], a BERT-based model designed and pre-trained for answering questions about tabular data. We use 60% of the annotated data to train these models. Finally, we report for *Reranking* (RANK*), a supervised algorithm that learns to rank using a pairwise loss [39].

For our unsupervised approach, we used *Word2vec* (W-RW) on the random walks (RW) generated on the graph. In the default configuration, we generate 100 random walks of length 30 for every node. For the text to data task, we use *Skip-gram* with a window of size three as in the data to data match [1], while for text oriented tasks we use *CBOW* with a window of size 15. We report for our method with (W-RW-EX) and without (W-RW) applying the expansion technique in Section III. We use ConceptNet as our default external resource for expanding graphs, except for *IMDB* where we employed DBpedia as this relation contains mostly entities.

### A. Text to Data

For the *text to data* matching we use two datasets. We created a first scenario from the Internet Movie Database (*IMDb*) website with a corpus of movies reviews and a database of movies. We also report results for the *CoronaCheck* scenario, which matches COVID-19 claims to the official datasets [7]. For both scenarios, the task is to find tuples related to each sentence. For example, a sentence "Number of cases in US is higher than China" required to match two rows of a table to verify the claim.

**Datasets.** As the task is novel, we release two new scenarios:

**1. IMDb.** We created the dataset by manually matching two reviews for every movie in "top 1K of all times" to a sample of 50k tuples from the official IMDb dataset. The 2k reviews contain one to 207 sentences, sixteen on average. We created two versions of the target relation: an easier one with 13 attributes, including the title information (**WT**) and a more challenging one without title (**NT**).

**2. CoronaCheck.** This scenario contains a corpus of sentences about COVID-19 spread and effects, such as daily total death cases and new confirmed monthly cases, annotated w.r.t. the corresponding tuples in a dataset with 1.2k tuples about daily cases for all countries. We report for a dataset with 7k

| | Method | MRR | MAP@k | | | HasPositive@k | | |
|---|---|---|---|---|---|---|---|---|
| | | | 1 | 5 | 20 | 1 | 5 | 20 |
| WT | S-BE | .254 | .088 | .142 | .159 | .171 | .339 | .510 |
| | W-RW | .853 | .400 | .678 | .682 | .802 | .919 | .942 |
| | W-RW-EX | **.868** | **.410** | **.691** | **.706** | **.820** | **.926** | **.955** |
| | RANK* | .535 | .218 | .351 | .376 | .438 | .645 | .797 |
| | DITTO* | .759 | .349 | .549 | .553 | .699 | .839 | .877 |
| | TAPAS* | .722 | .375 | .525 | .526 | .802 | .849 | .929 |
| NT | S-BE | .218 | .067 | .118 | .139 | .136 | .301 | .454 |
| | W-RW | .780 | .362 | .574 | .589 | .727 | .841 | .906 |
| | W-RW-EX | **.792** | **.371** | **.587** | **.598** | **.749** | **.854** | **.911** |
| | RANK* | .404 | .156 | .236 | .260 | .312 | .494 | .688 |
| | DITTO* | .560 | .265 | .386 | .410 | .428 | .689 | .814 |
| | TAPAS* | .643 | .327 | .354 | .358 | .678 | .706 | .715 |

TABLE I: Quality of match results for *IMDb* scenario.

| | Method | MRR | MAP@k | | | HasPositive@k | | |
|---|---|---|---|---|---|---|---|---|
| | | | 1 | 5 | 20 | 1 | 5 | 20 |
| Gen | S-BE | .486 | .294 | .463 | .483 | .295 | .752 | .916 |
| | W-RW | .728 | .575 | .718 | .725 | .578 | .945 | .995 |
| | W-RW-EX | **.755** | **.601** | **.746** | **.752** | **.611** | **.959** | **.996** |
| | RANK* | .460 | .287 | .438 | .455 | .289 | .703 | .845 |
| | DEEP-M* | .376 | .347 | .368 | .374 | .349 | .395 | .439 |
| | DITTO* | .160 | .030 | .161 | .203 | .066 | .283 | .518 |
| | TAPAS* | .394 | .325 | .367 | .389 | .330 | .452 | .723 |
| Usr | S-BE | .354 | .177 | .284 | .320 | .200 | .620 | .860 |
| | W-RW | .518 | .296 | .427 | .472 | .306 | .755 | .979 |
| | W-RW-EX | **.538** | **.329** | **.451** | **.496** | **.371** | **.771** | **1** |
| | RANK* | .332 | .137 | .256 | .303 | .160 | .600 | .880 |
| | DEEP-M* | .321 | .200 | .200 | .248 | .280 | .280 | .600 |
| | DITTO* | .153 | .020 | .100 | .123 | .040 | .281 | .407 |
| | TAPAS* | .192 | .053 | .053 | .077 | .053 | .053 | .474 |

TABLE II: Quality of match results for *CoronaCheck* scenario.

sentences created from the data (**Gen**) and a more challenging dataset with 50 sentences submitted by users on the website https://coronacheck.eurecom.fr (**Usr**) [7].

**Evaluation Measures.** Mean Reciprocal Rank (*MRR*) is the average of reciprocal ranks of queries, i.e., the multiplicative inverse of the rank of the first correct answer. Mean Average Precision (MAP) is the mean of the precision scores after each relevant document is retrieved and we report MAP truncated at rank $k$ (*MAP@k*). We also report HasPositive@k for determining whether there is a true positive among the top-$k$ results.

**Matching results.** As the training-based methods (W2VEC, D2VEC) do not take tables as input, we serialize every tuple to a sentence using two special tokens (*[COL]* and *[VAL]* [2], e.g., the first row in Figure 1 starts with "[COL] title [VAL] The Sixth Sense [COL] director [VAL] Shyamalan". We then generate an embedding vector for every resulting sentence and match vectors for tuple and text metadata nodes. As results are poor for these baselines, we do not report them. For the pre-trained models, we report for S-BE, RANK*, DITTO*, DEEP-M*, and TAPAS*.

Table I and Table II show the results on the **IMDb** and **CoronaCheck** scenarios, respectively. Our method outperforms unsupervised S-BE in all scenarios and the techniques in Section III show a positive effect in both datasets. For **IMDb**, we observe an absolute increase of 0.45 for MRR with W-RW

in both datasets and at least 410x relative improvement for Positive@1. In **CoronaCheck**, the increase for **Gen** sentences is an absolute 0.24 for MRR and up to 0.30 for MAP@k and Positive@k. For the **Usr** sentences, increases are up to 0.2 for MRR and MAP and up to 0.18 for Positive@k. Our model clearly outperforms also supervised methods. Results show that pre-trained models fail short in this task and that the joint modeling enabled by our graph is needed to achieve good matches. We do not report *DeepMatcher* on *IMDB* because it failed due to the limited amount of memory in our machine.

### B. Text to Structured Text

In this task, we match taxonomy elements to a text document in a real enterprise scenario from an auditing company.

**Dataset.** This scenario contains 1622 audit text documents (containing one to 17 sentences, three on average) and a taxonomy containing 747 auditing concepts. Each path spans multiple nodes, e.g., $r_1 : a \rightarrow b \rightarrow c \rightarrow d$, where each variable is a concept. Right arrows show the hierarchical relations between terms, e.g., in $r_1$ c is a child of b and b is a child of a. The length of taxonomy paths are between two and five nodes (four on average). The final graph has 5.9k nodes and 164k edges. Text documents are manually matched to concept nodes by domain experts. About $40\%$ of documents are annotated with one concept, $10\%$ are matched to two concepts, and the rest are matched with three to 27 concepts (four on average).

| | Method | Exact Scores | | | Node Scores | | |
|---|---|---|---|---|---|---|---|
| | *Method* | P | R | F | P | R | F |
| K=1 | D2VEC | .254 | .217 | .234 | .554 | .503 | .527 |
| | S-BE | .094 | .071 | .081 | .379 | .358 | .368 |
| | W-RW | .346 | .265 | .300 | .593 | .530 | .560 |
| | W-RW-EX | **.367** | **.282** | **.319** | **.601** | **.545** | **.572** |
| K=1 | RANK* | .162 | .125 | .138 | .425 | .392 | .408 |
| | L-BE* | .381 | .304 | .338 | .626 | .567 | .595 |
| K=3 | D2VEC | .176 | .386 | .242 | .485 | .564 | .521 |
| | S-BE | .065 | .014 | .088 | .362 | .431 | .393 |
| | W-RW | 201 | .434 | .275 | .521 | .652 | .579 |
| | W-RW-EX | **.214** | **.475** | **.295** | **.528** | **.670** | **.594** |
| K=3 | RANK* | .162 | .125 | .138 | .425 | .392 | .408 |
| | L-BE* | .183 | .417 | .254 | .487 | .678 | .566 |
| K=5 | D2VEC | .132 | .470 | .206 | .457 | .679 | .546 |
| | S-BE | .052 | .179 | .080 | .356 | .473 | .406 |
| | W-RW | .145 | .508 | .222 | .478 | .699 | .568 |
| | W-RW-EX | **.151** | **.533** | **.236** | **.485** | .719 | **.580** |
| K=5 | RANK* | .072 | .242 | .110 | .365 | .522 | .429 |
| | L-BE* | .135 | .508 | .213 | .446 | .740 | .556 |
| K=10 | D2VEC | .087 | .587 | .152 | .42 | .758 | .541 |
| | S-BE | .038 | .253 | .066 | .347 | .541 | .423 |
| | W-RW | .092 | .613 | .160 | .437 | .768 | .557 |
| | W-RW-EX | **.094** | **.629** | **.164** | **.438** | **.783** | **.562** |
| K=10 | RANK* | .051 | .324 | .088 | .350 | .592 | .440 |
| | L-BE* | .081 | .584 | .141 | .393 | .797 | .526 |

TABLE III: Exact and Node scores for structured text matches.

**Evaluation Measures.** For this task, we change quality measures as we show results at different granularity. We report

Precision, Recall and F-score for concepts (in the taxonomy) assigned to every document w.r.t. the ground truth. As different taxonomy nodes can contain the same text, we compare the root to node path in the measures. With **Exact** matches, we consider a match in the top-k valid only if it is *equal* to the path in the ground truth. As two paths can be partially overlapping, we consider also partial matches with the **Node** score, which measures the intersection between the matched path(s) and the closest path(s) in the ground truth. For an accurate calculation, we exclude two most general levels of the taxonomy (root and first level under it) in the intersection and denote the new path with $p'$. We then use formula (1) below to calculate the **Node** score for two paths $p_1$ and $p_2$.

$$Node(p_1, p_2) = \frac{intersection((nodes(p'_1), nodes(p'_2))}{maximum((nodes(p'_1), nodes(p'_2))} \quad (1)$$

Consider $r_1$ and $r_2 : a \rightarrow b \rightarrow c$. After excluding the general nodes, we obtain $r_1 : c$ and $r_2 : c \rightarrow d$, thus Node($r_1$,$r_2$) = 0.5.

**Matching results.** Table III reports for both measures the precision, recall and F-score for matching top-$k$ paths to every document for different $k$ values. Results show that the task is very difficult. Indeed, different auditors have different opinions about the right matches for a given taxonomy node and the ground truth is constructed after a discussion to reach consensus. In this hard task, our methods outperform unsupervised methods with a large margin. This scenario contains some domain specific terms that are not covered by pre-trained models as we can observe by D2VEC (trained on the audit data) outperforming unsupervised S-BE. Only for the top-1 case the supervised *BERT_large* shows small margins for both measures. Supervised classifiers are effective for documents matched against one concept but do not have enough training data for the other cases. In Section V-F2, we show how our model combined with S-BE outperform supervised *BERT_large* solution also for k=1.

### C. Text to Text

We evaluate our framework in matching documents between three text corpora. While our solution is tailored towards structured data and text, its results are better than unsupervised state of the art baselines for this task and close to supervised ones. Two datasets come from the task of detecting previously fact-checked claims [39]. Given a check-worthy input claim and a set of verified claims, the goal is to rank the verified claims that help check the input claim it, above other claims. We also test the dataset from the STS (semantic textual similarity) GLUE task [40] as an unsupervised matching task. In this dataset, an original similarity score between text pairs is defined between 0 (completely dissimilar) and 5 (completely equivalent). We consider two snippets a true match when they have in the ground truth a score equal or greater than $k$.

**Datasets.** The *Snopes* dataset contains a set of 1k claims (tweets) and 11k verified claims (facts), while the *Politifact* dataset contains 768 claims (made by politicians) and 16.6k verified claims (facts). Text documents contain from one to

| Method | MRR | MAP@k | | | HasPositive@k | | |
|---|---|---|---|---|---|---|---|
| | | 1 | 5 | 20 | 1 | 5 | 20 |
| S-BE | .395 | .354 | .372 | .382 | .362 | .417 | .496 |
| W-RW | .489 | .346 | .396 | .401 | .409 | .579 | .702 |
| W-RW-EX | **.507** | **.358** | **.406** | **.418** | **.429** | **.600** | **.726** |
| RANK* | .608 | .531 | .588 | .599 | .535 | .688 | .787 |

TABLE IV: Quality of match results for *Politifact* scenario.

| Method | MRR | MAP@k | | | HasPositive@k | | |
|---|---|---|---|---|---|---|---|
| | | 1 | 5 | 20 | 1 | 5 | 20 |
| S-BE | .543 | .457 | .527 | .535 | .457 | .648 | .724 |
| W-RW | .695 | .586 | .688 | .693 | .587 | .820 | .886 |
| W-RW-EX | **.708** | **.613** | **.698** | **.706** | **.614** | **.843** | **.898** |
| RANK* | .788 | .691 | .782 | .784 | .693 | .894 | .925 |

TABLE V: Quality of match results for *Snopes* scenario.

nine sentences in *Snopes* and from one to 11 in *Politifact*. On average they have less than two sentences. We match top−$k$ verified claims (facts) for every claim. The *STS* dataset contains a set of 7k pairs of sentences and we report results for thresholds k=2 (5k pairs) and k=3 (3.7k pairs).

**Evaluation Measures.** We use Mean Reciprocal Rank (*MRR*), Mean Average Precision at k (*MAP@k*), and HasPositive@k.

**Baselines.** We use methods with good results reported for these datasets [39]: unsupervised S-BE and supervised RANK.

**Matching results.** Results in Tables IV, V and VI show that our method is the best unsupervised solution, outperforming S-BE in all measures and scenarios. Table VI shows that as sentences with higher similarities share more details, increasing $k$ improves the performance for all methods. Our approach sits between the best unsupervised baseline and the supervised method. One explanation is that these datasets contain generic textual claims with common terms, which is the best scenario for pre-trained models trained on very large corpora. Also, long natural language English sentences are nicely modelled by the attention mechanism in transformers. As we discuss in Section V-F2, by combining our embeddings with pre-trained language models, we can improve our performance.

### D. Compression Results

We report the performance of the compressing technique introduced in Section III. As a baseline technique, we report also for *SSuM*, a state of the art method that employs node merging and edge sparsifying to generate a super-graph as output [41]. Table VIII compares the performance of compression methods in terms of *size* (number of nodes and edges) of the compressed graph and of *quality* in the matching task (MRR).

For *MSP*, we report results for iterations equal to half (*MSP (0.5)*) and a quarter (*MSP (0.25)*) of the expanded graph's nodes. *SSuM (0.1)* is set with a compression ratio of 0.9 as this is the value generating the best quality results (MRR) in our experiments.

In terms of size reduction, *MSP (0.25)* is the compression method with the best results in four cases and it is second to *SSuM* only for IMDB. However, it shows an higher decrease in match quality results w.r.t. *MSP (0.5)*, which is the

| | Method | MRR | MAP@k | | | HasPositive@k | | |
|---|---|---|---|---|---|---|---|---|
| | | | 1 | 5 | 20 | 1 | 5 | 20 |
| k=2 | S-BE | .739 | .649 | .723 | .733 | .657 | .836 | .920 |
| | W-RW | .780 | .691 | .765 | .772 | .703 | .872 | .947 |
| | W-RW-EX | **.796** | **.707** | **.785** | **.788** | **.716** | **.887** | **.962** |
| | RANK* | .798 | .714 | .789 | .796 | .717 | .899 | .967 |
| k=3 | S-BE | .842 | .767 | .832 | .838 | .773 | .925 | .999 |
| | W-RW | .841 | .766 | .832 | .836 | .775 | .926 | .974 |
| | W-RW-EX | **.858** | **.787** | **.848** | **.854** | **.795** | **.887** | **.988** |
| | RANK* | .890 | .830 | .887 | .890 | .830 | .969 | .994 |

TABLE VI: Quality of match results for *STS* scenario.

| *Method* | Text to data | | Structured text | | Text to text | |
|---|---|---|---|---|---|---|
| | *Train* | *Test* | *Train* | *Test* | *Train* | *Test* |
| W2VEC | 13.9 | 239 | 3.5 | 11.82 | 5.0 | 107 |
| D2VEC | 47.7 | 17.2 | 8.5 | 1.30 | 14.9 | 21.95 |
| S-BE | - | 2.6 | - | 1.16 | - | 7.5 |
| W-RW | 152 | 0.07 | 207 | 0.05 | 189 | 0.41 |
| RANK* | 3206 | 0.09 | 3916 | 0.05 | 6918 | 1.2 |
| L-BE* | 3616 | 0.25 | 3280 | 0.69 | 251 | 2.6 |
| TAPAS* | 7301 | 7.2 | - | - | - | - |
| DEEP-M* | 3492 | 0.68 | - | - | - | - |
| DITTO* | 34528 | 2.28 | - | - | - | - |

TABLE VII: Train and test execution times (sec).

compression method with the best results in all cases. For Corona it even does better than the expanded graph. This dataset contains many numerical values (about 25% of its data nodes in the expanded graph), which are misleading in some cases, as they are more likely to raise spurious connections in the graph. In general, *MSP* performs better than *SSuM*. For higher compression, *MSP (0.25)* produces smaller graphs with better match accuracy in most cases. For *MSP (0.5)*, we observe better quality in the matches in all cases except Audit, and comparable size in the compressed graphs. The results show the benefit of considering shortest paths among metadata nodes. For *MSP*, the graph size and the match accuracy follow the expected behavior w.r.t. the compression ratio.

*MSP*, in both executions, shows its best results for scenarios with at least one relational table. In these cases, the compressed graph is smaller than the original one (and much smaller than the expanded one), with better or very close matching quality. For text-only scenarios, the size reduction is remarkable, and better than *SSuM*, but a significant drop in matching quality can be observed. The conclusion is that the use graph compression depends on the kind of data and the requirements for the target application at hand.

### E. Execution Times.

Table VII reports execution times for all methods averaged over the experiments for every task. For training time, embeddings methods (W2VEC, D2VEC) and transformer-based methods (RANK*, L-BE*) are trained (fine tuned) on a smaller corpus than our method (W-RW). This is because we create 100 walks for each node in the graph, which leads to bigger corpora in general. Due to this difference, *Word2Vec* and *Doc2Vec* are faster than other methods in training, while our method has execution times smaller than those taken to fine tune transformers. S-BE has no training.
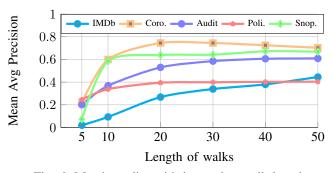


Fig. 6: Match quality with increasing walk length.

We report the average execution time for a single match (test). Our solution is the fastest. Document based methods, like *Doc2Vec*, are faster than word based embedding solutions. This is because for the latter methods we generate vectors for *all* tokens in the document and aggregate them. Classifiers are faster than document embeddings based matching, but slower than our method. In the training step of our method, expansion and compression take less than 3k seconds in all cases, with the exception of IMDB (by far the largest) with 79k seconds for expansion (with DBpedia) and 51k seconds for compression with *MSP (0.5)*.

### F. Ablation Study

We first report on the impact of different parameters on the performance of our method W-RW. We then evaluate the impact of the improvements proposed in Section II.

*1) Impact of parameters:* We examine the impact of length and number of random walks, followed by number of tokens in data nodes. We also report the impact of the graph size (in terms of number of nodes) on the execution time of our model. For *CoronaCheck*, we report results for the union of the *Generated* and *User* sentences.

**Length of random walks.** Figure 6 shows the mean average precision results for all scenarios when increasing the length of the random walks. Increasing the walk length increases the performance for all scenarios up to size 20. The increase is higher at lower values and then stabilizes or gradually decreases for most scenarios. We explain the different behavior for *IMDb* and *Audit* with the fact that they have the biggest and most dense graphs. *IMDb* graph is the biggest both in terms of nodes (107k nodes vs 35k node for *Snopes*) and edges (1m vs 168k edges for *Politifact*). Because of their size, larger graphs benefit of walks longer than 20.

**Number of walks.** Figure 7 shows the mean average precision when increasing the number of walks. The performance for all datasets improve with more walks, but with diminishing results. Results also confirm that graphs with more edges per node need more walks to obtain the best results. After 20 walks per node, results for *IMDb* keep improving, while for *CoronaCheck*, which is the most sparse graph with an average of four edges per node, there is no improvement.

**Number of tokens in terms.** Allowing more tokens in data nodes (terms) increases the mean average precision in all

| Dataset | Original Graph | | | Expanded Graph | | | MSP (0.5) | | | MSP (0.25) | | | SSuM (0.1) | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | #N | #E | MRR | #N | #E | MRR | #N | #E | MRR | #N | #E | MRR | #N | #E | MRR |
| IMDB | 107k | 1m | .780 | 237k | 1.5m | **.792** | 82k | 887k | .779 | 75k | 840k | .755 | **27k** | **540k** | .601 |
| Corona | 10k | 43k | .728 | 15k | 56k | .755 | 10k | 40k | **.769** | **8.5k** | **32.5k** | .757 | 10k | 33k | .610 |
| Snopes | 35k | 129k | .695 | 142k | 622k | **.708** | 84k | 479k | .647 | **49k** | **292k** | .586 | 83k | 470k | .590 |
| Politi | 24k | 168k | .489 | 62k | 317k | **.507** | 37k | 242k | .500 | **33k** | **225k** | .484 | 52k | 257k | .397 |
| Audit | 6k | 164k | .421 | 17k | 202k | **.452** | 7k | 161k | .389 | **5.5k** | **144k** | .362 | 14k | 150k | .392 |

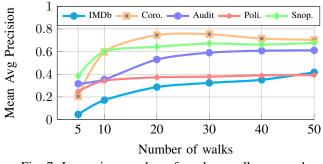TABLE VIII: Compression performance: number of graph nodes (#N) and edges (#E) compared with matching quality MRR.

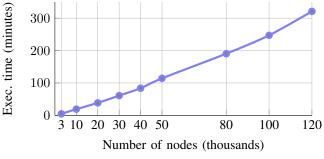

Fig. 7: Increasing number of random walks per node.



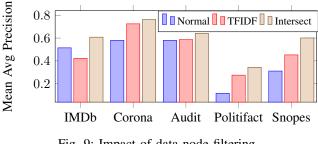Fig. 8: Increasing number of nodes in the graph.



Fig. 9: Impact of data node filtering.

size of the graphs, we expand them using Algorithm 2 and *ConceptNet*. For every graph, we generate 100 random walks of length 30 for each node and report in Figure 8 the total time to generate random walks and to train the word embeddings. Results show that by increasing the number of nodes in the graph, the execution time increases linearly.

*2) Improving graph generation:* We discuss here the impact of the techniques introduced in Section II to improve embeddings and matching quality.
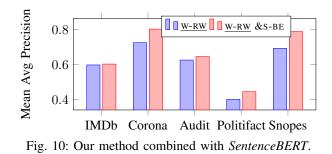
**Connecting metadata nodes.**

For measuring the impact of edges between related metadata nodes in a structured text, we run the same experiments in Section V-B for W-RW without edges among metadata nodes. The quality of the matches is negatively affected, with the Node F-scores for increasing values of $K$ (1, 3, 5, 10) dropping by .08, 0.04, 0.02, and 0.01 in absolute values.

**Filtering data nodes.** Merging tokens create shorter paths between related metadata nodes. In graph creation, we filter tokens of the second corpus based on the nodes from the first one. In this experiment, we compare the performance of our technique (Intersect) against a solution based on *TF-IDF* [2], which retains non-stopword tokens with high *TF-IDF* scores. It has been shown that in a data to data matching task this technique increases the quality performance on a text-heavy dataset from 41% to 93% [2]. For implementing this method, for each document we keep $k$ tokens with highest TF-IDF scores for different $k$ values. For each scenario, we run $k =3$, 5, 10 and 20 and report the best result.

Figure 9 shows the performance of these two techniques in terms of mean average precision for all scenarios. The results show that, except for *IMDb* dataset and *TF-IDF*, both summarizing techniques improve the mean average precision of matching. It also shows that our technique works better than *TF-IDF* in all scenarios.

**Combining matching scores.** Figure 10 depicts the results

scenarios. There is a significant increase in quality going from one to two tokens and the impact is smaller with higher values. The highest increase is for *Snopes* and *Audit* datasets with an increase up to 0.24 and 0.07, respectively. For text to data scenarios, *IMDb* has an increase up to 0.05 and *CoronaCheck* up to 0.03. The amount of increase for a scenario is related to the number of new nodes added to the graph. For *Snopes*, an increase in the number of tokens in a term adds an average of 24k new nodes to the graph, which is close to the initial number of nodes (35K). For *Audit*, the increase is 1.7k, which is 62% of the initial graph's token nodes. In *IMDb*, 23K new nodes are generated (22% of initial graph) at every increase. For all scenarios, except *Snopes*, the number of nodes does not increase drastically after allowing three tokens in a term. The behaviour of a dataset is determined by the length of documents in its first corpus as tokens in the documents of the second corpus get filtered. *Snopes* has the biggest documents in its first corpus in comparison to the other datasets: it has claims of 43 tokens on average while *PolitiFact* has 18.

**Number of nodes in the graph.** To show that our method can scale on a commodity machine, we generate graphs of increasing size. The graphs are generated from the *STS* dataset with $k$ values between 0 and 5. For increasing further the

Fig. 10: Our method combined with *SentenceBERT*.

of averaging the cosine similarity scores from our solution with those from the pre-trained S-BE. Our solution already outperforms S-BE in the original experiments, but averaging the two methods improves the matching quality in all scenarios even above the results from our methods alone. The biggest improvements are for *Snopes* and *CoronaCheck* with 0.9 and 0.8 increase, respectively. This simple combination show the benefit of exploiting domain specific embeddings and pre-trained ones from large, generic corpora.

**Merging nodes.** For all scenarios, we employ different techniques and resources to improve their quality results (Section II-C). In *CoronaCheck*, 17% of the nodes are numeric values and merging with equal-width buckets decreases the number of graph nodes. We had the best results with equal-width buckets of size 7, which increased the mean average precision from 0.72 to 0.76. In *Audit*, *Politifact*, and *Snopes* datasets, there are few numeric values and bucketing had no effects. In *IMDB*, we observe a small loss because numeric values are release dates, for which is better to avoid merging.

We also use *Wikipedia2Vec* to merge similar nodes. *IMDb* contains variations for the same entity (e.g., director names) and merging them with a threshold $\gamma = 0.57$ increases the performance by 2.5%. For *Snopes* and *Politifact*, using the same value, the increase is by 1.7% and 1.5%, respectively. As *CoronaCheck* contains typos in user sentences (e.g., country names), merging such typos leads to a 3.4% increase. *Auditing* does not show improvements by merging nodes with pre-trained resources because of the difference between the general meaning of a term and its meaning in a specific domain. Models pre-trained on general corpora do not help much in a domain specific scenario. For example, the *Wikipedia2Vec* similarity is high between *Financial statement* and *Financial reporting*, same for *Auditor* and *Risk control*. However, these terms have different meanings in audit documents.

## VI. RELATED WORK

Entity resolution for relational data [3] has been recently studied with deep learning solutions [1], [8], [2]. As they rely on the presence of a schema, one way to use them in this setting is to treat the text paragraphs as tuples within a single column, but this leads to poor results. Even approaches that consider schema-agnostic attribute matching, assume a semi-structured input [42], [43]. In contrast, our solution models texts and tables in a unified representation for learning embeddings, with benefits in the matching tasks even with very lightweight algorithms. Compared to [1], we extended the graph generation to handle free text (without schema) with new merging and tokenization techniques and with a new expansion and compression approach. We model text matching both as a binary classification task and as a multi-label classification task to use SOTA baselines based on fine tuning language models [39], [44]. These methods outperform traditional IR approaches, such as BM25, but do not focus on the problem of matching text and relational datasets. This latter problem has been studied in settings that do not cover our use cases, either because they assume supervision or because they are no domain specific [45], [7], [6]. Other approaches for text and data matching assume a very expensive training over large document corpora and millions of tables [13], [12], [14]. We report results for the fine-tuning on top of these pre-trained model (TAPAS*), but do not report results based on pre-training as they are significantly worse in terms of quality w.r.t. our solution. Indeed, the pre-training is not designed for a setting with a small table corpus. Our setting is also different from the problem of entity linking, as we are matching long text to tuples in relational data [21], but it could be used as a pre-processing tool to provide correspondences to knowledge translation solutions [46].

Our default method to generate data and metadata representations is *Word2Vec* [47] as it is powerful in discovering relationships in the corpus as well as similarity between tokens [48]. While we generate embeddings, some of our baselines use pre-trained ones [15], [10]. Document embedding methods model longer text sequences by aggregating the vectors of words in the given sentence or paragraph [49], [38] or by learning the document vector with special tokens [44], [50]. Given our graph, it is also possible to generate embeddings directly for its nodes [51], [36], [37] with comparable results w.r.t. the random walks followed by word embedding generation in the data to data matching task [1]. Our study confirms that they do not bring clear benefit, but are more resources intensive than *Word2Vec*. Finally, our work can be seen as a new instance of the general approaches of using deep learning for data integration [52], [53], [17] and of improving pre-trained embeddings w.r.t. relational data [54].

## VII. CONCLUSION

We presented a new generic matching task that allows both structured text documents and relational data. Results show that lightweight embeddings effectively model the similarity between heterogeneous corpora. Our proposal outperforms in quality and test time all unsupervised baselines and it is competitive to supervised solutions. Our graph expansion always leads to the best matching quality, while the compression is effective in reducing the graph size, but comes with a trade off in the performance of the matching for text-only corpora. We plan to extend our framework to support a richer graph with typed edges and blocking to speed up performance. However, while word embeddings are effective and efficient to generate, a more complex architecture should be considered to enable the benefits of fine tuning when examples are available.

REFERENCES

[1] R. Cappuzzo, P. Papotti, and S. Thirumuruganathan, "Creating embeddings of heterogeneous relational datasets for data integration tasks," in *SIGMOD*, 2020.

[2] Y. Li, J. Li, Y. Suhara, A. Doan, and W. Tan, "Deep entity matching with pre-trained language models," *PVLDB*, 2021.

[3] L. Getoor and A. Machanavajjhala, "Entity resolution for big data," in *ACM SIGKDD*, 2013, pp. 1527–1527.

[4] P. Rajpurkar, R. Jia, and P. Liang, "Know what you don't know: Unanswerable questions for squad," in *ACL*, 2018, pp. 784–789.

[5] J. Guo, Y. Fan, X. Ji, and X. Cheng, "Matchzoo: a learning, practicing, and developing system for neural text matching," in *SIGIR*, 2019, pp. 1297–1300.

[6] W. Chen, H. Wang, J. Chen, Y. Zhang, H. Wang, S. Li, X. Zhou, and W. Y. Wang, "Tabfact: A large-scale dataset for table-based fact verification," in *ICLR*, 2020.

[7] G. Karagiannis, M. Saeed, P. Papotti, and I. Trummer, "Scrutinizer: A mixed-initiative approach to large-scale, data-driven claim verification," *Proc. VLDB Endow.*, vol. 13, no. 11, pp. 2508–2521, 2020.

[8] M. Ebraheem, S. Thirumuruganathan, S. R. Joty, M. Ouzzani, and N. Tang, "Distributed representations of tuples for entity resolution," *Proc. VLDB Endow.*, vol. 11, no. 11, pp. 1454–1467, 2018.

[9] S. Mudgal, H. Li, T. Rekatsinas, A. Doan, Y. Park, G. Krishnan, R. Deep, E. Arcaute, and V. Raghavendra, "Deep learning for entity matching: A design space exploration," in *Proceedings of the 2018 International Conference on Management of Data*, 2018, pp. 19–34.

[10] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova, "Bert: Pre-training of deep bidirectional transformers for language understanding," *arXiv preprint arXiv:1810.04805*, 2018.

[11] M. Zaheer, G. Guruganesh, A. Dubey, J. Ainslie, C. Alberti, S. Ontañón, P. Pham, A. Ravula, Q. Wang, L. Yang, and A. Ahmed, "Big Bird: Transformers for longer sequences," *CoRR*, vol. abs/2007.14062, 2020. [Online]. Available: https://arxiv.org/abs/2007.14062

[12] P. Yin, G. Neubig, W. Yih, and S. Riedel, "Tabert: Pretraining for joint understanding of textual and tabular data," in *ACL*, 2020, pp. 8413–8426.

[13] J. Herzig, P. K. Nowak, T. Müller, F. Piccinno, and J. M. Eisenschlos, "Tapas: Weakly supervised table parsing via pre-training," in *ACL*, 2020, pp. 4320–4333.

[14] X. Deng, H. Sun, A. Lees, Y. Wu, and C. Yu, "TURL: table understanding through representation learning," *Proc. VLDB Endow.*, vol. 14, no. 3, pp. 307–319, 2020.

[15] J. Pennington, R. Socher, and C. D. Manning, "Glove: Global vectors for word representation," in *EMNLP*, 2014, pp. 1532–1543.

[16] R. Speer, J. Chin, and C. Havasi, "Conceptnet 5.5: An open multilingual graph of general knowledge," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 31, no. 1, 2017.

[17] C. Koutras, G. Siachamis, A. Ionescu, K. Psarakis, J. Brons, M. Fragkoulis, C. Lofi, A. Bonifati, and A. Katsifodimos, "Valentine: Evaluating matching techniques for dataset discovery," *CoRR*, vol. abs/2010.07386, 2020. [Online]. Available: https://arxiv.org/abs/2010.07386

[18] D. Freedman and P. Diaconis, "On the histogram as a density estimator: L2 theory," *Z. Wahrscheinlichkeitstheorie und verwandte Gebiete*, vol. 57, no. 4, pp. 453–476, 1981.

[19] G. A. Miller, *WordNet: An electronic lexical database.* MIT press, 1998.

[20] I. Yamada, H. Shindo, H. Takeda, and Y. Takefuji, "Joint learning of the embedding of words and entities for named entity disambiguation," in *SIGNLL*, 2016.

[21] H. Chen, S. Wadhwa, X. D. Li, and A. Zukov-Gregoric, "Yelm: End-to-end contextualized entity linking," *arXiv preprint arXiv:1911.03834*, 2019.

[22] E. J. Gerritse, F. Hasibi, and A. P. de Vries, "Graph-embedding empowered entity retrieval," in *European Conference on Information Retrieval*, 2020, pp. 97–110.

[23] P. Mcnamee and J. Mayfield, "Character n-gram tokenization for european language text retrieval," *Information retrieval*, vol. 7, no. 1-2, pp. 73–97, 2004.

[24] C. Bizer, J. Lehmann, G. Kobilarov, S. Auer, C. Becker, R. Cyganiak, and S. Hellmann, "DBpedia-A crystallization point for the web of data," *Web Semantics: science, services and agents on the WWW*, vol. 7, no. 3, pp. 154–165, 2009.

[25] R. Speer, J. Chin, and C. Havasi, "Conceptnet 5.5: An open multilingual graph of general knowledge," in *AAAI*. AAAI Press, 2017, pp. 4444–4451.

[26] L. A. Adamic, R. M. Lukose, A. R. Puniyani, and B. A. Huberman, "Search in power-law networks," *Physical review E*, vol. 64, no. 4, p. 046135, 2001.

[27] Y. Liu, T. Safavi, A. Dighe, and D. Koutra, "Graph summarization methods and applications: A survey," *ACM Comput. Surv.*, vol. 51, no. 3, pp. 62:1–62:34, 2018.

[28] M. P. Stumpf, C. Wiuf, and R. M. May, "Subnets of scale-free networks are not scale-free: sampling properties of networks," *Proceedings of the National Academy of Sciences*, vol. 102, no. 12, pp. 4221–4224, 2005.

[29] J. Leskovec and C. Faloutsos, "Sampling from large graphs," in *Proceedings of the 12th ACM SIGKDD international conference on Knowledge discovery and data mining*, 2006, pp. 631–636.

[30] V. Krishnamurthy, M. Faloutsos, M. Chrobak, L. Lao, J.-H. Cui, and A. G. Percus, "Reducing large internet topologies for faster simulations," in *International Conference on Research in Networking*. Springer, 2005, pp. 328–341.

[31] N. K. Ahmed, J. Neville, and R. Kompella, "Network sampling: From static to streaming graphs," *ACM Transactions on Knowledge Discovery from Data (TKDD)*, vol. 8, no. 2, pp. 1–56, 2013.

[32] J. Leskovec, J. Kleinberg, and C. Faloutsos, "Graphs over time: densification laws, shrinking diameters and possible explanations," in *Proceedings of the eleventh ACM SIGKDD international conference on Knowledge discovery in data mining*, 2005, pp. 177–187.

[33] A. Rezvanian and M. R. Meybodi, "Sampling social networks using shortest paths," *Physica A: Statistical Mechanics and its Applications*, vol. 424, pp. 254–268, 2015.

[34] Y. Li, Z. Wu, S. Lin, H. Xie, M. Lv, Y. Xu, and J. C. Lui, "Walking with perception: Efficient random walk sampling via common neighbor awareness," in *2019 IEEE 35th International Conference on Data Engineering (ICDE)*. IEEE, 2019, pp. 962–973.

[35] W. L. Hamilton, R. Ying, and J. Leskovec, "Representation learning on graphs: Methods and applications," *arXiv preprint arXiv:1709.05584*, 2017.

[36] B. Perozzi, R. Al-Rfou, and S. Skiena, "Deepwalk: Online learning of social representations," in *SIGKDD*, 2014, pp. 701–710.

[37] A. Grover and J. Leskovec, "node2vec: Scalable feature learning for networks," in *SIGKDD*, 2016, pp. 855–864.

[38] C. De Boom, S. Van Canneyt, T. Demeester, and B. Dhoedt, "Representation learning for very short texts using weighted word embedding aggregation," *Pattern Recognition Letters*, vol. 80, pp. 150–156, 2016.

[39] S. Shaar, N. Babulkov, G. D. S. Martino, and P. Nakov, "That is a known lie: Detecting previously fact-checked claims," in *ACL*, 2020, pp. 3607–3618.

[40] D. Cer, M. Diab, E. Agirre, I. Lopez-Gazpio, and L. Specia, "Semeval-2017 task 1: Semantic textual similarity-multilingual and cross-lingual focused evaluation," *arXiv preprint arXiv:1708.00055*, 2017.

[41] K. Lee, H. Jo, J. Ko, S. Lim, and K. Shin, "Ssumm: Sparse summarization of massive graphs," in *SIGKDD*, 2020, pp. 144–154.

[42] G. M. Mandilaras, G. Papadakis, L. Gagliardelli, G. Simonini, E. Thanos, G. Giannakopoulos, S. Bergamaschi, T. Palpanas, M. Koubarakis, A. Lara-Clares, and A. Fariña, "Reproducible experiments on three-dimensional entity resolution with JedAI," *Inf. Syst.*, vol. 102, 2021.

[43] R. Singh, V. V. Meduri, A. K. Elmagarmid, S. Madden, P. Papotti, J. Quiané-Ruiz, A. Solar-Lezama, and N. Tang, "Synthesizing entity matching rules by examples," *Proc. VLDB Endow.*, vol. 11, no. 2, pp. 189–202, 2017.

[44] A. Adhikari, A. Ram, R. Tang, and J. Lin, "Docbert: Bert for document classification," *arXiv preprint arXiv:1904.08398*, 2019.

[45] Y. Ibrahim, M. Riedewald, G. Weikum, and D. Zeinalipour-Yazti, "Bridging quantities in tables and text," in *ICDE*. IEEE, 2019, pp. 1010–1021.

[46] B. G. Bashardoost, R. J. Miller, K. A. Lyons, and F. Nargesian, "Knowledge translation," *Proc. VLDB Endow.*, vol. 13, no. 11, pp. 2018–2032, 2020.

[47] T. Mikolov, K. Chen, G. Corrado, and J. Dean, "Efficient estimation of word representations in vector space," *arXiv preprint arXiv:1301.3781*, 2013.

[48] K. Kowsari, K. Jafari Meimandi, M. Heidarysafa, S. Mendu, L. Barnes, and D. Brown, "Text classification algorithms: A survey," *Information*, vol. 10, no. 4, p. 150, 2019.

[49] R. Kiros, Y. Zhu, R. R. Salakhutdinov, R. Zemel, R. Urtasun, A. Tor-ralba, and S. Fidler, "Skip-thought vectors," in *NIPS*, 2015, pp. 3294–3302.

[50] N. Reimers and I. Gurevych, "Sentence-bert: Sentence embeddings using siamese bert-networks," *arXiv preprint arXiv:1908.10084*, 2019.

[51] P. Goyal and E. Ferrara, "Graph embedding techniques, applications, and performance: A survey," *Knowledge-Based Systems*, vol. 151, pp. 78–94, 2018.

[52] C. Fu, X. Han, L. Sun, B. Chen, W. Zhang, S. Wu, and H. Kong, "End-to-end multi-perspective matching for entity resolution," in *IJCAI*, 2019.

[53] R. C. Fernandez and S. Madden, "Termite: a system for tunneling through heterogeneous data," in *aiDM*. ACM, 2019, pp. 7:1–7:8.

[54] M. Günther, P. Oehme, M. Thiele, and W. Lehner, "Learning from textual data in database systems," in *CIKM*, 2020, pp. 375–384.