# MM-Rec: Multimodal News Recommendation

Chuhan Wu[1], Fangzhao Wu[2], Tao Qi[1], Yongfeng Huang[1]

[1]Department of Electronic Engineering & BNRist, Tsinghua University, Beijing 100084

[2]Microsoft Research Asia, Beijing 100080, China

{wuchuhan15,wufangzhao,taoqi.qt}@gmail.com,yfhuang@tsinghua.edu.cn

## ABSTRACT

Accurate news representation is critical for news recommendation. Most of existing news representation methods learn news representations only from news texts while ignoring the visual information in news like images. In fact, users may click news not only because of the interest in news titles but also due to the attraction of news images. Thus, images are useful for representing news and predicting user behaviors. In this paper, we propose a multimodal news recommendation method, which can incorporate both textual and visual information of news to learn multimodal news representations. We first extract region-of-interests (ROIs) from news images via objective detection. Then we use a pre-trained visiolinguistic model to encode both news texts and news image ROIs and model their inherent relatedness using co-attentional Transformers. In addition, we propose a crossmodal candidate-aware attention network to select relevant historical clicked news for accurate user modeling by measuring the crossmodal relatedness between clicked news and candidate news. Experiments validate that incorporating multimodal news information can effectively improve news recommendation.

## KEYWORDS

News recommendation, multimodality, news understanding

## 1 INTRODUCTION

Learning accurate news representations is the backbone of news recommendation [11, 26]. Most of existing news representations methods learn news representations merely from news texts [1, 4, 9, 16, 21–23, 25]. For example, Okura et al. [16] used autoencoders to learn news representations from the news content. In NPA [24] convolutional neural network (CNN) and personalized attention network are used to learn news representations from news titles. NRMS [25] uses multi-head self-attention networks to model news from titles. In fact, besides the titles, many news articles are associated with images when displayed on news websites, as shown

Figure 1: News with images for news recommendation.

in Fig. 1. Users may click news articles to read not only because of the interest in the content of the news title, but also due to the fascination of the corresponding image [2]. For example, in Fig. 1 the image of the second clicked news displays a highlight moment in an NFL game, which may be attractive for users interested in American football. Thus, the visual information of news images can provide rich complementary information for news understanding and user behavior prediction.

In this paper we study how to incorporate visual information of news to enhance news recommendation. Our work is motivated by the following observations. Fist, news titles and images usually have some inherent relatedness in describing news content and attracting user attention. For example, in the second news of Fig. 1, the word "Cowboys" in news title is related to the players shown in the news image. Modeling the relatedness between news titles and images can help better model news and infer user interest for news recommendation. Second, a user may have multiple interests, and a candidate news may only be related to a specific interest encoded in part of the historical clicked news. For example, in Fig. 1 the candidate news is only related to the second clicked news. Thus, modeling the relevance between clicked news and candidate news can help predict users' interest in a specific candidate news more accurately. In addition, candidate news may have crossmodal relatedness with clicked news. In Fig. 1 the image of candidate news is related to both the image and title of the second clicked news, because both images show the same football team and its name is mentioned in the title of the second clicked news. Modeling the crossmodal relations between candidate news and clicked news is helpful for measuring their relevance accurately.

In this paper, we present a multimodal news recommendation method named MM-Rec, which can leverage both textual and visual news information for news recommendation. In our approach, we first extract region-of-interests (ROIs) of news images via a pre-trained Mask R-CNN model [7] for objective detection. Then we use a pre-trained visiolinguistic model [14] to encode both news texts and news image ROIs and model their inherent crossmodal relatedness via a co-attentional Transformer network to learn accurate multimodal news representations. In addition, we propose a

crossmodal candidate-aware attention network to select relevant clicked news for user modeling by evaluating the crossmodal relevance between candidate news and clicked news, which can help better model users' specific interest in candidate news. Experiments on real-world dataset show that incorporating multimodal news information can help effectively improve news recommendation.

## 2 MM-REC

In this section, we introduce our *MM-Rec* method that incorporates both textual and visual news information for news recommendation. We first introduce its multimodal news encoder that learns multimodal news representations from news texts and images, and then introduce multimodal news recommendation based on crossmodal candidate-aware attention.

### 2.1 Multimodal News Encoder

The architecture of the multimodal news encoder is shown in Fig. 2. It takes the image and text of a news as input. On news websites, many news articles have an image displayed together with news title, as shown in Fig. 1. Users may click news not only due to their interests in news title, but also because of the attraction of news images [2]. Thus, modeling the visual content of news such as images is important for news representation learning. In addition, different regions in news image also have different informativeness. For instance, in the image of the first news in Fig. 1, the regions containing Fauci are more informative than the background regions. Thus, for the image part, we use a Mask-RCNN [7] model pre-trained in an objective detection task to extract ROIs of news images. We further use a ResNet-50 [8] model to extract features of ROIs, which forms a feature sequence $[\mathbf{e}_1^p, \mathbf{e}_2^p, ..., \mathbf{e}_K^p]$, where $K$ is the number of ROIs. In the text part, following previous works [22, 25] we use news titles to represent textual news content. We tokenize a news title into a word sequence $[w_1^t, w_2^t, ..., w_M^t]$, where $M$ is the number of words.

An intuitive way is modeling news texts and images with separate models. However, the title and image of the same news usually have some relations with each other. For instance, in the first news of Fig. 1, the word "Fauci" in the news title is related to his photo. Capturing the relatedness between news titles and images can help better understand their content and infer user interests. Visiolinguistic models are effective in modeling the crossmodal relations between texts and images [13, 14, 18, 19]. Thus, we propose to apply a pre-trained visiolinguistic model ViLBERT [14] to capture the inherent relatedness between news title and image when learning representations of them. The input of ViLBERT are the ROI and word sequences. It first models the contexts of words via several Transformer [20] networks, and then use several co-attentional transformers [15] to capture the crossmodal interactions between the image and title. The output is a hidden ROI representation sequence $\mathbf{H}^p = [\mathbf{h}_1^p, \mathbf{h}_2^p, ..., \mathbf{h}_K^p]$ and a hidden word representation sequence $\mathbf{H}^t = [\mathbf{h}_1^t, \mathbf{h}_2^t, ..., \mathbf{h}_M^t]$.

After the ViLBERT model, we use a word attention network to learn unified title representations and an image attention network to learn unified image representations. The attention weights of words in news title are computed as follows:

$$\mathbf{a}^t = \text{softmax}[(\mathbf{W}^t \mathbf{R}^t)^\top \mathbf{q}^t], \tag{1}$$

where $\mathbf{q}^t$ is an attention query vector and $\mathbf{W}^t$ is a parameter matrix. The final representation of news title is the summation of the hidden word representations weighted by their attention weights, i.e., $\mathbf{r}^t = \mathbf{R}^t \times \mathbf{a}^t$. The attention weights of ROIs are computed in a similar way as follows:

$$\mathbf{a}^p = \text{softmax}[(\mathbf{W}^p \mathbf{R}^p)^\top \mathbf{q}^p], \tag{2}$$

where $\mathbf{q}^p$ and $\mathbf{W}^p$ are parameters. The final representation of news image is the summation of hidden ROI representations weighted by their attention weights, i.e., $\mathbf{r}^p = \mathbf{R}^p \times \mathbf{a}^p$.

### 2.2 Multimodal News Recommendation

In this section, we introduce how to make news recommendations using the multimodal news representations learned by the multimodal news encoder. Since personalized news recommendation usually rely on the relevance between candidate news articles and users' personal interest to rank candidate news for a target user, we first introduce the method of user interest modeling. Following many previous works [16, 25], we model users' interest in news from the representations of their previously clicked news articles. We use the multimodal news encoder method to learn the text and image representations of previously clicked news from both news titles and images, which are denoted as $\mathbf{R}^t = [\mathbf{r}_1^t, \mathbf{r}_2^t, ..., \mathbf{r}_P^t]$ and $\mathbf{R}^p = [\mathbf{r}_1^p, \mathbf{r}_2^p, ..., \mathbf{r}_P^p]$, where $P$ is the number of clicked news. However, not all clicked news are informative for inferring user interests on a candidate news, because it may be relevant to a few clicked news only. For example, the candidate news in Fig. 1 is only related to the second clicked news. Thus, selecting clicked news according to their relevance to candidate news in user modeling may be useful for accurately matching candidate news with user interests. In addition, candidate news may have some crossmodal relations with the images and titles of clicked news. For example, in Fig. 1 the players in the candidate news image are related to the players in the images of the second clicked news and the word "Cowboys" in its title. Motivated by these observations, we propose a crossmodal candidate-aware attention network to measure the crossmodel relevance between clicked news and candidate news for better modeling user interests in candidate news. We denote the image and text representations of a candidate news $D_c$ as $\mathbf{r}_c^t$ and $\mathbf{r}_c^p$, respectively. We compute the text-text attention weights for clicked news that represent their text-text relevance to candidate news as follows:

$$\mathbf{a}^{t,t} = \text{softmax}(\mathbf{R}^t \times \mathbf{r}_c^t). \tag{3}$$

In a similar way, we compute the text-image, image-text and image-image attention weights of clicked news as follows:

$$\mathbf{a}^{t,p} = \text{softmax}(\mathbf{R}^p \times \mathbf{r}_c^t), \tag{4}$$

$$\mathbf{a}^{p,t} = \text{softmax}(\mathbf{R}^t \times \mathbf{r}_c^p), \tag{5}$$

$$\mathbf{a}^{p,p} = \text{softmax}(\mathbf{R}^p \times \mathbf{r}_c^p). \tag{6}$$

The final unified user embedding $\mathbf{u}$ is computed as follows:

$$\mathbf{u} = \mathbf{R}^p \times (\mathbf{a}^{t,p} + \mathbf{a}^{p,p}) + \mathbf{R}^t \times (\mathbf{a}^{t,t} + \mathbf{a}^{p,t}). \tag{7}$$

In our method, the score for personalized candidate news ranking is computed based on the multimodal representations $\mathbf{r}_c^t$ and $\mathbf{r}_c^p$ of candidate news and the representation $\mathbf{u}$ of a user. Motivated by [16], the news click score $\hat{y}$ is predicted by $\hat{y} = \mathbf{r}_c^t \times \mathbf{u} + \mathbf{r}_c^p \times \mathbf{u}$.
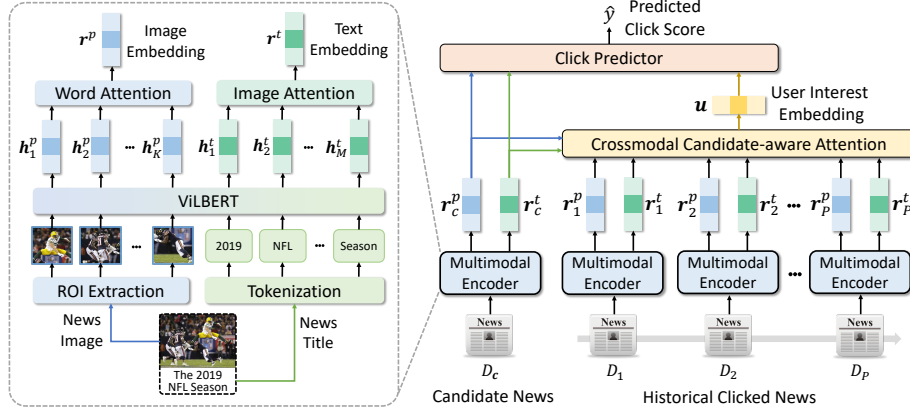
**Figure 2: The framework of *MM-Rec*.**

Following [25] we use negative sampling techniques to build labeled samples from news session logs. Concretely, for each clicked news, $N$ non-clicked news displayed in the same session are randomly selected, and we jointly predict the scores of the $N + 1$ news. These scores are further normalized via the softmax function, and we re-formulate the news click prediction problem as a multi-class classification task, i.e., predicting which news is clicked. We use cross-entropy as the loss function for model training.

## 3 EXPERIMENTS

### 3.1 Datasets and Experimental Settings

Since there is no high-quality dataset that contains multimodal information of news[1], we constructed a dataset based on the logs collected from the Microsoft News website during three weeks (from Feb. 25, 2020 to Mar. 16, 2020).[2] Logs in the first week were used to construct user histories and the rest sessions were used to form click and non-click samples. We sorted these sessions by time, and the first 1,000,000 sessions were used for training, the next 100,000 sessions for validation and the rest for test. The statistics of this dataset are shown in Table 1. We can see that the click-through rate (CTR) of news with images is slightly higher than those without images, showing that news images can better attract news clicks.

**Table 1: Statistics of our dataset.**

| # users | 536,289 | # clicked samples | 1,887,697 |
|---|---|---|---|
| # news | 152,723 | # non-clicked samples | 51,642,426 |
| # sessions | 1,200,000 | CTR of news w/ images | 0.0384 |
| # news w/ images | 111,312 | CTR of news w/o images | 0.0353 |

In our experiments, we finetuned the last three layers of ViL-BERT. We used Adam [12] as the optimizer (lr=1e-5). The negative sampling ratio $N$ was 4. These hyperparameters were tuned on the validation set. We repeated each experiment 5 times and reported the average AUC, MRR, nDCG@5 and nDCG@10 scores.

---

[1] Although there are several public datasets (e.g., Adressa [5]) containing news images, the ratio of news with images is small in these datasets and many URLs for downloading news images are not available now.

[2] Our dataset and codes will be publicly available soon.

### 3.2 Performance Comparison

We compare the proposed *MM-Rec* method with many baseline methods, including: (1) *LibFM* [17], a factorization machine-based recommendation method; (2) *DSSM* [10], deep structured semantic model. We regard clicked news as query and candidate news as documents; (3) *Wide&Deep* [3], a neural recommender with a wide linear part and a deep neural network part; (4) *DeepFM* [6], a neural recommender with factorization machines and deep neural networks; (5) *EBNR* [16], an embedding-based news recommendation method which learns news embeddings via autoencoder and learns user representations with a GRU network; (6) *DKN* [22], learning news representations via a knowledge-aware CNN model; (7) *DAN* [27], learning news representations with two parallel CNN networks from news title and entities; (8) *NAML* [23], using attentive multi-view learning to learn news representations; (9) *NRMS* [25], using multi-head self-attention networks to learn news representations. In these methods, only news texts are considered.

**Table 2: Performance comparison of different methods.**

| Methods | AUC | MRR | nDCG@5 | nDCG@10 |
|---|---|---|---|---|
| LibFM | 57.29±0.43 | 18.08±0.49 | 19.55±0.39 | 27.73±0.44 |
| DSSM | 60.27±0.25 | 20.73±0.22 | 22.34±0.23 | 30.66±0.26 |
| Wide&Deep | 58.83±0.30 | 19.29±0.33 | 21.06±0.31 | 29.31±0.28 |
| DeepFM | 58.96±0.32 | 19.44±0.29 | 21.20±0.34 | 29.45±0.31 |
| EBNR | 60.34±0.29 | 20.79±0.25 | 22.43±0.26 | 30.76±0.23 |
| DKN | 60.18±0.24 | 20.56±0.22 | 22.24±0.20 | 30.53±0.18 |
| DAN | 61.03±0.22 | 21.69±0.19 | 23.12±0.23 | 31.48±0.20 |
| NAML | 61.55±0.18 | 22.13±0.16 | 23.57±0.17 | 31.92±0.17 |
| NRMS | 62.01±0.13 | 22.68±0.15 | 24.08±0.15 | 32.38±0.15 |
| MM-Rec | **64.96**±0.12 | **25.22**±0.11 | **26.67**±0.12 | **34.23**±0.10 |

The experimental results are shown in Table 2. The results show that our *MM-Rec* approach that considers visual information of news outperforms other methods based on textual content only, such as *DAN* and *NRMS*. In addition, the results of t-test validate that the improvement of *MM-Rec* is significant ($p < 1e-3$). This is because users usually click news articles not only based on their interest in news texts, but also the attraction of news images. Thus, the visual information conveyed by news images is very useful for
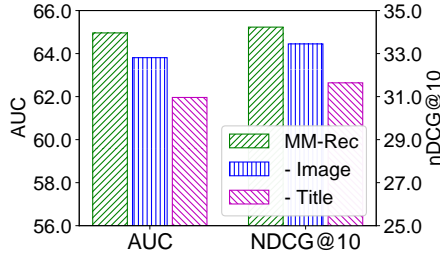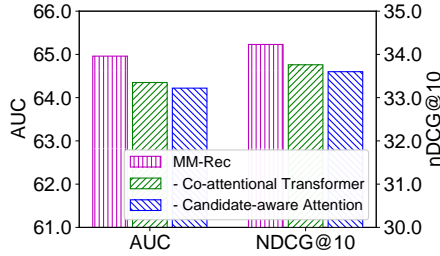
Figure 3: Effectiveness of multimodal information.



Figure 4: Effectiveness of the co-attentional Transformers in ViLBERT and crossmodal candidate-aware attention.



**Figure 5: The clicked news of a user and the rankings of candidate news given by *NRMS* and *MM-Rec*. The first candidate news is clicked while the others are not clicked.**

learning accurate news representations for recommendation. Our *MM-Rec* method can incorporate both textual and visual news information into news representation learning and meanwhile model their inherent relatedness for better news content understanding, while in existing news recommendation methods the image-related information is ignored. In addition, our approach can model the crossmodal relatedness between clicked news and candidate news for more accurate interest matching. Thus, *MM-Rec* can achieve better news recommendation performance.

## 3.3 Ablation Study

In this section, we first study the effectiveness of multimodal information for news representation. We compare the performance of our *MM-Rec* method with its two variants which only consider images or titles in news encoder. The results are shown in Fig. 3. We find that both news title and image are useful for learning news representations for recommendation. It shows that both textual and visual information of news are highly useful for understanding news content and inferring user interest. In addition, incorporating multimodal news information can further improve the recommendation performance, which shows that incorporating multimodal news information can help learn accurate news representations.

We also study the effectiveness of the co-attentional Transformers in the ViLBERT model and the crossmodal candidate-aware attention network for user interest modeling. We compare *MM-Rec* and its variants without co-attentional Transformers or replacing the crossmodal candidate-aware attention with the vanilla attention mechanism used in [23]. The results are shown in Fig. 4. We find that incorporating co-attentional Transformers network is useful. This is because there is inherent relatedness between news title

and image in representing the news content and attracting user attention. Thus, modeling their interactions can collaboratively enhance their representations. In addition, the candidate-aware attention network is also very useful. It is because different clicked news usually have different importance for modeling the user interests in a candidate news, and selecting them according to their crossmodal relatedness with candidate news can help better model user interests for candidate matching.

## 3.4 Case Study

Finally, we conduct several case studies to visually demonstrate the effectiveness of incorporating multimodal information into news recommendation. We show the clicked news of a random user and the rankings given by *NRMS* and *MM-Rec* in Fig. 5. We find that both *NRMS* and *UniRec* assign the last candidate news low rankings, because from its title we can easily infer that it is irrelevant to the user interests. However, the *NRMS* model fails to promote the first candidate news, which is highly related to the user's clicked news about NFL. This may be because it is difficult to measure their relevance solely based on their titles. Fortunately, our *MM-Rec* method ranks the first candidate news at the top position because it is easy to match it with user interests based on visual information. These results show the effectiveness of multimodal information in news recommendation.

## 4 CONCLUSION

In this paper, we present a multimodal news recommendation that can utilize both textual and visual news information for news modeling. We propose to use a visiolinguistic model to encode both news texts and images and capture their inherent crossmodal relatedness. In addition, we propose a crossmodal candidate-aware attention network to select relevant clicked news based on their crossmodal relevance to candidate news, which can better model users' specific interest in candidate news. Experimental results on a real-world dataset show that our method can effectively exploit multimodal news information to facilitate news recommendation.

# REFERENCES

[1] Mingxiao An, Fangzhao Wu, Chuhan Wu, Kun Zhang, Zheng Liu, and Xing Xie. 2019. Neural News Recommendation with Long-and Short-term User Representations. In *ACL*. 336–345.

[2] Monika Bednarek and Helen Caple. 2012. 'Value added': Language, image and news values. *Discourse, context & media* 1, 2-3 (2012), 103–113.

[3] Heng-Tze Cheng, Levent Koc, Jeremiah Harmsen, Tal Shaked, Tushar Chandra, Hrishi Aradhye, Glen Anderson, Greg Corrado, Wei Chai, Mustafa Ispir, et al. 2016. Wide & deep learning for recommender systems. In *DLRS*. ACM, 7–10.

[4] Suyu Ge, Chuhan Wu, Fangzhao Wu, Tao Qi, and Yongfeng Huang. 2020. Graph enhanced representation learning for news recommendation. In *WWW*. 2863–2869.

[5] Jon Atle Gulla, Lemei Zhang, Peng Liu, Özlem Özgöbek, and Xiaomeng Su. 2017. The Adressa dataset for news recommendation. In *WI*. ACM, 1042–1048.

[6] Huifeng Guo, Ruiming Tang, Yunming Ye, Zhenguo Li, and Xiuqiang He. 2017. DeepFM: a factorization-machine based neural network for CTR prediction. In *AAAI*. AAAI Press, 1725–1731.

[7] Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross Girshick. 2017. Mask r-cnn. In *ICCV*. 2961–2969.

[8] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 2016. Deep residual learning for image recognition. In *CVPR*. 770–778.

[9] Linmei Hu, Chen Li, Chuan Shi, Cheng Yang, and Chao Shao. 2020. Graph neural news recommendation with long-term and short-term interest modeling. *Information Processing & Management* 57, 2 (2020), 102142.

[10] 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 *CIKM*. ACM, 2333–2338.

[11] Mozhgan Karimi, Dietmar Jannach, and Michael Jugovac. 2018. News recommender systems–Survey and roads ahead. *Information Processing & Management* 54, 6 (2018), 1203–1227.

[12] Diederik P Kingma and Jimmy Ba. 2014. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980* (2014).

[13] Gen Li, Nan Duan, Yuejian Fang, Ming Gong, and Daxin Jiang. 2020. Unicoder-vl: A universal encoder for vision and language by cross-modal pre-training. In *AAAI*, Vol. 34. 11336–11344.

[14] Jiasen Lu, Dhruv Batra, Devi Parikh, and Stefan Lee. 2019. Vilbert: Pretraining task-agnostic visiolinguistic representations for vision-and-language tasks. In *NeurIPS*. 13–23.

[15] Jiasen Lu, Jianwei Yang, Dhruv Batra, and Devi Parikh. 2016. Hierarchical question-image co-attention for visual question answering. In *NIPS*. 289–297.

[16] Shumpei Okura, Yukihiro Tagami, Shingo Ono, and Akira Tajima. 2017. Embedding-based news recommendation for millions of users. In *KDD*. ACM, 1933–1942.

[17] Steffen Rendle. 2012. Factorization machines with libfm. *TIST* 3, 3 (2012), 57.

[18] Weijie Su, Xizhou Zhu, Yue Cao, Bin Li, Lewei Lu, Furu Wei, and Jifeng Dai. 2019. VL-BERT: Pre-training of Generic Visual-Linguistic Representations. In *ICLR*.

[19] Hao Tan and Mohit Bansal. 2019. LXMERT: Learning Cross-Modality Encoder Representations from Transformers. In *EMNLP-IJCNLP*. 5103–5114.

[20] 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.

[21] Heyuan Wang, Fangzhao Wu, Zheng Liu, and Xing Xie. 2020. Fine-grained Interest Matching for Neural News Recommendation. In *ACL*. 836–845.

[22] Hongwei Wang, Fuzheng Zhang, Xing Xie, and Minyi Guo. 2018. DKN: Deep Knowledge-Aware Network for News Recommendation. In *WWW*. 1835–1844.

[23] Chuhan Wu, Fangzhao Wu, Mingxiao An, Jianqiang Huang, Yongfeng Huang, and Xing Xie. 2019. Neural News Recommendation with Attentive Multi-View Learning. In *IJCAI*. 3863–3869.

[24] Chuhan Wu, Fangzhao Wu, Mingxiao An, Jianqiang Huang, Yongfeng Huang, and Xing Xie. 2019. NPA: Neural News Recommendation with Personalized Attention. In *KDD*. ACM, 2576–2584.

[25] Chuhan Wu, Fangzhao Wu, Suyu Ge, Tao Qi, Yongfeng Huang, and Xing Xie. 2019. Neural News Recommendation with Multi-Head Self-Attention. In *EMNLP-IJCNLP*. 6390–6395.

[26] Fangzhao Wu, Ying Qiao, Jiun-Hung Chen, Chuhan Wu, Tao Qi, Jianxun Lian, Danyang Liu, Xing Xie, Jianfeng Gao, Winnie Wu, et al. 2020. MIND: A Large-scale Dataset for News Recommendation. In *ACL*. 3597–3606.

[27] Qiannan Zhu, Xiaofei Zhou, Zeliang Song, Jianlong Tan, and Li Guo. 2019. DAN: Deep Attention Neural Network for News Recommendation. In *AAAI*, Vol. 33. 5973–5980.