

# Match-Ignition: Plugging PageRank into Transformer for Long-form Text Matching

Liang Pang<sup>1</sup>, Yanyan Lan<sup>2\*</sup>, Xueqi Cheng<sup>3</sup>

pangliang@ict.ac.cn, lanyanyan@tsinghua.edu.cn, cxq@ict.ac.cn

<sup>1</sup>Data Intelligence System Research Center,

Institute of Computing Technology, Chinese Academy of Sciences, Beijing, China

<sup>2</sup>Institute for AI Industry Research, Tsinghua University, Beijing, China

<sup>3</sup>CAS Key Lab of Network Data Science and Technology,

Institute of Computing Technology, Chinese Academy of Sciences, Beijing, China

## ABSTRACT

Neural text matching models have been widely used in community question answering, information retrieval, and dialogue. However, these models designed for short texts cannot well address the long-form text matching problem, because there are many contexts in long-form texts can not be directly aligned with each other, and it is difficult for existing models to capture the key matching signals from such noisy data. Besides, these models are computationally expensive for simply use all textual data indiscriminately. To tackle the effectiveness and efficiency problem, we propose a novel hierarchical noise filtering model, namely Match-Ignition. The main idea is to plug the well-known PageRank algorithm into the Transformer, to identify and filter both sentence and word level noisy information in the matching process. Noisy sentences are usually easy to detect because previous work has shown that their similarity can be explicitly evaluated by the word overlapping, so we directly use PageRank to filter such information based on a sentence similarity graph. Unlike sentences, words rely on their contexts to express concrete meanings, so we propose to jointly learn the filtering and matching process, to well capture the critical word-level matching signals. Specifically, a word graph is first built based on the attention scores in each self-attention block of Transformer, and key words are then selected by applying PageRank on this graph. In this way, noisy words will be filtered out layer by layer in the matching process. Experimental results show that Match-Ignition outperforms both SOTA short text matching models and recent long-form text matching models. We also conduct detailed analysis to show that Match-Ignition efficiently captures important sentences and words, to facilitate the long-form text matching process.

## CCS CONCEPTS

• Information systems → Retrieval models and ranking; • Computing methodologies → Neural networks.

\*Corresponding author

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from [permissions@acm.org](https://permissions.acm.org).

CIKM '21, November 1–5, 2021, Virtual Event, QLD, Australia

© 2021 Association for Computing Machinery.

ACM ISBN 978-1-4503-8446-9/21/11...\$15.00

<https://doi.org/10.1145/3459637.3482450>

### Short-form Text Matching

(Paraphrasing Identification):

Sentence1: *What makes a pizza the best?*

| | | X

Sentence2: *What makes a good pizza?*

### Long-form Text Matching

(Redundancy News Identification):

Doc1: ... it took unreal shooting by Tracy McGrady. **McGrady made four 3-point shots and scored 13 points in the final 35 seconds, including a 26-footer with 1.7 seconds remaining, to send the Spurs to San Antonio with a second loss in as ...**

Doc2: ... **McGrady had one of the most memorable performances of his career, the final 35 seconds win 13 points when against the San Antonio Spurs to secure a comeback victory. The sequence included four consecutive three-pointers ...**

Figure 1: The top example is a short-form text matching for the paraphrasing identification, and the lines indicate the alignments between words from two sentences. The bottom example is a long-form text matching for the redundancy news identification, and the highlights indicate the important matching signals for the identity event of two news.

## KEYWORDS

Text Matching; Long-form Text; PageRank Algorithm

### ACM Reference Format:

Liang Pang, Yanyan Lan, Xueqi Cheng. 2021. Match-Ignition: Plugging PageRank into Transformer for Long-form Text Matching. In *Proceedings of the 30th ACM International Conference on Information and Knowledge Management (CIKM '21)*, November 1–5, 2021, Virtual Event, QLD, Australia. ACM, New York, NY, USA, 10 pages. <https://doi.org/10.1145/3459637.3482450>

## 1 INTRODUCTION

Semantic text matching is an essential problem in many natural language applications, such as community question answering [42], information retrieval [17], and dialogue [25]. Many deep text matching models have been proposed and gain good performance, such as representation based models [17, 28, 37, 40], interaction based models [13, 16, 29, 41], and their combinations [9, 27, 44].

However, these models cannot be well applied to long-form text matching problems, which have attracted increasing attention in the field of news deduplication [22], citation recommendation [43], plagiarism detection [47] and attachment suggestion [19]. This is

mainly because long-form text matching is quite different from the short-form text matching problem. For short-form text matching, almost every term in the short texts is critical to the matching score, because short text matching tasks are just like finding a reasonable semantic alignment between two sentences [29]. For example, in paraphrasing identification, the major problem is to find the paraphrasing sentence for the given sentence. In this case, the matching score is mainly determined by the alignment between each word in the sentences, as shown in Figure 1.

Long-form text matching has its own characteristics. Firstly, long-form text matching cares more about the global semantic meanings rather than the bipartite alignment. The fine-grained matching signals between long-form texts are usually very sparse, which makes the existing short text matching models hard to figure out from huge noisy signals. For example, redundant news identification merely focuses on where/when the event happened and what the event is, instead of who posted this news and the detailed descriptions of the news. Secondly, long-form text intrinsically consists of a two-level structure, i.e. sentences and words. Most existing short text matching approaches can only process text word by word while missing the sentence-level structure. For example, one sentence should be ignored entirely if it is irrelevant to the current document, e.g. advertisement, even though some of its internal words are relevant. Thirdly, long-form text matching contains a very long text by nature, which makes the existing short text matching models computational expensive because they have to treat every word indiscriminately and emphasize the sufficient interactions between words [9, 44]. In practice, the long-form text often has to be truncated in the computation. For example, BERT only accepts text lengths of less than 512. These operations may hurt the final matching performance. From these discussions, we can see that noise is the main challenge in long-form text matching, affecting both performance and efficiency.

In this paper, we propose a novel hierarchical noise filtering model, namely Match-Ignition, to distill the significant matching signals via the well-known link analysis algorithm PageRank [5]. PageRank utilizes random walk on a graph to determine the importance of each node. In this way, the noises (i.e. less important nodes) can be eliminated and the algorithm will be accelerated. Considering the two-level structures in the long-form text matching problem, our model contains two hierarchies, i.e. sentence-level and word-level. In the sentence-level noise filtering process, the nodes are defined as sentences from a pair of long-form texts, and the link is defined as the similarities between each pair of sentences. That is to say, the similarities inside each long-form text and between the two long-form texts are both captured in our graph. Then the noisy sentences could be identified by PageRank score, and be directly removed. The word-level noise filtering process is jointly learned with the matching process. That is because each word relies on its context to express its concrete meanings, so noisy words need to be estimated dynamically during the matching process. To this end, we first apply the state-of-the-art Transformer to the texts, which well captures the contextual information among words. It turns out that the attention matrix in the self-attention block, the key component of Transformer, could be treated as a fully connected word-level similarity graph [8, 15, 46]. PageRank is then applied to filter out noise words at each layer. We can see this technique is different

from previous works [8, 15, 46] which focus on eliminating links in the graph because our model focuses on filtering noisy words, i.e., nodes in the graph. Please note that attention weights could be directly applied to filter noisy words. The reason why we still use PageRank here is that the attention weights have been proven not reliable in explaining the importance of words in [6]. Furthermore, PageRank has the ability to consider the global importance of each word by value propagation on a graph, which is more thorough than attention weights.

We experiment on three long-form text matching tasks, news deduplication, citation recommendation, and plagiarism detection, including seven public datasets, e.g. CNSE, CNSS, AAN-Abs, AAN-Body, OC, S2ORC, and PAN. The experimental results show that Match-Ignition outperforms all baseline methods, including both short text matching models and recent long-form text matching models. The further detailed analysis demonstrates that Match-Ignition efficiently captures important matching signals in long-form text, which helps understand the matching process. Besides, we compare different noisy filtering methods to show the superiority of using PageRank.

## 2 RELATED WORK

In this section, we first introduce the text matching models designed for short-form text matching, then review the most recent works for long-form text matching.

**Short-form Text Matching** Existing text matching models fall into representation-based approaches, interaction-based approaches, and their combinations [14].

Representation-based matching approaches are inspired by the Siamese architecture [7]. This kind of approach aims at encoding each input text in a pair into the high-level representations respectively based on a specific neural network encoder, and then the matching score is obtained by calculating the similarity between the two corresponding representation vectors. DSSM [17], C-DSSM [37], ARC-I [16], RNN-LSTM [28] and MV-LSTM [40] belong to this category. Interaction-based matching approaches are closer to the nature of the matching task to some extent since they aim at directly capturing the local matching patterns between two input text, rather than focusing on the text representations. The pioneering work includes ARC-II [16], MatchPyramid [29], and Match-SRNN [41]. Recently, there has been a trend that the two aforementioned branches of matching models should complement each other, rather than being viewed separately as two different approaches. DUET [27] is composed of two separated modules, one in a representations-based way, and another in an interaction-based way, the final matching score is just their weighted-sum result. The attention mechanism is another way to combine the above two approaches, such as RE2 [44] and BERT [9]. However, these existing approaches for short-form text matching have limited success in long-form text matching settings, due to their inability to capture and distill the main information from long documents. Besides, these models are computationally expensive because they simply use all textual data indiscriminately in the matching process.

**Long-form Text Matching** Due to the lack of the public datasets and the efficient algorithms, few work directly focuses on the long-form text matching and further explores the application scenarios

of it. In recent years, since the pioneering work SMASH proposed by Jiang et al. [19], they are the first to point out that long-form text matching (e.g. source text and target text both are long-form text) has a wide range of application scenarios, such as attachment suggestion, article recommendation, and citation recommendation. They propose a hierarchical recurrent neural network under Siamese architecture which is a kind of representation-based matching approach. It synthesizes information from different document structure levels, including paragraphs, sentences, and words. SMITH model [43] follows the SMASH’s settings, then utilizes powerful pre-trained language model BERT [9] as their key component and breaks the 512 tokens limitation to build a representation-based matching approach. Instead of using BERT as a component, TransformerXL [8] and Longformer [2] try to directly extent the Transformer structure towards the long-form text by introducing a sliding window or memory strategy into it. However, they are designed for general natural language understanding, not for text matching tasks. Another work on long-form text matching is Concept Interaction Graph (CIG) [22], which concerns modeling the relation between two documents, e.g. same event or story. It can be treated as an interaction-based matching approach, which selects a pair of sentences based on their concepts and similarities. Besides, they also construct two types of duplicate news detection datasets, which are labeled by professional editors.

All the previous works ignore the fact that long-form text provides overabundant information for matching, that is to say, there are usually many noises in the setting of long-form text matching. This phenomenon also is discussed in query-document matching tasks [11, 13, 18, 30], which is a short to long text matching because a query is a short-form text and a document is a long-form text. DeepRank [30] treats query and document differently, in their model, each query term acts as a filter that picks out text spans in the document which contain this query term. That is to say, query irrelevant text spans are the noise that can be ignored in the matching process. PACRR [18] also has similar findings, they filter document words using two kinds of processes, 1) keep first  $k$  terms in the document or 2) retain only the text that is highly relevant to the given query. These previous works provide strong evidence that our noise filtering motivation can be effective for long-form text matching problems.

### 3 MATCH-IGNITION

In this section, we first introduce the two components of Match-Ignition. They are sentence-level noise filter and word-level noise filter, shown in Figure 2(a) and Figure 2(c) respectively. After that, the model training details are described in the last subsection.

#### 3.1 Sentence-level Noise Filtering

To enable the application of graph-based ranking algorithms PageRank to natural languages, such as documents, a graph is needed to build that represents the relation between sentences. TextRank [26] makes it possible to form a sentence extraction algorithm, which can identify key sentences in a given document. It becomes a mature approach in automatic summarization. The most direct way is to apply the TextRank algorithm on each long-form text independently, so that we can reduce the length of each long-form text by

summarizing. However, the goal of long-form text matching is to find the matching signals between a pair of texts, which is different from summarization task that extracts key information from one text. Therefore, directly applying the TextRank algorithm to each text independently leads to the problem of loss of matching signals.

Inspired by the previous works [18, 30], they tell us that two texts can help each other for noise detection, so that both long-form texts should be represented in one graph to involve the matching information across two texts. Firstly, sentences in both long-form texts are collected together to form a united sentence collection. Formally, two long-form texts are first split into sentences, denoted as  $d_s = [s_1^1, s_2^1, \dots, s_{L_1}^1]$  and  $d_t = [s_1^2, s_2^2, \dots, s_{L_2}^2]$ , where  $L_1$  and  $L_2$  are the number of sentences in  $d_s$  and  $d_t$  respectively. The united sentence collection  $\mathcal{S} = \{s_1^1, s_2^1, \dots, s_{L_1}^1, s_1^2, s_2^2, \dots, s_{L_2}^2\}$  then have  $L_1 + L_2$  elements. Thus, the sentence similarity graph can be constructed by evaluating the sentence pair similarities in the united sentence collection  $\mathcal{S}$ . The sentence similarity is defined as the same as in TextRank [26], to measures the overlapping word ratio between two sentences:

$$Sim(s_i, s_j) = \frac{|\{w_k | w_k \in s_i, w_k \in s_j\}|}{\log(|s_i|) + \log(|s_j|)}, \quad s_i, s_j \in \mathcal{S}, \quad (1)$$

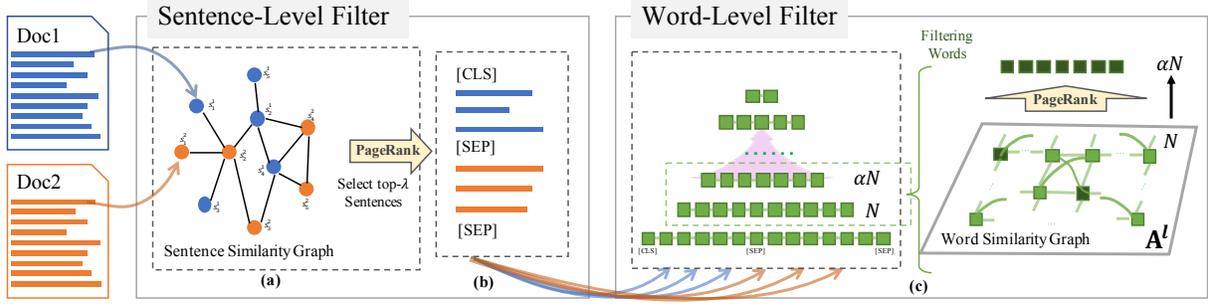
where  $w_k$  denotes the word in the sentence,  $|\cdot|$  denotes the length of the sentence or word set, and  $s_i, s_j$  are two sentences in the united sentence collection  $\mathcal{S}$ . To make sentence similarity sparsity e.g. returns 0 most of the time, we remove the stopwords in the sentences before we calculate the similarities. Thus, the final sentence similarity graph has sparse links. Finally, a PageRank algorithm is applied to this constructed sentence similarity graph, to get the important score of each sentence. To balance the information coming from different long-form texts for the following step, the top  $\lambda$  sentences are extracted for each long-form texts respectively. Thus, both texts contain  $\lambda$  sentences as their digestion, which we called a sentence-level filter.

As shown in Figure 2(b), the selected sentences are concatenated as a text sequence, which starts with [CLS] token and separates by [SEP] token. It is then treated as the input of the model in the word-level filter. Note that the hyper-parameter  $\lambda$  should be neither too small to lose a lot of information, nor too large to make text extremely long. A suitable  $\lambda$  can yield a moderate text sequence, which length is just less than the BERT max input length.

PageRank algorithm can also be used at the word level if we can define a word-by-word relation graph. However, sentences are adjectives from each other, noise in this level is discrete than an entire sentence can be removed in an unsupervised way, while a word relies on its context to express concrete meanings, noise in this level is continuous that should be estimated during the model training. Therefore, we need to construct a graph within the Transformer model structures.

#### 3.2 Word-level Noise Filtering

To filter the noise in the word level, a word-level graph needs to be constructed first in the inherent transformer structure (Sec 3.2.1). After that, the traditional PageRank algorithm is required to implement as a tensor version, for better to embed into the transformer structure (Sec 3.2.2). Finally, we propose our plug PageRank to the Transformer model for word-level noise filtering (Sec 3.2.3).



**Figure 2: The overall architecture of Match-Ignition. (a) represents the sentence-level filter, (b) represents the outputs of the sentence-level filter, and (c) represents the word-level filter.**

**3.2.1 Transformer as a Graph.** Transformer architecture [39] boosts the natural language processing a lot, where most well-known models are a member of this family, such as BERT [10], RoBERTa [23], and GPT2 [34]. They achieve state-of-the-art performance in almost all NLP tasks, e.g. named entity recognition, text classification, machine translation, and also text semantic matching. For long-form text matching, we also adopt this architecture.

The self-attention block is the main component in Transformer architecture, which figure out how important all the other words in the sentence are for the contextual word around it. Thus, the self-attention block builds the relations between words, which can be viewed as a fully connected graph among words [8, 15, 46]. Knowing that the updated word representations are simply the sum of linear transformations of representations across all the words, weighted by their importance. It makes full use of the attention mechanism in deep neural networks to update word representations. As have shown in [39], the attention function can be formalized as a scaled dot-product attention with inputs  $\mathbf{H}^l$ :

$$\mathbf{H}^{l+1} = \text{Attn}(\mathbf{Q}^l, \mathbf{K}^l, \mathbf{V}^l) = \text{Softmax}\left(\frac{\mathbf{Q}^l(\mathbf{K}^l)^T}{\sqrt{E}}\right)\mathbf{V}^l = \mathbf{A}^l\mathbf{V}^l, \quad (2)$$

where  $\mathbf{Q}^l = \mathbf{W}_Q^l\mathbf{H}^l \in \mathbb{R}^{N \times E}$  denote the attention query matrices,  $\mathbf{K}^l = \mathbf{W}_K^l\mathbf{H}^l \in \mathbb{R}^{N \times E}$  the key matrix, and  $\mathbf{V}^l = \mathbf{W}_V^l\mathbf{H}^l \in \mathbb{R}^{N \times E}$  the value matrix.  $N$  denotes the number of words in a text, and  $E$  denotes the dimensions of the representation. The attention mechanism can be explained as: for each attention query vector in  $\mathbf{Q}$ , it first computes the dot products of the attention query with all keys, aiming to evaluate the similarity between the attention query and each key. Then, it is divided each by  $\sqrt{E}$ , and applies a softmax function to obtain the weights on the values, denotes as  $\mathbf{A}^l$ . Finally, the new representation of the attention query vector is calculated as weighed sum of values. Getting this dot-product attention mechanism to work proves to be tricky bad random initializations can de-stabilize the learning process. It can be overcome by performing multiple ‘heads’ of attention and concatenating the result:

$$\begin{aligned} \mathbf{H}^{l+1} &= \text{Concat}(\text{head}_1, \dots, \text{head}_H)\mathbf{O}^l, \\ \text{head}_k &= \text{Attention}(\mathbf{Q}^{kl}, \mathbf{K}^{kl}, \mathbf{V}^{kl}) = \mathbf{A}^{kl}\mathbf{V}^{kl}, \end{aligned} \quad (3)$$

where  $\mathbf{Q}^{kl}$ ,  $\mathbf{K}^{kl}$  and  $\mathbf{V}^{kl}$  are of the  $k$ -th attention head at layer  $l$  with different learnable weights,  $\mathbf{O}^l$  down-projection to match the dimensions across layers,  $H$  is the number of the heads in each layer and  $L$  is the number of the layers.

If we treat each word as a node in a graph, they update their representations by aggregating all other contextual word representations, just like messages passing from other neighbor nodes in graph neural network [36]. Thus, for self-attention block, it can be treated as a fully-connected word graph, where its adjacency matrix is the transpose of the word-by-word similarity matrix  $\mathbf{A}^{kl}$ .

**3.2.2 PageRank in A Tensor View.** PageRank [5] is a graph-based ranking algorithm, which is essentially a way of deciding the importance of a vertex within a graph, and recursively attracts global information from the entire graph. Formally, given a graph  $G(V, E)$ , where  $V = \{v_1, v_2, \dots, v_N\}$  is a set of nodes and  $E$  is the links between these nodes. The goal is to determine the order of these nodes that the more important node has a higher rank. The PageRank value on each node  $v_i$ , denotes as  $u_i$ , is used to indicate the importance of the node  $v_i$ . For convenience, we define  $\mathbf{A}$  as the adjacency matrix, that  $\mathbf{A}_{ij}$  denotes the  $v_i$  has a link from  $v_j$  with weight  $\mathbf{A}_{ij}$ .  $\mathbf{A}$  is also a stochastic matrix because each column sums up to 1. At the initial step all  $u_i$  have the same value  $1/N$ , denotes that all nodes are equally important. At each following step, then PageRank value  $u_i$  is updated using other nodes and links pointed to it,

$$u_i = \sum_{v_j \in V} \mathbf{A}_{ij} \cdot u_j. \quad (4)$$

After several iterations, the PageRank values  $u_i$  will converge to a set of stable values  $u_i$ , and that is the solution of PageRank.

To implement PageRank in a tensor-based computational framework, such as TensorFlow [1] or PyTorch [31], we need a tensor version of PageRank algorithm. Let  $\mathbf{u}^t = [u_1^t, u_2^t, \dots, u_n^t]$  to be a vector of length  $N$ , that obtains all nodes PageRank values at step  $t$ . Then, PageRank can be rewritten as,

$$\mathbf{u}^{t+1} = \mathbf{A}\mathbf{u}^t. \quad (5)$$

To solve the problem of isolated nodes, a stable version of PageRank is proposed [5] and adopted by our work,

$$\mathbf{u}^{t+1} = d\mathbf{A}\mathbf{u}^t + (1-d)/N \cdot \mathbb{1}, \quad (6)$$

where  $d \in [0, 1]$  is a real value to determine the ratio of the two parts, and  $\mathbb{1}$  is a vector of length  $N$  with all its values are 1. The

factor  $d$  is usually set to 0.85, and this is the value we are also using in our implementation.

In practice, the number of the iteration step  $T$  is set to a fixed value for computational efficiency. Thus,  $\mathbf{u}^t$  is the final PageRank scores for each  $v_i \in V$ , and the larger of PageRank denotes the more importance of this node in the current graph, thus we can filter out the nodes with small PageRank values.

**3.2.3 Plug PageRank in Transformer.** In this section, we propose a novel approach that plugs PageRank in the Transformer model to filter the noise at the word level. Notice that, word-level noise is composite, thus need to be estimated dynamically during the matching process. In each self-attention block, an inherent PageRank algorithm is utilized to dynamically filter the noisy words, which can reduce the sequence length layer by layer.

Standard Transformer structure, which has been selected as our base model structure, has  $L$  layers of multi-head self-attention blocks, stacked one after another, and maintains the same sequence length  $N$  at each layer. From the description in Section 3.2.1, we have known that self-attention block in Transformer can be treated as a word-by-word graph, which can be specified using an adjacency matrix  $(\mathbf{A}^{kl})^\top$  at  $k$ -th head and  $l$ -th layer in Eq 3. The word-level noise filtering process is once per layer, thus we need to average the effects of all adjacency matrices across different heads in the  $l$ -th layer,

$$\mathbf{A}^l = \frac{1}{H} \sum_{k=1}^H \mathbf{A}^{kl}. \quad (7)$$

Because  $\mathbf{A}^l$  is the output of row-wise Softmax function, each row of  $\mathbf{A}^l$  sum to 1. Thus,  $(\mathbf{A}^l)^\top$  is a stochastic matrix, which can be treated as the adjacency matrix in a graph. With above observation, we substitute  $(\mathbf{A}^l)^\top$  into Eq 5 and yield:

$$\mathbf{u}^{t+1} = d(\mathbf{A}^l)^\top \mathbf{u}^t + (1-d)/N \cdot \mathbb{I}. \quad (8)$$

Iteratively solving the equation above, we then get the PageRank values for all words/nodes in the  $(l-1)$ -th layer, denote as  $\mathbf{u}$ . Thus,  $\mathbf{u}$  represents the importance of the words in the  $(l-1)$ -th layer. After applying the attention mechanism to the words in the  $(l-1)$ -th layer, we get a list of new word representations as to the input of  $l$ -th layer. To filter noisy words, we have to estimate the importance of the words/nodes in  $l$ -th layer, which can be evaluated by redistributing the importance of the word in  $(l-1)$ -th layer under the distribution  $\mathbf{A}^l$ , thus the word importance scores are  $\mathbf{r} = \mathbf{A}^l \mathbf{u}$ . Finally, we can reduce the sequence length at  $l$ -th layer by removing the nodes which have the small values in  $\mathbf{r}$ .

In this work, we design a strategy that remove the percentage  $\alpha \in [0\%, 100\%]$  nodes per layer, so that the  $l$ -th layer has  $(\alpha)^{l-1} \cdot N$  nodes. The hyper-parameter  $\alpha$  is called a word reduction ratio. For example, let  $L = 12$ ,  $N = 400$ , if we set  $\alpha$  to 10%, the numbers of nodes at each layer are 400, 360, 324, 291, 262, 236, 212, 191, 172, 154, 139, 125.

For the BERT model, some words are too special to be removed, such as [CLS] token and [SEP] token. If the model occasionally removes these tokens during the training, it will lead to an unstable training process. It also affects the overall performance. Therefore, a token mask is designed to keep these tokens across all the layers.

**Discussions:** Many previous works [8, 15, 46] have also noticed the relation between Transformer and graph. Star-Transformer [15]

**Table 1: Description of evaluation datasets, AvgWPerD denotes average number of words per document and AvgSPerD denotes the average number of sentences per document.**

Dataset	AvgWPerD	AvgSPerD	Train	Dev	Test
CNSE	982.7	20.1	17,438	5,813	5,812
CNSS	996.6	20.4	20,102	6,701	6,700
AAN-Abs	122.7	4.9	106,592	13,324	13,324
AAN-Body	3270.1	111.6	104,371	12,818	12,696
OC	190.4	7.0	240,000	30,000	30,000
S2ORC	263.7	9.3	152,000	19,000	19,000
PAN	1569.7	47.4	17,968	2,908	2,906

adds a hub node to model the long-distance dependence and eliminates the links far from 3-term steps. TransformerXL [8] uses a segment-level recurrence with a state reuse strategy to remove all the links between words in different segments, so that can break the fixed-length limitation. Sparse-Transformer [46] explicitly eliminate links in which attention scores are lower than the threshold to make the attention matrix sparse. All of these previous works focus on eliminating links in the graph, while in this work, we focus on filtering noise words, as well as nodes, in the graph.

### 3.3 Model Training

The sentence-level filter is the heuristic approach that does not need a training process. Thus, in this section, we only consider model training for the word-level filter component.

For the model training of word-level filter, we adopt the “pre-training + fine-tuning” paradigm as in BERT. In this paradigm, the pre-trained Transformer is firstly obtained using a large unlabeled plain text in an unsupervised learning fashion. Then, the Transformer plugging PageRank at each layer is fine-tuned using the supervised downstream task. Note that word-level filters do not change the parameters in the original Transformer, due to all the parameters in the Transformer are input sequence length independent. Therefore, change the sequence length layer by layer does not affect the structure of the Transformer. Benefit from the good property of PageRank-Transformer, we can directly adopt a publicly released Transformer model, such as BERT or RoBERTa trained on a large corpus, as our pre-trained model.

In the fine-tuning step, we add the PageRank module in each self-attention layer, without introducing any additional parameters. The objective function for long-form text matching task is a binary cross-entropy loss:

$$\mathcal{L} = - \sum_i y_i \log p_i + (1 - y_i) \log(1 - p_i), \quad (9)$$

where  $p_i$  is the probability represents the matching score, generated by the representation of [CLS], and  $y_i$  is the ground-truth label.

## 4 EXPERIMENTS

In this section, we conduct experiments and in-depth analysis on three long-form text matching tasks to demonstrate the effectiveness and efficiency of our proposed model.

## 4.1 Datasets

We use seven public datasets in our experiments, and their detailed statistics are shown in Table 1.

**News Deduplication:** For this task, we use two datasets, i.e. Chinese News Same Event dataset (CNSE) and the Chinese News Same Story dataset (CNSS) released in [22]. They are constructed based on large Chinese news articles collected from major Internet news, which cover diverse topics in the open domain<sup>0</sup>. CNSE is designed to identify whether a pair of news articles report the same breaking news (or event), and CNSS is used to identify whether they belong to the same series of news stories, labeled by professional editors. The negative samples in the two datasets are not randomly generated. Document pairs that contain similar keywords are selected and exclude samples with TF-IDF similarity below a certain threshold. Finally, we follow the settings in [22] and split either dataset into training, development, and testing set with the portion of instances 6:2:2.

**Citation Recommendation:** Citation recommendations can help researchers find related works and finish paper writing more efficiently. Given the content of a research paper and a candidate citation paper, the task aims to predict whether the candidate should be cited by the paper. In our experiment, four datasets are used for this task, i.e. AAN-Abs, AAN-Body, OC, and S2ORC. AAN-Abs and AAN-Body are both constructed from the AAN dataset [33], which contains computational linguistics papers published on ACL Anthology from 2001 to 2014, along with their metadata. For the AAN-Abs dataset released in [47], each paper’s abstract and its citations’ abstracts are extracted and treated as positive pairs, and negative instances are sampled from uncited papers. For the AAN-Body dataset, we follow the same setting described in the previous work [19, 38, 43], where they remove the reference sections to prevent the leakage of ground-truth and remove the abstract sections to increase the difficulty of the task. Besides, we also use the same training, development, and testing splitting released in [38]. The OC dataset [3] contains about 7.1M papers major in computer science and neuroscience. The S2ORC dataset [24] is a large contextual citation graph of 8.1M open access papers across broad domains of science. The papers in S2ORC are divided into sections and linked by citation edges. AAN-Abs, OC, and S2ORC are pre-processed, split, and released by [47], for a fair comparison, we adopt the same settings in our experiments.

**Plagiarism Detection:** The detection of plagiarism has received considerable attention to protect the copyright of publications. It is a typical long-form text matching problem because even partial text reuse will identify plagiarism between two documents. For this task, we use the PAN dataset [32], which collects web documents with various plagiarism phenomena. Human annotations are employed to indicate the text segments that are relevant to the plagiarism both in the source and suspicious documents. Follow the settings of [45], the positive pairs are constructed by the segment of the source document and the suspicious document annotated as plagiarism, while the negative pairs are subsequently constructed by replacing the source segment in the positive pair with a segment from the corresponding source documents which is not annotated as being

plagiarised. Note that the aforementioned AAN-Abs, OC, S2ORC, and PAN datasets can be directly downloaded<sup>1</sup>.

**Evaluation Metrics:** Since all the above tasks are binary classification, we use accuracy and F1 to act as the evaluation measures, similar to [22, 47]. Specifically for each method, we perform training for 10 epochs and then choose the epoch with the best validation performance to evaluate on the test set.

## 4.2 Baselines and Experimental Settings

We adopt three types of baseline methods for comparison, including traditional term-based retrieval methods, deep learning methods for short-form text matching, and recent deep learning methods for long-form text matching.

We select three traditional term-based methods for comparison, i.e. BM25 [35], LDA [4] and SimNet [22]. BM25 is one of the most popular traditional term-based methods. The experimental results on CNSE and CNSS datasets are directly brought from [22], while others are implemented by ourselves. LDA is a famous topic model, which is used here to demote the matching method by computing similarities between the two texts represented by topic modeling vectors. We do not implement it by ourselves, and the results are directly from [22]. SimNet first extracts five text-pair similarities and conducts classification by a multi-layer neural network, whose results are brought from [22].

Considering deep learning methods for short-form text matching, we compare four types of models, including representation-based models, i.e. DSSM [17], C-DSSM [37], and ARC-I [16]; interaction-based models, i.e. ARC-II [16] and MatchPyramid [29]; hybrid models, i.e. DUET [27] and RE2 [44], and the pretraining matching model BERT-Finetuning [10]. The results of DSSM, C-DSSM, ARC-I, ARC-II, MatchPyramid, and DUET on CNSE and CNSS dataset, are directly borrowed from the previous work [22], which uses the implementations from MatchZoo [12]<sup>2</sup>. RE2 [44] is implemented using released code by the author<sup>3</sup> with the default configuration, e.g. 300-dimensions pre-trained word vectors provided by *Glove.840B* and 30-epochs training. BERT-Finetuning [10] is fine-tuned on text matching tasks based on a large-scale pre-training language model, e.g. BERT for Chinese ‘bert-base-chinese’ and BERT for English ‘bert-base-uncased’ in the Transformers library<sup>4</sup>. For each of the pretraining models, we finetune it 10 epochs on the training set.

For the methods specially designed for the long-form text matching problem, we first focus on traditional hierarchical models, e.g. HAN [45] and its variance GRU-HAN [47], GRU-HAN-CDA [47] and SMASH [20]. Then, we compare our model with a group of BERT-based hierarchical models, e.g. MatchBERT [47], BERT-HAN [47], BERT-HAN-CDA [47] and SMITH [43]. Finally, we consider some matching models by representing each text by a pre-trained model specifically designed for the long-form text, like TransformerXL [8] and Longformer [2]. The results of GRU-HAN, GRU-HAN-CDA, BERT-HAN, and BERT-HAN-CDA on AAN-Abs, OC, S2ORC, and PAN datasets are from the previous work [47].

<sup>0</sup>Datasets are available at <https://github.com/BangLiu/ArticlePairMatching>

<sup>1</sup><https://xuhuizhou.github.io/Multilevel-Text-Alignment/>

<sup>2</sup><https://github.com/NTMC-Community/MatchZoo>

<sup>3</sup><https://github.com/alibaba-edu/simple-effective-text-matching>

<sup>4</sup><https://github.com/huggingface/transformers>

**Table 2: Experimental results on the news deduplication task, e.g. CNSE and CNSS datasets. Significant performance degradation with respect to Match-Ignition is denoted as (-) with p-value  $\leq 0.05$ . We only do significant test on the models reimplemented from the source code, while the results bring from [22] do not test due to the lack of the detailed predictions.**

Model	CNSE Dataset		CNSS Dataset	
	Acc	F1	Acc	F1
I BM25 [35]	69.63	66.60	67.77	70.40
LDA [4]	63.81	62.44	67.77	70.40
SimNet [22]	71.05	69.26	70.78	74.50
ARC-I [16]	53.84	48.68	50.10	66.58
ARC-II [16]	54.37	36.77	52.00	53.83
DSSM [17]	58.08	64.68	61.09	70.58
C-DSSM [37]	60.17	48.57	52.96	56.75
MatchPyramid [29]	66.36	54.01	54.01	62.52
DUET [27]	55.63	51.94	52.33	60.67
RE2 [44]	80.59 <sup>-</sup>	78.27 <sup>-</sup>	84.84 <sup>-</sup>	85.28 <sup>-</sup>
BERT-Finetuning [9]	81.30 <sup>-</sup>	79.20 <sup>-</sup>	86.64 <sup>-</sup>	87.08 <sup>-</sup>
III CIG-Siam-GCN [22]	74.58 <sup>-</sup>	73.69 <sup>-</sup>	78.91 <sup>-</sup>	80.72 <sup>-</sup>
CIG-Sim&Siam-GCN [22]	84.64 <sup>-</sup>	82.75 <sup>-</sup>	89.77 <sup>-</sup>	90.07 <sup>-</sup>
CIG-Sim&Siam-GCN-Sim <sup>g</sup> [22]	84.21 <sup>-</sup>	82.46 <sup>-</sup>	90.03 <sup>-</sup>	90.29 <sup>-</sup>
IV Match-Ignition	<b>86.32</b>	<b>84.55</b>	<b>91.28</b>	<b>91.39</b>

The results of HAN, SMASH, MatchBERT and SMITH on the AAN-Body dataset are from the previous work [43]. TransformerXL-Finetuning and Longformer-Finetuning are implemented using the Transformers library. The pretrained model we selected are ‘*transformer-xl-wt103*’ for TransformerXL and ‘*allenai/longformer-base-4096*’ for Longformer. For each of the pre-trained models, we finetune it 10 epochs on the training set.

Specially, for CNSE and CNSS datasets, Concept Interaction Graph (CIG) model [22] is the state-of-the-art approach, which generates the representation for each vertex, and then uses a GCN to obtain the matching score. So we compare with three representative models of this approach, i.e. CIG-Siam-GCN, CIG-Sim&Siam-GCN, and CIG-Sim&Siam-GCN-Sim<sup>g</sup>. The results are obtained by implementations based on their released code <sup>5</sup>.

The hyper-parameters of our Match-Ignition model are listed below. For the sentence-level filter, the number of selected sentences per text  $\lambda$  is set to 5, the  $d$  in PageRank algorithm defined in Eq 5 is set to 0.85. For the word-level filter, we adopt a pre-trained BERT model for Chinese, e.g. ‘*bert-base-chinese*’, which contains 12 heads and 12 layers. The words filtering ratio  $\alpha$  is set to 10%, that is to say, we remove 10% words per layer. The fine-tuning optimizer is Adam [21] with the learning rate  $10^{-5}$ .  $\beta_1 = 0.9$ ,  $\beta_2 = 0.999$ ,  $\epsilon = 10^{-8}$ , and batch size is set to 8. The model is built based on the Transformers library using PyTorch [31]. The source code will be released at <https://github.com/pl8787/Match-Ignition>.

<sup>5</sup><https://github.com/BangLiu/ArticlePairMatching>

### 4.3 Experimental Results

The performance comparison results of Match-Ignition against baseline models are shown in Table 2 and Table 3. From these experimental results, we can obtain the following summaries:

1) Comparing Match-Ignition with existing deep learning models for short-form text matching, we can see that Match-Ignition outperforms all types of the existing short-form text matching models on both CNSE and CNSS datasets. Specially, it performs significantly better than two strong baselines, e.g. the hybrid model RE2 and the pretraining model BERT-Finetuning. For the tasks of citation recommendation and plagiarism, we compare the Match-Ignition model with the two strongest short-form text matching methods in the news deduplication task. The results in Table 3 II show that Match-Ignition also significantly outperforms RE2 and BERT-Finetuning, which demonstrates the superiority of Match-Ignition against existing short-form text matching models.

2) Then we Compare Match-Ignition with the current state-of-the-art methods in the three tasks, including graph-based method CIG, hierarchical methods HAN and its variants. For the news deduplication task, the Match-Ignition model significantly outperforms the state-of-the-art method CIG-Sim&Siam-GCN-Sim<sup>g</sup>, as shown in Table 2 III. That is because CIG is usually affected by noisy concept terms, while our model has the ability to filter noisy information in the learning process. For the other two tasks, as we can see in Table 3 III, Match-Ignition outperforms the state-of-the-art hierarchical methods HAN and HAN-CDA. Note that another newly proposed state-of-the-arts model SMITH-WP+SP is also an extension of the hierarchical method to tackle long-form text matching problems. However, they do not use the AAN-Abs dataset but conduct their experiments on the context-only version of the AAN dataset, namely AAN-Body. To compare with it, we apply our model to AAN-Body dataset. As we can see from the results in Table 4 II, the Match-Ignition model also outperforms the SMITH-WP+SP model.

3) Comparing the Match-Ignition model with the matching models based on pretrained long-form text representation models, e.g. TransformerXL and Longformer. The experimental results in Table 3 IV show that comparing with the finetuned version of TransformerXL and Longformer, the Match-Ignition model achieves better performances. Especially in the plagiarism detection task, noises affect the performances a lot, since the goal of TransformerXL and Longformer is to preserve the information in the long-form text as much as possible. Note that TransformerXL and Longformer only release their English versions, which are not applicable in the Chinese dataset, thus we do not list their results for the news deduplication task on CNSE and CNSS datasets.

Furthermore, we compare with another branch of the state-of-the-art hierarchical methods, e.g. SMASH and SMITH, on the AAN-Body dataset in Table 4. We do not implement them on other datasets because 1) SMASH does not provide code for model construction and training; 2) SMITH needs to pretrain on large-scale data and then finetune. So we implement our model on the dataset they used, e.g. the AAN-Body dataset, to achieve a fair comparison. The experimental results show that Match-Ignition also outperforms these baseline methods.

**Table 3: Experimental results on the citation recommendation task, e.g. AAN-Abs, OC, and S2ORC datasets and the plagiarism detection task, e.g. PAN dataset. Significant performance degradation with respect to Match-Ignition is denoted as (-) with p-value  $\leq 0.05$ . We only do significant test on the models reimplemented from the source code, while the results bring from [47] do not test due to the lack of the detailed predictions.**

		AAN-Abs Dataset		OC Dataset		S2ORC Dataset		PAN Dataset	
Model		Acc	F1	Acc	F1	Acc	F1	Acc	F1
I	BM25 [35]	67.60 <sup>-</sup>	68.00 <sup>-</sup>	80.32 <sup>-</sup>	80.38 <sup>-</sup>	76.47 <sup>-</sup>	76.53 <sup>-</sup>	61.59 <sup>-</sup>	62.47 <sup>-</sup>
II	RE2 [44]	87.81 <sup>-</sup>	88.04 <sup>-</sup>	94.53 <sup>-</sup>	94.57 <sup>-</sup>	95.27 <sup>-</sup>	95.34 <sup>-</sup>	61.97 <sup>-</sup>	58.30 <sup>-</sup>
	BERT-Finetune [9]	88.10 <sup>-</sup>	88.02 <sup>-</sup>	94.87 <sup>-</sup>	94.87 <sup>-</sup>	96.32 <sup>-</sup>	96.29 <sup>-</sup>	59.11 <sup>-</sup>	69.66 <sup>-</sup>
III	GRU-HAN [47]	68.01	67.23	84.46	82.26	82.36	83.28	75.70	75.88
	GRU-HAN-CDA [47]	75.08	75.18	89.79	89.92	91.59	91.61	75.77	76.71
	BERT-HAN [47]	73.36	73.51	86.31	86.81	90.67	90.76	87.57	87.36
	BERT-HAN-CDA [47]	82.03	82.08	90.60	90.81	91.92	92.07	86.23	86.19
IV	TransformerXL-Finetune [8]	83.85 <sup>-</sup>	83.24 <sup>-</sup>	91.61 <sup>-</sup>	91.79 <sup>-</sup>	92.50 <sup>-</sup>	92.39 <sup>-</sup>	58.25 <sup>-</sup>	69.07 <sup>-</sup>
	Longformer-Finetune [2]	88.06 <sup>-</sup>	88.41 <sup>-</sup>	94.76 <sup>-</sup>	94.74 <sup>-</sup>	96.31 <sup>-</sup>	96.29 <sup>-</sup>	56.61 <sup>-</sup>	69.74 <sup>-</sup>
V	Match-Ignition	<b>89.62</b>	<b>89.64</b>	<b>95.70</b>	<b>95.71</b>	<b>96.97</b>	<b>96.97</b>	<b>89.37</b>	<b>89.42</b>

**Table 4: Experimental results on AAN-Body dataset make a fair comparison for our Match-Ignition model and the long-form text matching models proposed in [43].**

		AAN-Body Dataset	
Model		Acc	F1
I	BM25 [35]	59.66	59.90
II	RE2 [44]	80.15	79.25
III	HAN [45]	82.19	82.57
	SMASH [19]	83.75	82.78
	MatchBERT [43]	83.55	82.93
	SMITH-WP+SP [43]	85.36	85.43
IV	Match-Ignition	<b>89.92</b>	<b>89.91</b>

**Table 5: Ablation study of the two-level noise filtering mechanisms in Match-Ignition.**

Model	CNSE		CNSS	
	Acc	F1	Acc	F1
Match-Ignition	86.32	84.55	91.28	91.39
· Sentence-level Filter Only	84.11	82.17	91.04	91.07
· Word-level Filter Only	80.31	76.39	91.10	91.18
BERT-Finetune	81.30	79.20	86.64	87.08

## 4.4 Ablation Study

Now we conduct an ablation study to investigate the two-level noise filtering strategies in Match-Ignition on news deduplication.

**4.4.1 Investigations on the sentence-level filtering.** In Table 5, the ‘Sentence-level Filter Only’ and ‘Word-level Filter Only’ model denotes the result by only using the sentence-level and word-level filter, respectively. Therefore, the sentence-level filter is critical on both CNSE and CNSS datasets: 1) without it, the accuracy on CNSE will degrade from 86.32%/91.28% to 80.31%/91.10 on CNSE and CNSS, respectively; 2) Match-Ignition with sentence-level filter

**Table 6: The impact of word reduction ratio  $\alpha$  and the execution time of these models.**

Words Reduction Ratio $\alpha$	CNSE		CNSS		Time per batch	
	Acc	F1	Acc	F1	Train	Eval
0%	84.11	82.17	91.04	91.07	1.73s	0.42s
5%	85.68	83.65	90.70	90.73	1.58s	0.37s
10%	<b>86.32</b>	<b>84.55</b>	<b>91.28</b>	<b>91.39</b>	1.33s	0.31s
20%	82.55	79.66	90.25	90.21	1.07s	0.21s

**Table 7: Comparison results with other word-level noise filtering strategies.**

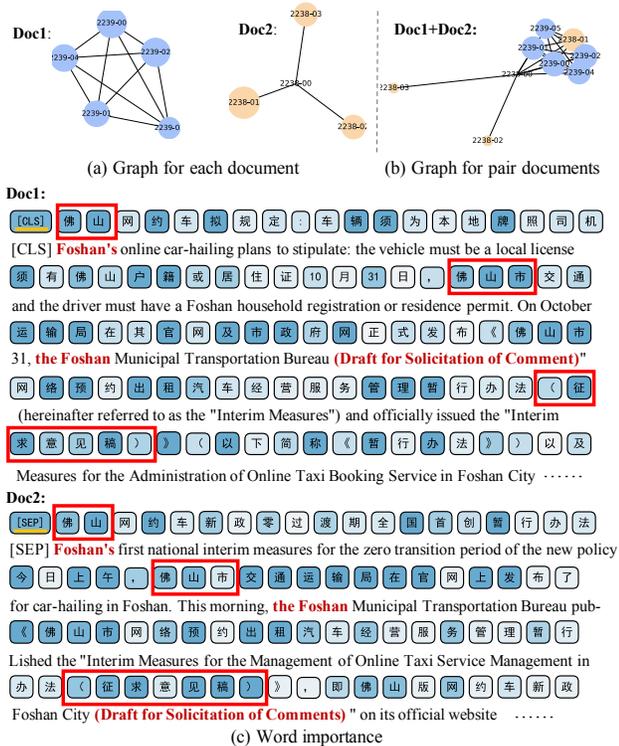
Word-level Filter	CNSE		CNSS	
	Acc	F1	Acc	F1
Random	80.38	79.19	87.68	88.20
Embedding Norm	80.54	78.15	84.92	85.03
Attention Weight	85.52	83.34	89.91	89.88
PageRank	86.32	84.55	91.28	91.39

only outperforms the strong baseline ‘BERT-Finetune’ by 3.5% and 5.1% w.r.t. accuracy on CNSE and CNSS, respectively.

## 4.5 Case Study

**4.5.1 Investigations on word-level filtering.** Still, from Table 5, we can see the effect of the word-level filtering: 1) the performance increases 2.6% and 0.3% from ‘Sentence-level Filter Only’ model to the full version Match-Ignition model on CNSE and CNSS, respectively; 2) though Match-Ignition with word-level filter only cannot beat BERT-Finetune on CNSE, it outperforms BERT-Finetune by 5.1% on CNSS. Therefore, word-level filtering also plays an important role in Match-Ignition.

Then we study the impact of word reduction ratio  $\alpha$ , which is a major hyper-parameter in the word-level filter because it determines how many words/nodes should be deleted in each layer. Specifically, we evaluate four words reduction ratio, where  $\alpha = 0\%$



**Figure 3: (a) sentence graph for each document using TextRank, (b) sentence graph built in Match-Ignition, each sentence is a node in the graph, its color represents the document it belongs to and its size represents the importance (PageRank value). (c) illustrates the word importances, and the darker color means the more important word.**

means the word-level filter is turned off. From the results shown in Table 6, we can see that too small or large a value of  $\alpha$  leads to bad performance, and  $\alpha = 10\%$  yields the best performances on both CNSE and CNSS datasets.

We also demonstrate the efficiency of the Match-Ignition model with different  $\alpha$  as in Table 6. Please note that the sentence-level filter executes very fast, comparing to evaluating Transformer. So the efficiency on sentence-level filter can be ignored. Theoretically, the major time cost in Match-Ignition is computing the word-by-word similarity matrices in self attention blocks in the Transformer model. Let  $N$  denotes the text length and  $L$  denotes the number of layers, the computation cost can be approximated by  $\text{TimeCost}(\alpha) \approx \sum_{l=0}^{L-1} (1 - \alpha)^{2l}$  where  $\text{TimeCost}(0\%) \approx 12$  and  $\text{TimeCost}(20\%) \approx 2.76$ , thus  $\alpha = 20\%$  is 4 times faster than  $\alpha = 0\%$  in theory. In our experiments, we use a single 12G Nvidia K80 GPU with batch size 8. The efficiency results in Table 6 show that  $\alpha = 20\%$  is 1.6 times faster than  $\alpha = 0\%$  at the training stage and 2 times faster at the evaluation stage.

Furthermore, we compare different types of word-level filtering strategies. *Random* stands for the method which randomly selects words at each layer, and *Embedding Norm* selects words depending on its embedding norms. *Attention Weight* uses the attention weight

to determine the importance of the word, which has been proven to be a special case of *PageRank*, i.e. without propagation on the graph. From the results in Table 7, *PageRank* achieves the best results, demonstrating the importance of word selecting strategies in word-level filtering.

To illustrate the Match-Ignition model more intuitively, we give an example from the CNSE dataset to visualize the sentence-level graph (Fig 3 (a)(b)) and word importance (Fig 3 (c)).

Specifically, Figure 3 (a) demonstrates the graph by directly applying TextRank on each document separately, and Figure 3 (b) shows the constructed sentence-level graph built-in Match-Ignition. The difference indicates the rationality of our model. For example, sentence '2238-01' are equally important as '2238-02' and '2238-03' in Figure 3 (a). While it becomes much more important in Figure 3 (b) because it has more connections with the sentences in Doc1. Therefore, our model is capable to capture the key sentences in the matching process, by considering connections both inside and between two documents.

Fig 3 (c) show the word importance in different colors, where darker color indicates a higher importance score. Specifically, the importance score is computed based on the number of layers retaining the word, which shows the importance of each word in the whole matching process. The results are accordant with human understanding. For example, the location 'the Foshan' and the name of the policy 'Draft for Solicitation of Comment' is important for determining the matching degree of the news, which indeed obtain a higher importance score in the model, as highlighted with rectangles. Furthermore, the results show that special tokens like [CLS] and [SEP] are also important for long-form text matching. That is because [CLS] token acts as the global information aggregator and [SEP] token acts as a separator of the two texts, which are two crucial indicators in long-form text matching.

## 5 CONCLUSION

In this paper, we propose a novel hierarchical noise filtering approach for the long-form text matching problem. The novelty lies in the employment of the well-known PageRank algorithm to identify and filter both sentence-level and word-level noisy information, which can be viewed as a generalized version of using attention weight with propagation on the graph. We conduct extensive experiments on three typical long-form text matching tasks including seven public datasets, and the results show that our proposed model significantly outperforms both short-form text matching models and recent state-of-the-arts long-form text matching models.

In the future, we plan to investigate how to jointly learn the sentence-level and word-level noise filter in Match-Ignition. In addition, we would like to study the relation between Match-Ignition and graph neural network, and whether there exists a graph neural network-based model to achieve the two-level noise filtering in long-form text matching.

## ACKNOWLEDGMENTS

This work was supported by National Natural Science Foundation of China (NSFC) under Grants No. 61906180, No. 61773362 and No. 91746301, National Key R&D Program of China under Grants 2020AAA0105200.

## REFERENCES

- [1] Martín Abadi, Paul Barham, Jianmin Chen, Zhifeng Chen, Andy Davis, Jeffrey Dean, Matthieu Devin, Sanjay Ghemawat, Geoffrey Irving, Michael Isard, et al. 2016. Tensorflow: A system for large-scale machine learning. In *12th {USENIX} Symposium on Operating Systems Design and Implementation*. 265–283.
- [2] Iz Beltagy, Matthew E. Peters, and Arman Cohan. 2020. Longformer: The Long-Document Transformer. arXiv:2004.05150 [cs.CL]
- [3] Chandra Bhagavatula, Sergey Feldman, Russell Power, and Waleed Ammar. 2018. Content-Based Citation Recommendation. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*. 238–251.
- [4] David M Blei, Andrew Y Ng, and Michael I Jordan. 2003. Latent dirichlet allocation. *Journal of machine Learning research* 3, Jan (2003), 993–1022.
- [5] Sergey Brin and Lawrence Page. 1998. The anatomy of a large-scale hypertextual Web search engine. *Computer Networks and ISDN Systems* 30, 1-7 (1998), 107–117.
- [6] Gino Brunner, Yang Liu, Damian Pascual, Oliver Richter, Massimiliano Ciaramita, and Roger Wattenhofer. 2019. On Identifiability in Transformers. In *International Conference on Learning Representations*.
- [7] Sumit Chopra, Raia Hadsell, Yann LeCun, et al. 2005. Learning a similarity metric discriminatively, with application to face verification. In *CVPR (1)* (Boston, Massachusetts). 539–546.
- [8] Zihang Dai, Zhilin Yang, Yiming Yang, Jaime G Carbonell, Quoc Le, and Ruslan Salakhutdinov. 2019. Transformer-XL: Attentive Language Models beyond a Fixed-Length Context. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*. 2978–2988.
- [9] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2018. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805* (2018).
- [10] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics*. 4171–4186.
- [11] Yixing Fan, Jiafeng Guo, Yanyan Lan, Jun Xu, Chengxiang Zhai, and Xueqi Cheng. 2018. Modeling Diverse Relevance Patterns in Ad-Hoc Retrieval. In *The 41st International ACM SIGIR Conference on Research Development in Information Retrieval* (Ann Arbor, MI, USA) (SIGIR '18). Association for Computing Machinery, New York, NY, USA, 375–384.
- [12] Yixing Fan, Liang Pang, JianPeng Hou, Jiafeng Guo, Yanyan Lan, and Xueqi Cheng. 2017. Matchzo: A toolkit for deep text matching. *arXiv preprint arXiv:1707.07270* (2017).
- [13] Jiafeng Guo, Yixing Fan, Qingyao Ai, and W Bruce Croft. 2016. A deep relevance matching model for ad-hoc retrieval. In *Proceedings of the 25th ACM CIKM*. ACM, 55–64.
- [14] Jiafeng Guo, Yixing Fan, Liang Pang, Liu Yang, Qingyao Ai, Hamed Zamani, Chen Wu, W Bruce Croft, and Xueqi Cheng. 2019. A deep look into neural ranking models for information retrieval. *Information Processing & Management* (2019), 102067.
- [15] Qipeng Guo, Xipeng Qiu, Pengfei Liu, Yunfan Shao, Xiangyang Xue, and Zheng Zhang. 2019. Star-Transformer. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*. 1315–1325.
- [16] Baotian Hu, Zhengdong Lu, Hang Li, and Qingcai Chen. 2014. Convolutional neural network architectures for matching natural language sentences. In *Advances in neural information processing systems*. 2042–2050.
- [17] Po-Sen Huang, Xiaodong He, Jianfeng Gao, Li Deng, Alex Acero, and Larry Heck. 2013. Learning deep structured semantic models for web search using clickthrough data. In *Proceedings of the 22nd ACM international conference on Information & Knowledge Management*. ACM, 2333–2338.
- [18] Kai Hui, Andrew Yates, Klaus Berberich, and Gerard de Melo. 2017. PACRR: A Position-Aware Neural IR Model for Relevance Matching. In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*. 1049–1058.
- [19] Jyun-Yu Jiang, Mingyang Zhang, Cheng Li, Michael Bendersky, Nadav Golbandi, and Marc Najork. 2019. Semantic Text Matching for Long-Form Documents. In *The World Wide Web Conference* (San Francisco, CA, USA) (WWW '19). Association for Computing Machinery, New York, NY, USA, 795–806.
- [20] Ray Jiang, Sven Gowal, Timothy A Mann, and Danilo J Rezende. 2018. Beyond greedy ranking: Slate optimization via List-CVAE. *arXiv preprint arXiv:1803.01682* (2018).
- [21] Diederik P Kingma and Jimmy Ba. 2014. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980* (2014).
- [22] Bang Liu, Di Niu, Haojie Wei, Jinghong Lin, Yan Cheng He, Kunfeng Lai, and Yu Xu. 2019. Matching Article Pairs with Graphical Decomposition and Convolutions. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*. Association for Computational Linguistics, Florence, Italy, 6284–6294.
- [23] Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. Roberta: A robustly optimized bert pretraining approach. *arXiv preprint arXiv:1907.11692* (2019).
- [24] Kyle Lo, Lucy Lu Wang, Mark Neumann, Rodney Kinney, and Daniel S Weld. 2020. S2ORC: The Semantic Scholar Open Research Corpus. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*. 4969–4983.
- [25] Zhengdong Lu and Hang Li. 2013. A deep architecture for matching short texts. In *Advances in Neural Information Processing Systems*. 1367–1375.
- [26] Rada Mihalcea and Paul Tarau. 2004. TextRank: Bringing order into text. In *Proceedings of the 2004 conference on empirical methods in natural language processing*. 404–411.
- [27] Bhaskar Mitra, Fernando Diaz, and Nick Craswell. 2017. Learning to match using local and distributed representations of text for web search. In *Proceedings of the 26th International Conference on World Wide Web*. International World Wide Web Conferences Steering Committee, 1291–1299.
- [28] Hamid Palangi, Li Deng, Yelong Shen, Jianfeng Gao, Xiaodong He, Jianshu Chen, Kinyong Song, and Rabab Ward. 2016. Deep sentence embedding using long short-term memory networks: Analysis and application to information retrieval. *IEEE/ACM Transactions on Audio, Speech and Language Processing (TASLP)* 24, 4 (2016), 694–707.
- [29] Liang Pang, Yanyan Lan, Jiafeng Guo, Jun Xu, Shengxian Wan, and Xueqi Cheng. 2016. Text matching as image recognition. In *Thirtieth AAAI Conference on Artificial Intelligence*.
- [30] Liang Pang, Yanyan Lan, Jiafeng Guo, Jun Xu, Jingfang Xu, and Xueqi Cheng. 2017. DeepRank: A new deep architecture for relevance ranking in information retrieval. In *Proceedings of the 2017 ACM CIKM*. ACM, 257–266.
- [31] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. 2019. PyTorch: An imperative style, high-performance deep learning library. In *Advances in Neural Information Processing Systems*. 8024–8035.
- [32] Martin Potthast, Matthias Hagen, Tim Gollub, Martin Tippmann, Johannes Kiesel, Paolo Rosso, Efstathios Stamatatos, and Benno Stein. 2013. Overview of the 5th international competition on plagiarism detection. In *CLEF Conference on Multilingual and Multimodal Information Access Evaluation*. CELCT, 301–331.
- [33] Dragomir R Radev, Pradeep Muthukrishnan, Vahed Qazvinian, and Amjad Abu-Jbara. 2013. The ACL anthology network corpus. *Language Resources and Evaluation* 47, 4 (2013), 919–944.
- [34] Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. [n.d.]. Language models are unsupervised multitask learners. ([n. d.]).
- [35] Stephen Robertson and Hugo Zaragoza. 2009. The Probabilistic Relevance Framework: BM25 and Beyond. *Information Retrieval* 3, 4 (2009), 333–389.
- [36] Franco Scarselli, Marco Gori, Ah Chung Tsoi, Markus Hagenbuchner, and Gabriele Monfardini. 2008. The graph neural network model. *IEEE Transactions on Neural Networks* 20, 1 (2008), 61–80.
- [37] Yelong Shen, Xiaodong He, Jianfeng Gao, Li Deng, and Grégoire Mesnil. 2014. A latent semantic model with convolutional-pooling structure for information retrieval. In *Proceedings of the 23rd ACM international conference on conference on information and knowledge management*. ACM, 101–110.
- [38] Yi Tay, Mostafa Dehghani, Samira Abnar, Yikang Shen, Dara Bahri, Philip Pham, Jinfeng Rao, Liu Yang, Sebastian Ruder, and Donald Metzler. 2020. Long Range Arena: A Benchmark for Efficient Transformers. arXiv:2011.04006 [cs.LG]
- [39] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In *NIPS*. 5998–6008.
- [40] Shengxian Wan, Yanyan Lan, Jiafeng Guo, Jun Xu, Liang Pang, and Xueqi Cheng. 2016. A deep architecture for semantic matching with multiple positional sentence representations. In *Thirtieth AAAI Conference on Artificial Intelligence*.
- [41] Shengxian Wan, Yanyan Lan, Jun Xu, Jiafeng Guo, Liang Pang, and Xueqi Cheng. 2016. Match-srnn: Modeling the recursive matching structure with spatial rnn. *arXiv preprint arXiv:1604.04378* (2016).
- [42] Zhiguo Wang, Wael Hamza, and Radu Florian. 2017. Bilateral multi-perspective matching for natural language sentences. *arXiv preprint arXiv:1702.03814* (2017).
- [43] Liu Yang, Mingyang Zhang, Cheng Li, Michael Bendersky, and Marc Najork. 2020. Beyond 512 Tokens: Siamese Multi-depth Transformer-based Hierarchical Encoder for Document Matching. arXiv:2004.12297 [cs.IR]
- [44] Runqi Yang, Jianhai Zhang, Xing Gao, Feng Ji, and Haiqing Chen. 2019. Simple and Effective Text Matching with Richer Alignment Features. *arXiv preprint arXiv:1908.00300* (2019).
- [45] Zichao Yang, Diyi Yang, Chris Dyer, Xiaodong He, Alex Smola, and Eduard Hovy. 2016. Hierarchical Attention Networks for Document Classification. In *Proceedings of the 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*. Association for Computational Linguistics, San Diego, California, 1480–1489.
- [46] Guangxiang Zhao, Junyang Lin, Zhiyuan Zhang, Xuancheng Ren, and Xu Sun. 2019. Sparse Transformer: Concentrated Attention Through Explicit Selection. (2019).
- [47] Xuhui Zhou, Nikolaos Pappas, and Noah A. Smith. 2020. Multilevel Text Alignment with Cross-Document Attention. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*. Association for Computational Linguistics, Online, 5012–5025.