# D-HAN: Dynamic News Recommendation with Hierarchical Attention Network

Qinghua Zhao, Xu Chen, Hui Zhang, Shuai Ma

**Abstract**—News recommendation is an effective information dissemination solution in modern society. While recent years have witnessed many promising news recommendation models, they mostly capture the user-news interactions on the document-level in a static manner. However, in real-world scenarios, the news can be quite complex and diverse, blindly squeezing all the contents into an embedding vector can be less effective in extracting information compatible with the personalized preference of the users. In addition, user preferences in the news recommendation scenario can be highly dynamic, and a tailored dynamic mechanism should be designed for better recommendation performance. In this paper, we propose a novel dynamic news recommender model. For better understanding the news content, we leverage the attention mechanism to represent the news from the sentence-, element- and document-levels, respectively. For capturing users' dynamic preferences, the continuous time information is seamlessly incorporated into the computing of the attention weights. More specifically, we design a hierarchical attention network, where the lower layer learns the importance of different sentences and elements, and the upper layer captures the correlations between the previously interacted and the target news. To comprehensively model the dynamic characters, we firstly enhance the traditional attention mechanism by incorporating both absolute and relative time information, and then we propose a dynamic negative sampling method to optimize the users' implicit feedback. We conduct extensive experiments based on three real-world datasets to demonstrate our model's effectiveness. Our source code and pre-trained representations are available at https://github.com/lshowway/D-HAN.

**Index Terms**—News Recommendation, Hierarchical Attention Network, Dynamic Negative Sampling

---------------- ✦ ----------------

## 1 INTRODUCTION

The ever-prospering of Internet technologies has gradually shifted how people receive information. Online news applications have rapidly replaced traditional printed media, which collect content from multiple publishers and receive a considerable volume of news articles. Despite many advantages, these applications usually have to face the information overloading problem, which motivates the development of news recommendation systems.

In the field of news recommendation, a key observation is that people's news-reading behaviors are not independent [7]. Previously interacted news has a substantial impact on the following reading choice. Along this line, a number of news recommendation models have been built [15], [16], [27], which capture people's sequential reading patterns. The contents interacted previously usually have different impacts on choosing the next one to recommend in practice. Specifically, news articles are composed of several sentences, and different sentences of people's previously interacted news have other impacts on their subsequent actions. Second, news articles have basic components called news elements, which are known as five W and one H (5W1H), i.e.,

- *Xu Chen and Shuai Ma are corresponding authors.*
- *Qinghua Zhao, Hui Zhang and Shuai Ma are with Beijing Advanced Innovation Center for Big Data and Brain Computing, Beihang University, Beijing, China. E-mail: {zhaoqh, zhangh17, mashuai}@buaa.edu.cn.*
- *Xu Chen is with Beijing Key Laboratory of Big Data Management and Analysis Methods. Gaoling School of Artificial Intelligence, Renmin University of China, Beijing 100872, China. E-mail: xu.chen@ruc.edu.cn.*

who, when, where, what, why and how [18]. The 5W1H elements clearly describe the critical information of news explicitly. The six elements are the basic principle of news writing, generally followed by the world press. For example, for a news article that a user read before, the first 4W elements are "Warriors, Cavaliers", "June 1-9, 2018", "Cleveland, Oakland" and "NBA finals", respectively. Third, think of each news as a whole, news articles in people's reading logs have different influences on their decisions about candidate news, i.e., to read or not. Continue the above example, and the user will read sports news primarily because of the NBA news rather than economic news in his reading logs. Besides, the current period and the time difference from the last clicked news significantly impact the following news. For example, users may browse news about the stock market at around 10 am on weekdays, while constellations news at midnight. Finally, treating users' news-reading behaviors as a sequence, since the order of this sequence or the interactions between the sequence hides much information. Therefore, given such various information of news reading sequences, it is desired to use them from sentence-, element-, document- and sequence-level. Also, it is necessary to incorporate dynamism in news recommendations.

Motivated by the above observations, in this paper, we propose **D-HAN**, a Dynamic news recommendation model based on Hierarchical Attention Network. The main building block of our multi-granularity model is a two-layer attention network. Specifically, the lower layer distinguishes sentence-level and element-level impacts. It automatically determines the attention weights between sen-

tences and elements. The upper layer discriminates the various correlations in document-level. It determines the attention weights between history-candidate news pairs. To comprehensively model the dynamic characters, we further incorporate news clicked timestamps and relative time intervals in the document-level layer. The history summarization layer consists of several Transformer encoders [31] for efficiently learning sequence-level information from users' history news-reading. The upper left dynamic negative sampling layer selects negative samples according to the output of document-level attention layer, and D-HAN model parameters dominate this process. The upper right prediction layer consists of a fully-connected network and takes the negative and positive samples as input.

**Contributions.** In summary, the contributions can be concluded as follows: (1) We propose to simultaneously capture different granular information, i.e., sentence-, element-, document- and sequence-level information for news recommendation. (2) We propose recommending news dynamically by a time-aware document-level attention layer, which incorporates the absolute and relative time information. (3) We propose incorporating negative sampling into the training process to optimize the model.

This study extends our earlier work [39] as follows. (1) On sentence- and element-level attention layers, we replace both the original additive attention and the cosine similarity with the scaled dot-product attention [31]. We replace the original additive attention and convolutional neural network (CNN) with a Transformer encoder on the document-level attention and history summarization layers. (2) We design a novel dynamic negative sampling layer such that the selected negative samples are more informative. (3) We replace the hard time-decaying factor with the absolute and relative time embedding on the time utilization method, which models the dynamic characteristics more comprehensively. (4) On the experiments, we compare the performance of each attention layer and add more experiments on the performance of the history summarization layer, the relative and absolute time information, and dynamic negative sampling.

## 2 RELATED WORK

News recommendation has previously attracted much attention, aiming to provide personalized news articles for users. Traditional news recommendation methods can be divided into content-based, collaborative filtering, and hybrid. Content-based methods recommend news solely based on content similarity [14], [24]. Collaborative filtering methods utilize users' feedback to news articles to make recommendation [5], suffering from serious cold-start problems. And hybrid methods combine the two strategies to achieve better recommendation performance [21], [23].

Recently, the neural recommendation has shown its superior performance. GRU4Rec [13] and its variants [19] and GRU4Rec++ [12] apply RNN to session-based recommendation. Caser uses horizontal and vertical convolution filters to capture sequential patterns [30]. Based on Caser, [37] further models the long-range dependence in the sequence.

[40] proposes a reinforcement learning framework, aiming at online news recommendation. Attention mechanism has shown effective results in machine translation [2], image captioning [36] and so on. It has also shown surprising potential in the field of recommendation. [32] enhances news recommendation with knowledge graphs, applying an attention network to get users' representations. [22] designs a deep fusion model which leverages various levels of interaction by inception module and merges information from different channels by attention mechanism. Deep Interest Network [41] designs an attention unit for learning the representations of users adaptively. [35] improves factorization machines by discriminating the importance of feature interactions via an attention network. Attention Collaborative Filtering model [4] introduces an attention mechanism into collaborative filtering to model item- and component-level implicit feedback in the multimedia recommendation.

With the great success of Bert [6] in other NLP tasks, self-attention mechanism or Transformer architecture has gradually been used for news recommendations. [8] applies a Transformer architecture with multi-head self-attention to obtain news content representation from news titles and topics. BERT4Rec [29] concatenates all history news articles into a document as the input of Transformer encoder. NRMS [34] proposes a news encoder and a user encoder, where multi-head self-attention is used to learn representations.

Dynamic recommendation considers the time-dependent effect in recommendation, which can help us to understand users' behaviors. For instance, [38] observes that people tend to visit different locations at different time in a day and utilize the absolute time for location recommendation. With time gates, Time-LSTM [42] can not only model sequential information, but also capture well the time interval information. To model the dynamics of sequence recommendation, TiSASRec [20] replaces the position encoding of Transformer with the absolute position of items and the time intervals between any two items in a news sequence.

## 3 OVERVIEW

### 3.1 Problem Definition

Assume there is an online news platform offering news services to users. Once the platform receives a new piece of news, it estimates the click rate for each user based on the news articles the user had read earlier. Formally, let $\mathcal{C}_i = [c_1, c_2, ..., c_L]$ and $\mathcal{A}_i = [a_1, a_2, ..., a_L]$ denote the sequence of the most recent $L$ pieces of news read by user $i$ and the corresponding click time, where $L$ is the number of news articles that we consider to estimate the click rate. Each piece of news $c_j$ consists of a sequence of sentences, i.e., $[s_{j1}, s_{j2}, ..., s_{jK}]$, where $s_{jk}$ is the $k$-th sentence in $c_j$ and $K$ is the maximum number of sentences we consider. Each piece of news $c_j$ is also represented by a set of news elements, detailed in the following subsection. Given the news sequence $\mathcal{C}_i$ and candidate news $c^*$, we aim to predict the click rate of $c^*$ by user $i$.
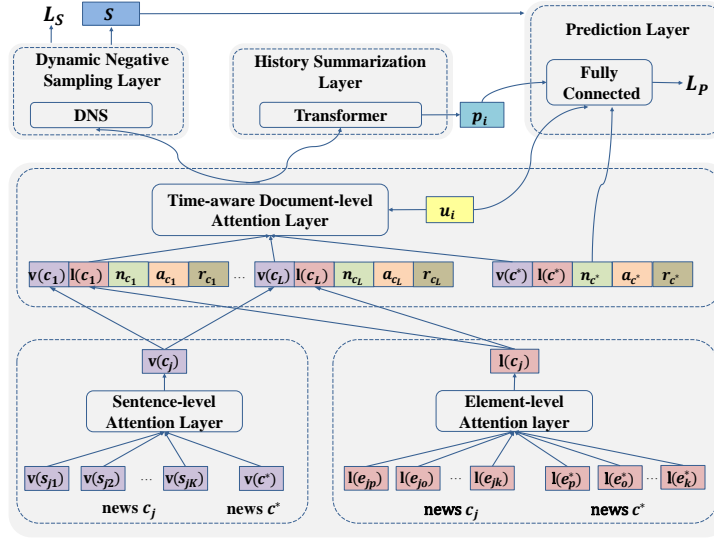
Fig. 1. The architecture of D-HAN model. The left bottom and right bottom are sentence-level attention layer and element-level attention layer, respectively. The middle is time-ware document-level attention layer, and the upper are dynamic negative sampling layer, history summarization layer and prediction layer, respectively.

## 3.2 News Elements

Extracting 5W1H elements is an intractable natural language processing (NLP) problem and only few works extract them from news contents [10], [33]. However, we can not employ them due to the language difference for [10] and the lack of specific news corpus for [33]. This part of NLP is not the focus of our model, so we define news elements that can be easily extracted by the named entity recognition and keywords extraction modules of NLP tools. Specifically, they are <u>person</u>, <u>organization</u>, <u>time</u>, <u>location</u> and <u>keywords</u>, corresponding to <u>who</u>, <u>who</u>, <u>when</u>, <u>where</u> and <u>what</u> elements of 5W1H, respectively. With news elements we define, each news can be summed up into a sentence, i.e., persons or organizations do something at a specific time and location. Formally, each piece of news $c_j$ is represented by a set of elements, i.e., $\{e_{jp}, e_{jo}, e_{jt}, e_{jl}, e_{jk}\}$, corresponding to the above elements we define, respectively.

## 3.3 D-HAN Framework

Figure 1 illustrates the architecture of our proposed D-HAN model. It has four main components: the core hierarchical attention layer, i.e., sentence-level, element-level, and time-aware document-level attention layer, the history summarization layer, the dynamic negative sampling layer, and the prediction layer. Our multi-granularity model takes news sequence $\mathcal{C}_i$ of user $i$, candidate news $c^*$ and their corresponding click time $\mathcal{A}_i$, $\mathcal{A}^*$ as input. Sentence-level attention layer first computes weights of interactions between sentences of $c_j$ and $c^*$ and gets the content vector $\mathbf{v}(c_j)$ of news $c_j$ and $\mathbf{v}(c^*)$ of news $c^*$. Second, element-level attention layer computes weights of interactions between elements of $c_j$ and $c^*$ and gets the element vector $\mathbf{l}(c_j)$ and $\mathbf{l}(c^*)$. Then time-aware document-level attention layer takes content vector, element vector, user structural embedding $\mathbf{u}_i$

and click time embedding as input, and obtains candidate-dependent representation. Next, these candidate-dependent representations are fed into dynamic negative sampling and history summarization layers simultaneously, based on which we obtain the negative samples $S$ and sequence vector $\mathbf{p}_i$ for user $i$, respectively. Finally, the prediction layer inputs the output of the other components to compute the click rate of $c^*$ by user $i$.

## 4 D-HAN MODEL

### 4.1 Sentence-level Attention Layer

When predicting candidate news $c^*$, sentences of news $c_j$ unequally affect user $i$'s choice, i.e., to read or not. Intuitively, sentences content-relevant to news $c^*$ have more significant impacts on reading. The sentence-level attention layer aims to discriminate various influences of sentences.

We first need content vectors of sentences of news $c_j$ and content vector of candidate news $c^*$. There exist many sentence embedding methods. We adopt Paragraph Vector [17] due to its consideration of the ordering and semantics of the words in sentences. The content vector $\mathbf{v}(s_{jk})$ of sentence $s_{jk}$ is embedded in a $d$-dimensional space $\mathbb{R}^d$, and the content vector $\mathbf{v}(c_j) \in \mathbb{R}^{K \times d}$ of news $c_j$ is stacked by $\mathbf{v}(s_{jk})$, and the content vector $\mathbf{v}(c^*) \in \mathbb{R}^d$ of candidate news $c^*$ is calculated by averaging the sentences vectors.

To model the interactions among news $c_j$ and candidate news $c^*$, we concatenate $\mathbf{u}_i$, $\mathbf{v}(c_j)$ and $\mathbf{v}(c^*)$ along the first dimension as $[\mathbf{u}_i\mathbf{v}(c_j)\mathbf{v}(c^*)] \in \mathbb{R}^{(K+2) \times d}$, and adopt scaled dot-product attention [31] to learn its representation. The self attention weights among $c_j$ and $c^*$ are:

$$b_j = \frac{[\mathbf{u}_i\mathbf{v}(c_j)\mathbf{v}(c^*)]\mathbf{W}_1 \cdot (\mathbf{W}_2[\mathbf{u}_i\mathbf{v}(c_j)\mathbf{v}(c^*)]^T)}{\sqrt{d}}, \quad (1)$$

where $\mathbf{W}_1 \in \mathbb{R}^{d \times d}$, $\mathbf{W}_2 \in \mathbb{R}^{d \times d}$ are the parameters of model, $\cdot$ denotes the dot product operation, $\sqrt{d}$ is a scale factor to prevent dot products growing large in magnitude, $b_j \in \mathbb{R}^{(K+2) \times (K+2)}$. Note that we omit the bias for brevity here. These self attention weights are further normalized by softmax function $\beta_{jk} = \frac{\exp(b_{jk})}{\sum_{j'=1}^{K} \exp(b_{jk'})}$. Here attention weight $\beta_{jk}$ can be interpreted as the content relevance among sentence $s_{jk}$ and candidate news $c^*$, user embedding $\mathbf{u}_i$. With these weights, we compute the content vector $\mathbf{v}(c_j)$ of news $c_j$ with respect to candidate news $c^*$ as the dot product between $\beta_j$ and linearly transformed $[\mathbf{u}_i \mathbf{v}(c_j) \mathbf{v}(c^*)]$ by model parameter $\mathbf{W}_3 \in \mathbb{R}^{\mathbf{d} \times \mathbf{d}}$, that is, $\mathbf{v}(c_j) = \beta_j \cdot [\mathbf{u}_i \mathbf{v}(c_j) \mathbf{v}(c^*)] \mathbf{W}_3 \in \mathbb{R}^{(K+2) \times d}$.

## 4.2 Element-level Attention Layer

Given a pair of news $c_j$ and candidate news $c^*$, different elements play a different role in users' decisions, and the goal of the element-level attention layer is to discriminate various impacts of other elements.

With the named entity recognition and keywords extraction modules of NLP tools, we can extract elements we define, i.e., underline{person}, underline{organization}, underline{time}, underline{location} and underline{keywords}, for news $c_j$ and candidate news $c^*$. Each element is extracted in the form of one or more words. In this paper, each word is embedded in a $d$-dimensional space $\mathbb{R}^d$ by Word2vec [25] which is successful in capturing semantics relatedness. For instance, let $\mathbf{l}'c_j = \{\mathbf{l}'e_{jp}, \mathbf{l}'e_{jo}, \mathbf{l}'e_{jt}, \mathbf{l}'e_{jl}, \mathbf{l}'e_{jk}\} \in \mathbb{R}^{5 \times d}$, and $\mathbf{l}'e_{jp} \in \mathbb{R}^d$ is the vector of element $p$ of news $c_j$, which is obtained by averaging the vectors of words that represent $e_{jp}$. We concatenate the corresponding element vector of news $c_j$ and $c^*$ along $d$-dimension as $[\mathbf{l}'c_j \mathbf{l}'c^*] \in \mathbb{R}^{5 \times 2d}$, and then to smooth their internal differences, we apply three one-layer feed-forward neural networks to transform them into $d$-dimension space:

$$q_{c_j} = [\mathbf{l}'c_j \mathbf{l}'c^*]\mathbf{W}_4,$$
$$k_{c_j} = [\mathbf{l}'c_j \mathbf{l}'c^*]\mathbf{W}_5, \qquad (2)$$
$$v_{c_j} = [\mathbf{l}'c_j \mathbf{l}'c^*]\mathbf{W}_6,$$

where $\mathbf{W}_4 \in \mathbb{R}^{2d \times d}$, $\mathbf{W}_5 \in \mathbb{R}^{2d \times d}$ and $\mathbf{W}_6 \in \mathbb{R}^{2d \times d}$ are model parameters, $q_{c_j}$, $k_{c_j}$ and $v_{c_j}$ represent different linear transformations of $[\mathbf{l}'c_j \mathbf{l}'c^*]$, which will be used to compute self-attended vector representation.

To model the internal and external interactions among element vector of $c_j$ and $c^*$, we take dot product between $q_{c_j}$ and $k_{c_j}$ to get the raw self-attention scores:

$$\gamma_j = \frac{q_{c_j} \cdot k_{c_j}{}^T}{\sqrt{d}} \in \mathbb{R}^{5 \times 5}, \qquad (3)$$

where $\sqrt{d}$ is a scale factor, and the raw self-attention scores are normalized by softmax function to probabilities followed by a dropout operation on entire elements to attend to.

Here self-attention scores $\gamma_j$ can be interpreted as the relevance of elements under the influence of both historical news $c_j$ and candidate news $c^*$. With these weights, we compute the element vector $\mathbf{l}(c_j)$ of news $c_j$ as dot product between $\gamma_j$ and $v_{c_j}$: $\mathbf{l}(c_j) = \gamma_j \cdot v_{c_j} \in \mathbb{R}^{5 \times d}$.

## 4.3 Time-aware Document-level Attention Layer

Given news-reading sequence $\mathcal{C}_i$ of user $i$, news articles unequally influence whether he reads candidate news $c^*$ or not, and different periods are also crucial for making decisions. Two pieces of news frequently co-clicked by people tend to be similar. To preserve the structural information, we learn a news id embedding in a $d$-dimensional space for each news according to its id, i.e., $\mathbf{n}_{c_j} \in \mathbb{R}^d$ for news $c_j$. News id embeddings are randomly initialized and automatically learned in the training phase. Similarly, two users with similar history news sequences tend to be similar. Users' structural information partly reflects user preferences [4]. Therefore, we learn an embedding in a $d$-dimensional space for each user, i.e., $\mathbf{u}_i \in \mathbb{R}^d$ for user $i$. User embeddings are also randomly initialized and updated during training. For news $c_j$, we concatenate its sentence, element and structural embedding along this $d$-dimension to obtain representation $\mathbf{x}'_{\mathbf{c_j}} = [\,\mathbf{v}(c_j)\,\mathbf{l}(c_j)\,\mathbf{n}_{c_j}] \in \mathbb{R}^{3d}$. For candidate news $c^*$, this operation leads to $\mathbf{x}^* = [\,\mathbf{v}(c^*)\,\mathbf{l}(c^*)\,\mathbf{n}_{c^*}] \in \mathbb{R}^{3d}$.

An observation is that users tend to read different news at different periods with a distinct time difference from the last clicked time. To model this dynamics of news recommendation, we propose to model news content information and timestamp information simultaneously. Specifically, we adopt the year, month, week, day, hour, minutes as the absolute time and the absolute time difference between news $c_j$ and candidate news $c^*$ as the relative time interval. To preserve the time information, we learn $d$-dimensional absolute time embedding and time interval embedding for news sequence $\mathcal{C}_i$, i.e., $\mathbf{a}_{\mathcal{C}_i} \in \mathbb{R}^{(L+1) \times d}$ for absolute time and $\mathbf{r}_{\mathcal{C}_i} \in \mathbb{R}^{L \times d}$ for relative time interval.

• **Embed with $\mathbf{r}_{\mathcal{C}_i}$.** Given a news sequence $\mathcal{C}_i$, its representation $\mathbf{x}'_{\mathcal{C}_i}$ is concatenated with its corresponding relative time interval embedding $\mathbf{r}_{\mathcal{C}_i}$ and candidate news representation $\mathbf{x}^*$ as $\mathbf{z}_{\mathcal{C}_i} = [\mathbf{x}'_{\mathcal{C}_i} \mathbf{r}_{\mathcal{C}_i} \mathbf{x}^*] \in \mathbb{R}^{L \times 7d}$, where $\mathbf{x}'_{\mathcal{C}_i} = \{\mathbf{x}'_{c_1}, \mathbf{x}'_{c_2}, ..., \mathbf{x}'_{c_L}\} \in \mathbb{R}^{L \times 3d}$. The time-aware representation of candidate news is $\mathbf{z}^* = \mathbf{x}^*$.

• **Embed with $\mathbf{a}_{\mathcal{C}_i}$.** History news representation and candidate news representation are concatenated with their corresponding absolute time embedding as $\mathbf{z}_{\mathcal{C}_i} = [\mathbf{x}'_{\mathcal{C}_i} \mathbf{a}_{\mathcal{C}_i}^{1:L}] \in \mathbb{R}^{L \times 4d}$, $\mathbf{z}^* = [\mathbf{x}^* \mathbf{a}_{\mathcal{C}_i}^{L+1}] \in \mathbb{R}^{4d}$, respectively.

• **Embed with both $\mathbf{r}_{\mathcal{C}_i}$ and $\mathbf{a}_{\mathcal{C}_i}$.** First, to use the absolute timestamp and relative time interval, we concatenate $\mathbf{x}'_{\mathcal{C}_i}$, $\mathbf{a}_{\mathcal{C}_i}^{1:L}$ and $\mathbf{r}_{\mathcal{C}_i}$ as $\mathbf{z}_{\mathcal{C}_i} = [\mathbf{x}'_{\mathcal{C}_i} \mathbf{a}_{\mathcal{C}_i}^{1:L} \mathbf{r}_{\mathcal{C}_i}] \in \mathbb{R}^{L \times 5d}$, Then, to embed time embedding into candidate news representation, we concatenate the representation $\mathbf{x}^*$ of candidate news $c^*$ and the representation $\mathbf{a}_{\mathcal{C}_i}^{L+1}$ of absolute time of candidate news $c^*$ as $\mathbf{z}^* = [\mathbf{x}^* \mathbf{a}_{\mathcal{C}_i}^{L+1}] \in \mathbb{R}^{4d}$.

Given the time-aware candidate news representation $\mathbf{z}^*$ and the time-aware news sequence representation $\mathbf{z}_{\mathcal{C}_i}$, concatenate them together and then perform a linear transformation to transform them into $d$-dimensional space. Specifically, take news $c_j$ for example, the time-aware representation $\mathbf{z}_{c_j}$ of news $c_j$ is concatenated with $\mathbf{z}^*$ and user embedding $\mathbf{u}_i$ and then linearly transformed as $\mathbf{t}_{c_j} = [\mathbf{z}_{c_j} \mathbf{z}^* \mathbf{u}_i] \mathbf{W}_{\mathbf{c}} \in \mathbb{R}^d$, where $\mathbf{W}_{\mathbf{c}} \in \mathbb{R}^{3d \times d}$. With this operation, the document-level attention layer can consider content rele-

vance, user preferences, and time dynamics simultaneously. To determine the representation of news $c_j$ concerning the candidate news $c^*$ and user $i$, we borrow the idea of Transformer, considering its strong capability in modeling the correlations between the events in a sequence of user behaviors. The original Transformer is designed for the NLP applications, where the sequential signals are captured by the word indexes. However, continuous time information can be essential in user behavior modeling, as mentioned before. Thus, we revise the traditional Transformer by replacing the position encoding with a continuous time embedding. The variables are firstly input into the self-attention layer, and then a position-wise feed-forward layer is leveraged to process the output. At last, we use a fully connected layer and a residual connection layer to predict the final results:

$$\begin{aligned} \mathbf{t}_{c_j} &= \text{Attention}(\mathbf{t}_{c_j}), \\ \mathbf{x}_{\mathbf{c_j}} &= \phi(\mathbf{t}_{c_j} \cdot \mathbf{W_a}), \\ \mathbf{x}_{\mathbf{c_j}} &= \text{LN}(\text{Dropout}(\mathbf{x}_{\mathbf{c_j}} \cdot \mathbf{W_b}) + \mathbf{t}_{c_j}). \end{aligned} \quad (4)$$

where $\phi$ is the activation function, $\mathbf{W_a} \in \mathbb{R}^{d \times d'}$, $\mathbf{W_b} \in \mathbb{R}^{d' \times d}$ are model parameters, and $d'$ is the intermediate size of position-wise feed-forward layer. 'Dropout' is an approach used for alleviating the overfitting problems, 'LN' represents layer normalization, which normalizes an input vector by its mean and variance for stable training.

### 4.4 Dynamic Negative Sampling Layer

Previous news recommendation models primarily leverage uniform negative sampling to optimize the users' implicit feedback. However, as discussed by [3], uniform negative sampling is not optimal since it selects too random samples, which can be less discriminative to the positive ones. We adopt a dynamic negative sampling (DNS) method to train our model more effectively. Our general idea is to build more informative item pairs, which are critical for model optimization. In our method, if an item pair is hard to separate, we regard it to be more informative since the model can learn more by optimizing it. To evaluate how difficult an item pair can be separated, we introduce a similarity function, that is: $f(\mathbf{y}, \mathbf{X}) = \mathbf{W}(\mathbf{X}^T \cdot \mathbf{y}) + \mathbf{b}$, where $\mathbf{y} \in \mathbb{R}^d$ and $\mathbf{X} \in \mathbb{R}^{d \times N}$ are the representations of the positive and the whole candidate news, respectively. $\mathbf{W} \in \mathbb{R}^{N \times N}$ and $\mathbf{d} \in \mathbb{R}^N$ are weighting parameters. Based on $f$, we select the most similar items in a greedy way: $S = \arg\max f(\mathbf{y}, \mathbf{X})$, where $S$ is the index set of the selected negative items. In our model, we further introduce a loss to constraint the parameters in $f$:

$$L_S = \log\left(y^T \text{Merge}(S, \mathbf{X})\right), \quad (5)$$

where Merge is a function projecting the embeddings of the selected negative samples into a vector. We would like to make the negative samples similar to the positive ones on different distance metrics by this loss function.

### 4.5 History Summarization Layer

To summarize the users' historical behaviors, we leverage Transformer to process the previously interacted news.

Compared with other sequential models like convolutional neural network (CNN) and recurrent neural network (RNN), Transformer directly captures the correlations between any two steps of events in the sequence, which is effective in many other machine learning tasks. In our model, to obtain the representation of news sequence, we stack the representation of $L$ news articles into a feature tensor $E \in \mathbb{R}^{L \times d}$. We adopt $M$ Transformer layers. The input of the first layer is $E$, and the previous layer's output is fed into the next layer. Multiple Transformer layers are capable of modeling abstract sequential patterns, and the output $\mathbf{p}_i$ of the last layer is used for model inference.

### 4.6 Prediction Layer

We concatenate the sequence vector $\mathbf{p}_i$, the representation $\mathbf{x}^*$ of candidate news $c^*$ and the user embedding $\mathbf{u}_i$, and feed them into a fully-connected layer to estimate the click rate of user $i$ click candidate news $c^*$: $\hat{y}_{i*} = \phi([\mathbf{p}_i \, \mathbf{x}^* \, \mathbf{u}_i]\mathbf{W}^{(1)} + \mathbf{b}^{(1)})\mathbf{W}^{(2)} + \mathbf{b}^{(2)}$, where $\mathbf{W}^{(1)} \in \mathbb{R}^{5d \times 2d}$, $\mathbf{b}^{(1)} \in \mathbb{R}^{2d}$, $\mathbf{W}^{(2)} \in \mathbb{R}^{2d \times 1}$ and $\mathbf{b}^{(2)} \in \mathbb{R}$ are the parameters of the prediction layer. Besides, the negative samples set $S$ output by dynamic negative sampling layer are combined with positive samples to compute loss. We adopt binary cross-entropy loss function:

$$L_P = \sum_{i,c^* \in D^+ \cup S} y_{i*} \log \sigma(\hat{y}_{i*}) + (1 - y_{i*}) \log(1 - \sigma(\hat{y}_{i*})), \quad (6)$$

where $y_{i*}$ is the instance label, $D^+$ represents the positive instance set. Therefore, the total loss of D-HAN is $L = L_P + \alpha L_S$, where $\alpha \in (0.0, 1.0]$ is a trade-off parameter.

## 5 EXPERIMENTAL STUDY

In this section, we conduct experiments to answer the following questions:

• **RQ1**: Does our proposed D-HAN model outperform several state-of-the-art models?

• **RQ2**: What are the impacts of each hierarchical attention layer on the performance of D-HAN model?

• **RQ3**: Is relative time interval or absolute timestamp of news items improve model performance?

• **RQ4**: Can the history summarization layer further improve the performance of hierarchical attention layers?

• **RQ5**: Are the attention weights able to learn meaningful information in sentence-, element- and document-level?

• **RQ6**: Can the dynamic negative sampling layer improve the recommendation performance?

### 5.1 Experimental Setup

**Datasets.** We conduct experiments on three datasets: Adressa, Cert and Caing. Each log of all datasets contains user ID, news ID, reading timestamp and news contents. Adressa (http://reclab.idi.ntnu.no/dataset/) and Caing (http://www.dcjingsai.com/) are available online.

• **Adressa**: This dataset is constructed by [9] for evaluating news recommendations. It contains reading logs of 10 weeks

TABLE 1
Statistics of the datasets.

| Datasets | #Interaction | #User | #News | Sparsity |
|---|---|---|---|---|
| Adressa | 1,604,879 | 66,649 | 12,034 | 99.79% |
| Cert | 1,573,959 | 199 | 588,907 | 98.65% |
| Caing | 61,615 | 1,947 | 5,275 | 99.40% |

from Adresseavisen, a Norwegian news portal. In this paper, we filter the first 15 days to do the experiments.

• **Cert**: This dataset is provided by the Computer Emergency Response Technical Team of China, **from March 2016 to April 2017**, containing each user's logs from various news portals, including people.com, cctv.com, et al.

• **Caing**: This dataset is from Caing, a famous news portal in China. It contains complete reading logs of 10,000 users during March 2016.

All the datasets are preprocessed to make sure that all users have at least 15 interactions. The statistical details of the three datasets are summarized in Table 1.

**Evaluation Protocols.** To evaluate the performance of news recommendation, we adopt the leave-one-out strategy, which has been widely used in [4], [15], [16]. For each user, we use a sliding window of $L + 1$ ($L$ historical news and 1 candidate news) length to slide over his interactions, and each window generates one instance. We hold out the latest instance for testing and utilize the remaining instances for training. For every user, we randomly sample 99 news articles during evaluation, which are not interacted by the user and rank the ground-truth news that the user has consumed among the 99 news. The performance of the ranked list is evaluated by Hit Ratio (HR) and Normalized Discounted Cumulative Gain (NDCG). The HR@N measures whether the ground-truth news is ranked on the top-N list, and the NDCG@N accounts for the position of hit by assigning higher scores to hits at top ranks [11].

**Baselines.** We compare D-HAN with the following methods, which can be divided into two groups. The first group only includes collaborative filtering methods:

• **BPR** [28]: This method optimizes the latent factor model with the pairwise ranking loss on implicit feedback data.

• **GRU4Rec** [13]: This method applies recurrent neural network for session-based recommendation.

• **Caser** [30]: This method utilizes horizontal and vertical convolution filters to simultaneously capture sequential patterns and model users' general preferences.

The second group contains recommendation methods that utilize content information.

• **GRU4Rec++** [12]: This method is an improved version of GRU4Rec, further considering news content information.

• **WE3CN** [16]: This method applies 3D convolution neural network for news recommendation, utilizing content and sequential information simultaneously.

• **DNA** [39]: This method is the preliminary work of this paper, where the same hierarchy architecture is adopted.

Many methods based on deep learning for news recommendation have been proposed [22], [26], [32], [40]. However, they require external information not available in our datasets, i.e., knowledge graphs by [32], news category

by [26], web browsing and searching records by [22] and online environment by [40]. Therefore, we do not include them as baselines. Moreover, in news recommendation, there exists a severe item cold-start problem for the collaborative filtering methods of the first group. In the testing phase, these models cannot use meaningful representations of unseen news articles. In this case, we use the average representation of the top-100 temporally-closest news articles instead.

**Implementation.** We use the named entity recognition and keywords extraction modules of the following NLP tools: For the Norwegian dataset Adressa, we utilize Polyglot (https://polyglot.readthedocs.io/en/latest/index.html) [1] which is a natural language pipeline that supports massive multilingual applications. For the two Chinese datasets Cert and Caing, we choose NLPIR (http://ictclas.nlpir.org) (also known as ICTCLAS), a tool that integrates many functions such as word segmentation, part-of-speech tagging, and named entity recognition. For BPR, Caser, and GRU4Rec, user embedding and news embedding dimensions are set as 20. For GRU4Rec++, the dimension of the content vector is set as 64. For WE3CN, the dimension of the word vector is set as 64. Other parameters in the baselines are set as default. We implement our D-HAN model and other deep learning models with Pytorch. The number of news $L$ of each sequence is set as 10. The number of sentences $K$ for each news is set as 20. The dimension $d$ is set as 64. The dimension of intermediate size $d'$ is set as 256. Merge$(x, y)$ looks up the corresponding vector of the first element in $x$ from $y$. The loss trade-off parameter $\alpha$ is set to 1.0. The number of the history summarization layer is set as 2. Adam optimizer is applied for training, and the learning rate, batch size, weight decay, and dropout are fixed to $10^{-3}$, 256, $10^{-4}$ and 0.2, respectively. Each experiment is repeated three times, and we report the average results.

## 5.2 Performance Comparison (RQ1)

The performance of our proposed D-HAN and two types of baselines on three datasets is shown in Table 2, in which we have the following observations:

Our proposed D-HAN model achieves the best performance on all the datasets, significantly outperform the best baseline. In specific, on the metrics of HR@10 and NDCG@10, D-HAN outperforms DNA by about 8.34% and 18.12% on Adressa dataset, 11.29% and 20.44% on Cert dataset and 4.85% and 2.92% on Caing dataset, respectively. We attribute the performance improvement to self-attention mechanism applied in each attention layer, the time information adopted in document-level attention layer, the Transformer encoder applied in modeling sequential information and the dynamic negative sampling method. We will detail the effectiveness of each component next sections.

Content information improves the performance of news recommendation. For example, in most cases, the models of the second group are superior to the models of the first group. More directly, GRU4Rec++ is merely added content information based on GRU4Rec and outperforms GRU4Rec on all datasets. However, Caser, a CF model, outperforms

TABLE 2
Performance comparison on three datasets for all methods.

| Adressa | HR@1 | HR@5 | HR@10 | NDCG@1 | NDCG@5 | NDCG@10 |
|---|---|---|---|---|---|---|
| BPR | 0.0777 | 0.2676 | 0.4278 | 0.0777 | 0.1724 | 0.2239 |
| GRU4Rec | 0.2969 | 0.6926 | 0.8535 | 0.2969 | 0.5026 | 0.5551 |
| Caser | 0.3327 | 0.7277 | 0.8684 | 0.3327 | 0.5397 | 0.5856 |
| GRU4Rec++ | 0.3423 | 0.7348 | 0.8773 | 0.3423 | 0.5474 | 0.5939 |
| WE3CN | 0.3070 | 0.6790 | 0.8340 | 0.3070 | 0.4965 | 0.5466 |
| DNA | 0.4528 | 0.8627 | 0.9505 | 0.4528 | 0.6726 | 0.7015 |
| D-HAN | **0.5685** | **0.9031** | **0.9653** | **0.5685** | **0.7486** | **0.7676** |
| Imp. | *25.55%* | *4.68%* | *1.56%* | *25.55%* | *11.30%* | *9.42%* |
| **Cert** | HR@1 | HR@5 | HR@10 | NDCG@1 | NDCG@5 | NDCG@10 |
| BPR | 0.2463 | 0.4539 | 0.5444 | 0.2463 | 0.3511 | 0.3801 |
| GRU4Rec | 0.3367 | 0.4623 | 0.5193 | 0.3367 | 0.4038 | 0.4222 |
| Caser | 0.4556 | 0.5812 | 0.6281 | 0.4556 | 0.5251 | 0.5403 |
| GRU4Rec++ | 0.3920 | 0.5343 | 0.5662 | 0.3920 | 0.4677 | 0.4781 |
| WE3CN | 0.3744 | 0.6884 | 0.8015 | 0.3744 | 0.5447 | 0.5818 |
| DNA | 0.5239 | 0.8116 | 0.8920 | 0.5239 | 0.6746 | 0.7007 |
| D-HAN | **0.7822** | **0.9447** | **0.9749** | **0.7822** | **0.8662** | **0.8726** |
| Imp. | *49.30%* | *16.40%* | *9.29%* | *49.30%* | *28.40%* | *24.53%* |
| **Caing** | HR@1 | HR@5 | HR@10 | NDCG@1 | NDCG@5 | NDCG@10 |
| BPR | 0.3546 | 0.5728 | 0.6774 | 0.3546 | 0.4699 | 0.5038 |
| GRU4Rec | 0.4633 | 0.7713 | 0.8541 | 0.4633 | 0.6314 | 0.6586 |
| Caser | 0.5999 | 0.7964 | 0.8514 | 0.5999 | 0.7057 | 0.7235 |
| GRU4Rec++ | 0.6150 | 0.8139 | 0.8615 | 0.6150 | 0.7237 | 0.7393 |
| WE3CN | 0.4992 | 0.6622 | 0.7329 | 0.4992 | 0.5861 | 0.6089 |
| DNA | 0.6220 | 0.8391 | 0.9033 | 0.6220 | 0.7401 | 0.7609 |
| D-HAN | **0.7026** | **0.8804** | **0.9279** | **0.7026** | **0.7979** | **0.8119** |
| Imp. | *12.96%* | *4.93%* | *2.72%* | *12.96%* | *7.81%* | *6.71%* |

The bolded numbers are the best results of each column. Relative improvements of D-HAN to the best baseline are presented in the last row of each child table. Note that to compare model performance more fairly, time information is not considered here.

GRU4Rec++ on Cert dataset. This is probably because many adjacent actions do not have apparent dependency relationships in the Cert dataset. The RNN-based GRU4Rec++ cannot handle datasets with this characteristics, while CNN-based Caser can still capture sequential patterns with convolutional filters. Caser outperforms WE3CN on Adressa and Caing datasets, this is probably because, WE3CN represents each news article with the first 50 words, which are not enough to express the contents of news articles very well.

Sequential information improves the performance of news recommendation. BPR only utilizes users' feedback information and performs the worst. In addition to users' feedback information, GRU4Rec and Caser both utilize sequential information of users' reading behaviors and perform better than BPR. This suggests the effectiveness of considering sequential information in news recommendation.

### 5.3 Impacts of Attention Layer (RQ2)

A unique design of our model is leveraging attention mechanism to process the news information on the sentence-, element- and document-level. Comparing the previous DNA model, we have improved the attention computing methods. In this section, we study the effectiveness of the revised attention methods. From Table 3, we can observe:

With sentence-level attention layer, D-HAN$_1$ improves the performance by about 1.81%, 13.49% and 3.63% on HR@1 for different datasets, respectively. This result demonstrates that the scaled dot-product attention used in D-HAN$_1$ is more effective than the additive attention mechanism used in DNA.

With element-level attention layer, D-HAN$_2$ improves the performance by about 7.84%, 20.21% and 4.63% on HR@1 for all the three datasets, respectively, which indicates that there is much space to explore for element-level information, and it is important to make full use of the internal relationship between two elements in one piece of news and the external relationship between two elements in pairwise news. Note that D-HAN$_2$ is better than D-HAN$_1$, this is contrary to our intuition. Generally, we think that news text plays a greater role in news recommendation than news elements. Probably this is because that sentence-level vector instead of the token-level vector is used as the input of sentence-level attention component for a fair comparison with baselines, and this operation may result in the loss of much information. Exploring the representation of sentences by stacking the representation of tokens is left to be done in the future.

With document-level attention layer, the performance of D-HAN$_3$ is better than D-HAN$_1$ and D-HAN$_2$. Specifically, for Cert dataset, D-HAN$_3$ gets improvements by 32.37%, 14.54%, 20.81% and 18.2% on HR@1, HR@5, NDCG@5 and NDCG@10, respectively. It demonstrates that the document-level information is pretty important for news sequence recommendation, and Transformer encoder adopted in document-level attention component is effective.

D-HAN$_4$ which adopts sentence-, element- and document-level attention components simultaneously can achieve the best performance. Specifically, D-HAN$_4$ get the most performance improvements, 21.86%, 38.77% and 12.93% on HR@1 for all three datasets.

TABLE 3
Comparison between sentence-, element-, and document-level attention layer.

| Adressa | Layer | HR@1 | HR@5 | HR@10 | NDCG@1 | NDCG@5 | NDCG@10 |
|---------|-------|------|------|-------|--------|--------|---------|
| DNA | - | 0.4528 | 0.8627 | 0.9505 | 0.4528 | 0.6726 | 0.7015 |
| D-HAN$_1$ | S | 0.4610 | 0.8624 | 0.9473 | 0.4610 | 0.6764 | 0.7035 |
| D-HAN$_2$ | E | 0.4883 | 0.8665 | 0.9488 | 0.4883 | 0.6904 | 0.7171 |
| D-HAN$_3$ | N | 0.5515 | 0.8997 | **0.9634** | 0.5515 | 0.7397 | 0.7601 |
| D-HAN$_4$ | S+E+N | **0.5518** | **0.9001** | 0.9630 | **0.5518** | **0.7398** | **0.7604** |
| **Cert** | Layer | HR@1 | HR@5 | HR@10 | NDCG@1 | NDCG@5 | NDCG@10 |
| DNA | - | 0.5239 | 0.8116 | 0.8920 | 0.5239 | 0.6746 | 0.7007 |
| D-HAN$_1$ | S | 0.5946 | 0.8894 | 0.9581 | 0.5946 | 0.7501 | 0.7682 |
| D-HAN$_2$ | E | 0.6298 | 0.8794 | 0.9564 | 0.6298 | 0.7565 | 0.7804 |
| D-HAN$_3$ | N | 0.6935 | 0.9296 | 0.9732 | 0.6935 | 0.8150 | 0.8282 |
| D-HAN$_4$ | S+E+N | **0.7270** | **0.9430** | **0.9782** | **0.7270** | **0.8381** | **0.8499** |
| **Caing** | Layer | HR@1 | HR@5 | HR@10 | NDCG@1 | NDCG@5 | NDCG@10 |
| DNA | - | 0.6220 | 0.8391 | 0.9033 | 0.622 | 0.7401 | 0.7609 |
| D-HAN$_1$ | S | 0.6446 | 0.8391 | 0.9021 | 0.6445 | 0.7468 | 0.7667 |
| D-HAN$_2$ | E | 0.6487 | 0.8408 | 0.9041 | 0.6487 | 0.7477 | 0.7681 |
| D-HAN$_3$ | N | 0.7023 | **0.8833** | **0.9334** | 0.7023 | 0.7949 | 0.8113 |
| D-HAN$_4$ | S+E+N | **0.7024** | 0.8759 | 0.9256 | **0.7024** | **0.7969** | **0.8115** |

S, E, N represents sentence-, element-, document-level attention component are adopted, respectively.
The bolded numbers are the best result of each column. Note that to demonstrate the effectiveness of
document-level attention component, time information is not adopted here.

TABLE 4
Time embedding performance comparison on three datasets.

| Adressa | HR@1 | HR@5 | HR@10 | NDCG@1 | NDCG@5 | NDCG@10 |
|---------|------|------|-------|--------|--------|---------|
| D-HAN$_{-RA}$ | 0.5685 | 0.9031 | 0.9653 | 0.5685 | 0.7486 | 0.7676 |
| D-HAN$_{-A}$ | 0.5990 | 0.9182 | 0.9698 | 0.5990 | 0.7718 | 0.7889 |
| D-HAN$_{-R}$ | 0.6327 | 0.9299 | 0.9752 | 0.6327 | 0.7947 | 0.8095 |
| D-HAN | **0.6355** | **0.9324** | **0.9761** | **0.6355** | **0.7977** | **0.8119** |
| **Cert** | HR@1 | HR@5 | HR@10 | NDCG@1 | NDCG@5 | NDCG@10 |
| D-HAN$_{-RA}$ | 0.7822 | 0.9447 | 0.9749 | 0.7822 | 0.8662 | 0.8726 |
| D-HAN$_{-A}$ | **0.7822** | **0.9548** | **0.9799** | **0.7822** | **0.8673** | **0.8754** |
| D-HAN$_{-R}$ | 0.6617 | 0.9112 | 0.9682 | 0.6617 | 0.7835 | 0.8046 |
| D-HAN | 0.7513 | 0.9346 | 0.9774 | 0.7513 | 0.8380 | 0.8515 |
| **Caing** | HR@1 | HR@5 | HR@10 | NDCG@1 | NDCG@5 | NDCG@10 |
| D-HAN$_{-RA}$ | 0.7026 | 0.8804 | 0.9279 | 0.7026 | 0.7979 | 0.8119 |
| D-HAN$_{-A}$ | 0.7182 | 0.8872 | 0.9336 | 0.7182 | 0.8078 | 0.8237 |
| D-HAN$_{-R}$ | 0.7191 | 0.8930 | 0.9346 | 0.7191 | 0.8090 | 0.8235 |
| D-HAN | **0.7244** | **0.8942** | **0.9361** | **0.7244** | **0.8153** | **0.8294** |

The bolded numbers are the best results of each column.

## 5.4 Influence of Time Embedding (RQ3)

Another character of our model is the comprehensive time information modeling. In specific, we incorporate both absolute and relative time information in the document-level attention layer. We compare our final model with its three variants: D-HAN$_{-RA}$ is a method without considering any continuous time information. D-HAN$_{-A}$ is a method by removing the absolute time information in the final model. D-HAN$_{-R}$ is a method by removing the relative time information in the final model. Relative time intervals can be computed by a) time difference between any two news items in a sequence, b) time difference between two adjacent news

in a sequence, c) time difference between candidate news and each news in a sequence. [20] adopts the first one to form a relative time interval matrix, and we adopt the last one since it is more effective in our experiments. The results are presented in Table 4. We can see:

Without relative time information, D-HAN$_{-R}$ decreases the model performance by 0.44% and 0.73% on HR@1 for Adressa and Caing, respectively. This result demonstrates the effectiveness of the relative time information for the user behavior modeling. Without absolute time information, D-HAN$_{-A}$ causes 5.7% and 0.86% drop in HR@1 for Adressa and Caing, respectively, which manifests that the absolute

TABLE 5
Comparison of history summarization layer.

| Adressa | HR@1 | HR@5 | HR@10 | NDCG@1 | NDCG@5 | NDCG@10 |
|---|---|---|---|---|---|---|
| D-HAN$_{-S}$ | 0.5518 | 0.9001 | 0.9630 | 0.5518 | 0.7398 | 0.7604 |
| CNN | 0.5551 | 0.9014 | 0.9638 | 0.5551 | 0.7419 | 0.7620 |
| LSTM | 0.5593 | 0.9036 | 0.9653 | 0.5593 | 0.7456 | 0.7659 |
| BiLSTM | 0.5621 | 0.9057 | 0.9658 | 0.5621 | 0.7473 | 0.7671 |
| SH | **0.5685** | 0.9031 | 0.9653 | **0.5685** | 0.7486 | 0.7676 |
| MH=2 | 0.5653 | **0.9059** | **0.9662** | 0.5653 | **0.7499** | **0.7677** |
| MH=4 | 0.5634 | 0.9054 | 0.9656 | 0.5633 | 0.7477 | 0.7673 |
| MH=8 | 0.5636 | 0.9041 | 0.9651 | 0.5635 | 0.7469 | 0.7665 |
| **Cert** | HR@1 | HR@5 | HR@10 | NDCG@1 | NDCG@5 | NDCG@10 |
| D-HAN$_{-S}$ | 0.7270 | 0.9430 | 0.9782 | 0.7270 | 0.8381 | 0.8499 |
| CNN | **0.7889** | 0.9464 | 0.9749 | **0.7889** | **0.8687** | **0.8791** |
| LSTM | 0.7236 | 0.9112 | 0.9698 | 0.7236 | 0.8172 | 0.8280 |
| BiLSTM | 0.7085 | 0.9028 | 0.9664 | 0.7085 | 0.8042 | 0.8228 |
| SH | 0.7822 | 0.9447 | 0.9749 | 0.7822 | 0.8662 | 0.8726 |
| MH=2 | 0.7672 | 0.9447 | 0.9732 | 0.7672 | 0.8531 | 0.8618 |
| MH=4 | 0.7822 | **0.9548** | 0.9799 | 0.7822 | 0.8677 | 0.8723 |
| MH=8 | 0.7571 | 0.9498 | **0.9832** | 0.7571 | 0.8560 | 0.8643 |
| **Caing** | HR@1 | HR@5 | HR@10 | NDCG@1 | NDCG@5 | NDCG@10 |
| D-HAN$_{-S}$ | 0.7024 | 0.8759 | 0.9256 | 0.7024 | 0.7969 | 0.8115 |
| CNN | 0.6968 | 0.8752 | 0.9276 | 0.6968 | 0.7925 | 0.8093 |
| LSTM | 0.6901 | 0.8639 | 0.9175 | 0.6901 | 0.7801 | 0.7978 |
| BiLSTM | 0.6976 | 0.8725 | 0.9245 | 0.6976 | 0.7902 | 0.8067 |
| SH | 0.7026 | 0.8804 | 0.9279 | 0.7026 | **0.7979** | 0.8119 |
| MH=2 | **0.7062** | 0.8805 | 0.9291 | **0.7062** | 0.7972 | 0.8126 |
| MH=4 | 0.7049 | 0.8818 | **0.9309** | 0.7049 | 0.7971 | 0.8129 |
| MH=8 | 0.7054 | **0.8819** | 0.9305 | 0.7054 | 0.7971 | **0.8131** |

MH=2 denotes the number of heads of multi-head Transformer is 2. SH denotes single head. The bolded numbers are the best result of each column. D-HAN$_{-S}$ represents D-HAN removing history summarization layer.

TABLE 6
Comparison of different negative sampling methods.

| Adressa | HR@1 | NDCG@5 | NDCG@10 |
|---|---|---|---|
| D-HAN | **0.5685** | 0.7486 | 0.7676 |
| D-HAN$_D$ | 0.5666 | **0.7508** | **0.7698** |
| **Cert** | HR@1 | NDCG@5 | NDCG@10 |
| D-HAN | 0.7822 | 0.8662 | 0.8726 |
| D-HAN$_D$ | **0.7915** | **0.8782** | **0.8856** |
| **Caing** | HR@1 | NDCG@5 | NDCG@10 |
| D-HAN | 0.7026 | 0.7979 | 0.8119 |
| D-HAN$_D$ | **0.7089** | **0.8031** | **0.8182** |

D-HAN$_D$ denotes D-HAN with DNS.

time information can also be useful in the news recommendation scenarios. Besides, D-HAN$_{-R}$ is better than D-HAN$_{-A}$ in most cases, indicating that absolute time is more informative than relative time interval. On Adressa and Caing dataset, D-HAN performs best, indicating that relative time interval and absolute time can complement each other, and the combination of these two time information is more effective than using either one of them alone.

## 5.5 Impacts of History Summarization Layer (RQ4)

In this section, we compare the history summarization layer performance by an ablation study. D-HAN$_{-S}$ removes the history summarization layer. Besides, we also exploit different sequential models for user history summarization. Specifically, we leverage CNN, LSTM, BiLSTM, single head Transformer and multi-head Transformer to process the previous user behaviors. The results are presented in Table 5, we can see: (a) Comparing the results of D-HAN$_{-S}$ with Transformer-based models, history summarization layer improves model performance for all three datasets. For Adressa dataset, setting CNN, LSTM or BiLSTM as history summarization layer can improve model performance and CNN performs the best on Cert dataset, while for Cert and Caing dataset, RNN-based models can not promote model performance. Such observation manifests that CNN and RNN-based summarization methods have their own advantages on some specific datasets, but it is encouraging to see that Transformer-based models can consistently achieve the best performance. We speculate that Transformer can directly model the correlations between any two steps in the news sequence, more sufficiently capturing the sequential patterns. (b) Comparing the results of SH with MHs indicates that more heads do not necessarily bring significant performance improvements. This may be because the dimension in our

experiment is small ($d = 64$), resulting in each head playing a small role [20]. In our experiments, we use a single-head Transformer for brevity.

## 5.6 Visualization (RQ5)

Figure 2 presents the case studies of the sentence-, element-, document- and time-aware document-level attention weights to show what each module learns for each level information. We randomly sample pairs of historical and candidate news from the testing set. The attention weights are extracted from the best performance epoch. Note that figure 2(c) and figure 2(d) use the same instance for fair comparison. Due to the space limitation, we only present the visualization on Caing dataset.

**Sentence-level attention weights:** Figure 2(a) shows a typical self-attention weights of a history-candidate news pair, where the candidate news representation is concatenated as the last sentence of the history news, as introduced in Section 4.1. Note that in our experiments, 20 sentences are extracted from each news item, and insufficient ones are padded to the left. The result shows that sentences tend to attend to the first few sentences, and sentences in the latter part are less attended. Specifically, the upper left is much brighter (larger weights), and the bottom right is much darker (smaller weights). This is possibly because the former is the ground-truth sentences of the news while the latter is pads. In addition, the result also shows that the first three sentences are significantly attended, which is consistent with the text characteristics in the news scene, and the first few sentences maintain the main content of the entire news. Besides, the last column and the last row have higher attention weights than the adjacent rows and columns, but the bottom right corner candidate news attention weight is lower. This is possibly because candidate news provides the reference information for historical news, but not for itself.

**Element-level attention weights:** Figure 2(b) shows a heatmap of attention weights between two elements of one history-candidate news pair. Different from Figure 2(a), this pair is concatenated along the $d$ dimension as detailed in Section 4.2. The x-axis and y-axis show the five elements of news we use, e.g., time, person, organization, location and keywords, respectively. As the figure shows, the person element plays the largest role, followed by the location element, demonstrating that these two elements are more important than others in this instance. Unlike the heatmap of sentence-level attention weights, which has fixed and typical patterns, the heatmap of element-level attention weights is different for different instances. This may be due to that these five artificially selected elements are already essential elements for news sequence recommendation.

**Document-level attention weights:** We firstly describe several possible cases demonstrating user behavior patterns:

**Case 1:** Click news with the relevant or same topic. This kind of case often happens when the read news is newly released, such as major events or hot news, and users tend to browse more relevant news to obtain more relevant information. When the event gets out of date quickly, users will no longer be interested in it. But if this event lasts for a long time, users will be interested in this kind of news for a long time. For instance, a piece of news about major casualties on the highway, when a user reads this news for the first time, his next step may tend to click relevant news. But if a piece of news is about the war between Azerbaijan and Armenia, since the war will not stop in a short time, users tend to track the progress of the war.

**Case 2:** Click other news instead of hot news read by many users. For two pieces of news, one of them is clicked many times and the other is not received by many people. If the user skips the former and reads the latter, it means that the user is not interested in news relevant to the former.

**Case 3:** Click news of a user's long-term interests. Although recent behaviors greatly influence a user's current news click behavior, her interest is often hidden in long-term historical records, it is a stable tendency.

**Case 4:** Click news in which a user is interested during a specific time period. The time period has an impact on the users' propensity, e.g., browsing news related to the stock market at around 10 am on weekdays, browsing news related to food cooking on weekends, and reading news related to constellations at midnight.

Figure 2(c) shows a typical heatmap of document-level attention weights between news items in one news sequence, it is an example of case 1. We can clearly see that the upper left of the figure has lower weights, while the bottom, left and bottom left have higher weights. Specifically, the $9th$ news items have the largest weights followed by the $6th$ news, reflecting the former most influence user decision followed by the latter in this instance. In addition, to compare the influence of more distant and recent news on users' decisions, we randomly sample 500 instances in the test set and check the heatmap of their document-level attention weights. We divide the influence into three categories: (a) recent news have larger influence, when the weights are concentrated in the $5th$ to $10th$ news. (b) distant news have larger influence, when the weights focus on the $1th$ to $5th$ news. (c) older news have the same influence as recent news, when the weights fall evenly or can not be clearly distinguished. The statistics show that these three categories account for 46.8%, 23.6% and 29.6%, respectively, indicating that user behavior patterns in news recommendation scenes are more affected by recent behavior.

**Time-aware document-level attention weights:** For the above cases and results, we can see that the time interval and timestamp are pretty key information, that is, one user's operations are different in different time interval scales or different periods. For example, if a user reads a newly released hot news, then in a short time, e.g., 30 minutes, he might be full of interest and curiosity about the news. But for a long time, e.g., two days, he might be no longer interested in this type of news, reflecting the timeliness of news scene.

Figure 2(d) shows the heatmap of document-level attention weights of Figure 2(c) embeded with time information. Comparing Figure 2(d) with Figure 2(c), we can see that time-aware document-level attention weights are more sparse, and weights are dominated by the $5th$ news item, indicating relative time interval and absolute timestamp

(a) Sentence-level $w$     (b) Element-level $w$     (c) Document-level $w$     (d) Time-aware document-level $w$
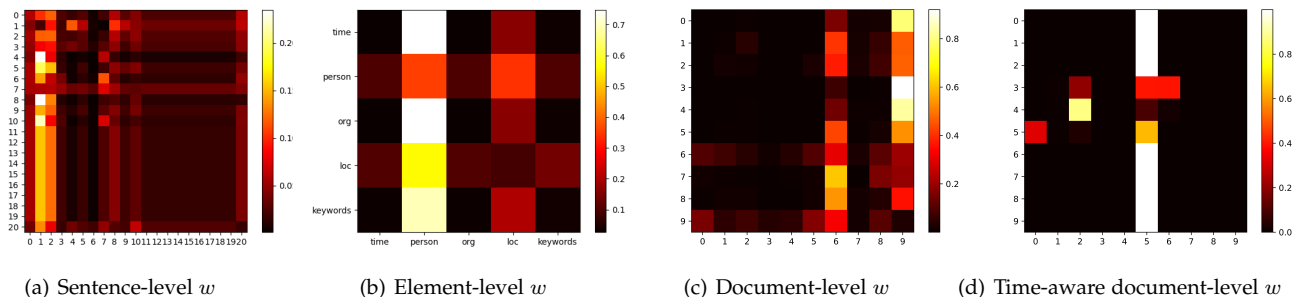
Fig. 2. Visualization of sentence-, element-, document- and time-aware document-level attention weights between historical news sequence and candidate news. $w$ represents attention weights.

embedding do bring great influence on the weights.

Too long sequence may cause historical information overload, i.e., distant information accounts for too high proportion for news recommendation, resulting in recommended news being limited to a certain range and lack diversity. To the opposite, too short sequence will lead to instant information overload, i.e., recent news account for too high proportion for news recommendation, resulting in recommended news being limited to the type of previous one. At present, however, the length of news sequence is set to a fixed value based on experiments, e.g., 10 in our experiment. We have an idea that the long-term and short-term sequences can be processed separately by different models and then combined with different proportions dynamically to use historical and instant information. We left this idea in the future.

### 5.7 Impacts of DNS Layer (RQ6)

In this section, we study the impacts of the DNS layer. For a fair comparison, we adopt DNS in the training phase, but uniform negative sampling in the test phase. Due to the space limitation, we only present the results on HR@1, NDCG@5 and NDCG@10 as shown in Table 6. We can observe that DNS improves the model performance in most cases. This results verify the effectiveness of our designed dynamic negative sampling method. Comparing with the uniform sampling, DNS can well discover more informative negative samples with higher discriminative abilities.

**Summary.** From these tests, we find the following.

(1) Our approach D-HAN is effective for news recommendation. HR@N and NDCG@N scores of D-HAN are consistently higher than compared methods in all datasets.

(2) Each component in D-HAN improves the performance. We provide insights into attention weights by experimentally demonstrating their reasonableness.

(3) Dynamic strategies dealing with negative sampling are more effective.

## 6 CONCLUSION

We propose to use sentence-, element- and document-level information simultaneously by incorporating them into a hierarchical attention network in the news recommendation scene and applying Transformer encoder to capture news

sequential information. To model the dynamics of news recommendation, we propose a time-aware document-level attention layer to incorporate relative time interval and absolute timestamp. To generate more informative negative samples, we propose a dynamic negative sampling method that generates negative instances dynamically while training and then guide model optimization.

## REFERENCES

[1] R. Al-Rfou, V. Kulkarni, B. Perozzi, and S. Skiena. Polyglot-ner: Massive multilingual named entity recognition. In Proceedings of the 2015 SIAM International Conference on Data Mining, pages 586–594. SIAM, 2015.

[2] D. Bahdanau, K. Cho, and Y. Bengio. Neural machine translation by jointly learning to align and translate. In ICLR, 2015.

[3] R. Bamler and S. Mandt. Extreme classification via adversarial softmax approximation. arXiv preprint arXiv:2002.06298, 2020.

[4] J. Chen, H. Zhang, X. He, L. Nie, W. Liu, and T. Chua. Attentive collaborative filtering: Multimedia recommendation with item- and component-level attention. In SIGIR, 2017.

[5] A. Das, M. Datar, A. Garg, and S. Rajaram. Google news personalization: scalable online collaborative filtering. In WWW, 2007.

[6] 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.

[7] F. Garcin, C. Dimitrakakis, and B. Faltings. Personalized news recommendation with context trees. In RecSys, 2013.

[8] S. Ge, C. Wu, F. Wu, T. Qi, and Y. Huang. Graph enhanced representation learning for news recommendation. In Proceedings of The Web Conference 2020, pages 2863–2869, 2020.

[9] J. A. Gulla, L. Zhang, P. Liu, Ö. Özgöbek, and X. Su. The adressa dataset for news recommendation. In Proceedings of the International Conference on Web Intelligence, Leipzig, Germany, August 23-26, 2017, pages 1042–1048, 2017.

[10] F. Hamborg, C. Breitinger, M. Schubotz, S. Lachnit, and B. Gipp. Extraction of main event descriptors from news articles by answering the journalistic five W and one H questions. In JCDL, 2018.

[11] X. He, L. Liao, H. Zhang, L. Nie, X. Hu, and T. Chua. Neural collaborative filtering. In WWW, 2017.

[12] B. Hidasi and A. Karatzoglou. Recurrent neural networks with top-k gains for session-based recommendations. In Proceedings of the 27th ACM International Conference on Information and Knowledge Management, pages 843–852, 2018.

[13] B. Hidasi, A. Karatzoglou, L. Baltrunas, and D. Tikk. Session-based recommendations with recurrent neural networks. CoRR, abs/1511.06939, 2015.

[14] C. Hsieh, L. Yang, H. Wei, M. Naaman, and D. Estrin. Immersive recommendation: News and event recommendations using personal digital traces. In WWW, 2016.

[15] D. Khattar, V. Kumar, V. Varma, and M. Gupta. HRAM: A hybrid recurrent attention machine for news recommendation. In CIKM, pages 1619–1622, 2018.

[16] D. Khattar, V. Kumar, V. Varma, and M. Gupta. Weave&rec: A word embedding based 3-d convolutional network for news recommendation. In CIKM, 2018.

[17] Q. V. Le and T. Mikolov. Distributed representations of sentences and documents. In ICML, pages 1188–1196, 2014.

[18] J. Li, J. Li, and J. Tang. A flexible topic-driven framework for news exploration. In Proceedings of KDD, volume 2007, 2007.

[19] J. Li, P. Ren, Z. Chen, Z. Ren, T. Lian, and J. Ma. Neural attentive session-based recommendation. In CIKM, 2017.

[20] J. Li, Y. Wang, and J. McAuley. Time interval aware self-attention for sequential recommendation. In Proceedings of the 13th International Conference on Web Search and Data Mining, pages 322–330, 2020.

[21] L. Li, D. Wang, T. Li, D. Knox, and B. Padmanabhan. SCENE: a scalable two-stage personalized news recommendation system. In SIGIR, 2011.

[22] J. Lian, F. Zhang, X. Xie, and G. Sun. Towards better representation learning for personalized news recommendation: a multi-channel deep fusion approach. In IJCAI, pages 3805–3811, 2018.

[23] P. Lv, X. Meng, and Y. Zhang. Fere: Exploiting influence of multi-dimensional features resided in news domain for recommendation. Inf. Process. Manage., 53(5):1215–1241, 2017.

[24] Y. Lv, T. Moon, P. Kolari, Z. Zheng, X. Wang, and Y. Chang. Learning to model relatedness for news recommendation. In WWW, 2011.

[25] T. Mikolov, K. Chen, G. Corrado, and J. Dean. Efficient estimation of word representations in vector space. CoRR, abs/1301.3781, 2013.

[26] S. Okura, Y. Tagami, S. Ono, and A. Tajima. Embedding-based news recommendation for millions of users. In SIGKDD, 2017.

[27] K. Park, J. Lee, and J. Choi. Deep neural networks for news recommendations. In CIKM, 2017.

[28] S. Rendle, C. Freudenthaler, Z. Gantner, and L. Schmidt-Thieme. BPR: bayesian personalized ranking from implicit feedback. In UAI, 2009.

[29] F. Sun, J. Liu, J. Wu, C. Pei, X. Lin, W. Ou, and P. Jiang. Bert4rec: Sequential recommendation with bidirectional encoder representations from transformer. In Proceedings of the 28th ACM International Conference on Information and Knowledge Management, pages 1441–1450, 2019.

[30] J. Tang and K. Wang. Personalized top-n sequential recommendation via convolutional sequence embedding. In WSDM, 2018.

[31] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin. Attention is all you need. In Advances in neural information processing systems, pages 5998–6008, 2017.

[32] H. Wang, F. Zhang, X. Xie, and M. Guo. DKN: deep knowledge-aware network for news recommendation. In WWW, 2018.

[33] W. Wang. Chinese news event 5w1h semantic elements extraction for event ontology population. In WWW(Companion Volume), pages 197–202, 2012.

[34] C. Wu, F. Wu, S. Ge, T. Qi, Y. Huang, and X. Xie. Neural news recommendation with multi-head self-attention. In Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP), pages 6390–6395, 2019.

[35] J. Xiao, H. Ye, X. He, H. Zhang, F. Wu, and T. Chua. Attentional factorization machines: Learning the weight of feature interactions via attention networks. In IJCAI, pages 3119–3125, 2017.

[36] K. Xu, J. Ba, R. Kiros, K. Cho, A. C. Courville, R. Salakhutdinov, R. S. Zemel, and Y. Bengio. Show, attend and tell: Neural image caption generation with visual attention. In ICML, pages 2048–2057, 2015.

[37] F. Yuan, A. Karatzoglou, I. Arapakis, J. M. Jose, and X. He. A simple convolutional generative network for next item recommendation. In WSDM, 2019.

[38] Q. Yuan, G. Cong, Z. Ma, A. Sun, and N. Magnenat-Thalmann. Time-aware point-of-interest recommendation. In SIGIR, 2013.

[39] H. Zhang, X. Chen, and S. Ma. Dynamic news recommendation with hierarchical attention network. In 2019 IEEE International Conference on Data Mining (ICDM), pages 1456–1461. IEEE, 2019.

[40] G. Zheng, F. Zhang, Z. Zheng, Y. Xiang, N. J. Yuan, X. Xie, and Z. Li. DRN: A deep reinforcement learning framework for news recommendation. In WWW, 2018.

[41] G. Zhou, X. Zhu, C. Song, Y. Fan, H. Zhu, X. Ma, Y. Yan, J. Jin, H. Li, and K. Gai. Deep interest network for click-through rate prediction. In SIGKDD, pages 1059–1068, 2018.

[42] Y. Zhu, H. Li, Y. Liao, B. Wang, Z. Guan, H. Liu, and D. Cai. What to do next: Modeling user behaviors by time-lstm. In IJCAI, 2017.

**Qinghua Zhao** is a PhD student at the School of Computer Science and Engineering, Beihang University, supervised by Prof. Shuai Ma. She received his M.S. degree in computer technology from University of Chinese Academy of Sciences in 2018. Her research focuses on news recommendation and sentiment analysis.



**Xu Chen** is currently an assistant professor at Gaoling School of Artificial Intelligence, Renmin University of China. He obtained his PhD degree from Tsinghua University, China. He has published about 20 papers on top-tier conferences and journals such as SIGIR, TOIS, WWW, WSDM, CIKM, AAAI. His papers have won the best paper honorable mention award on the Web Conference 2018 and the best paper award on AIRS 2017.



**Hui Zhang** is currently a software engineer at Agricultural Bank of China, Beijing, China. She obtained her M.S. degree in computer technology from Beihang University in 2020 and B.S. degree in computer science and technology from Southeast University in 2017. Her current research interests include news recommendation and data mining.



**Shuai Ma** is a professor at the School of Computer Science and Engineering, Beihang University, China. He obtained his PhD degrees from University of Edinburgh in 2010, and from Peking University in 2004, respectively. He is a recipient of the best paper award for VLDB 2010. He is an Associate Editor of VLDB Journal since 2017, Knowledge and Information Systems since 2020 and IEEE Transactions On Big Data Since 2020. His current research interests include database theory and systems.