# VICTOR: An Implicit Approach to Mitigate Misinformation via Continuous Verification Reading

Kuan-Chieh Lo
kclo7898@iis.sinica.edu.tw
Academia Sinica
Taipei, Taiwan

Shih-Chieh Dai*
sjdai@utexas.edu
The University of Texas at Austin
Austin, TX, USA

Aiping Xiong
axx29@psu.edu
Pennsylvania State University
University Park, PA, USA

Jing Jiang
jingjiang@smu.edu.sg
Singapore Management University
Singapore

Lun-Wei Ku
lwku@iis.sinica.edu.tw
Academia Sinica
Taipei, Taiwan

## ABSTRACT

We design and evaluate VICTOR, an easy-to-apply module on top of a recommender system to mitigate misinformation. VICTOR takes an elegant, implicit approach to deliver fake-news verifications, such that readers of fake news can continuously access more verified news articles about fake-news events without explicit correction. We frame fake-news intervention within VICTOR as a graph-based question-answering (QA) task, with Q as a fake-news article and A as the corresponding verified articles. Specifically, VICTOR adopts reinforcement learning: it first considers fake-news readers' preferences supported by underlying news recommender systems and then directs their reading sequence towards the verified news articles. To verify the performance of VICTOR, we collect and organize VERI, a new dataset consisting of real-news articles, user browsing logs, and fake-real news pairs for a large number of misinformation events. We evaluate zero-shot and few-shot VICTOR on VERI to simulate the never-exposed-ever and seen-before conditions of users while reading a piece of fake news. Results demonstrate that compared to baselines, VICTOR proactively delivers 6% more verified articles with a diversity increase of 7.5% to over 68% of at-risk users who have been exposed to fake news. Moreover, we conduct a field user study in which 165 participants evaluated fake news articles. Participants in the VICTOR condition show better exposure rates, proposal rates, and click rates on verified news articles than those in the other two conditions. Altogether, our work demonstrates the potentials of VICTOR, i.e., combat fake news by delivering verified information implicitly.

## CCS CONCEPTS

• **Social and professional topics**;

---

---

## KEYWORDS

misinformation, fake news intervention, user research

## 1 INTRODUCTION

With the vast dispatching power of the Internet, disseminating information has become inexpensive and rapid. Likewise, receiving news from the world has become instantaneous. However, users have ended up relying heavily on few choices proposed by recommendation systems that appear on their screen. Most of these recommenders are content-based, click-based, or a hybrid thereof, which point readers in the direction of sensational, popular topics and information similar to what they have seen. The former may spread fake news made to be attractive, and the latter may just reinforce the readers' own beliefs. Indeed, previous studies reveal that the use of news recommendation systems has contributed to the prevalence of misinformation dissemination [19].

When a piece of fake news has been identified, decision-aid interfaces are also proposed to discourage people's consumption or belief in the misinformation, such as on Facebook [24, 30], Google [16], and Twitter [1]. With these approaches, tags, indicators, or labels are typically attached to the information that was suspicious or fact-checked to be fake when proposed by the recommendation system. Nevertheless, empirical user studies reveal mixed results of such warnings in mitigating fake news. Although the effect of warnings is evident in some studies [4], it is limited or absent in others [8, 21, 25]. In addition, unintended or potential harmful effects of warnings have also been identified in the literature, including false negatives (i.e., the implied truth effect) [22], warning habituation [2], and adverse effects (e.g., backfiring) due to individual differences [9].

In the web environment, it is common that we know either the fake news or its real version, but linking them together requires considerable human effort. In addition, once a story is confirmed to be fake, it is quickly removed from the media or the platform[1].
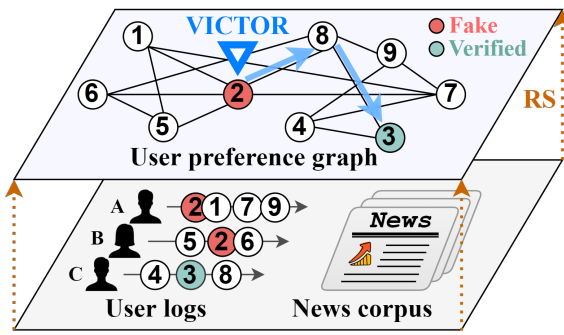
---

**Figure 1: The diagram shows a general picture of VICTOR. Users A and B are potentially at-risk users since they have read a piece of fake news article. To preserve users' reading preferences, we construct the user preference graph by the results of a recommender system, which is trained by the user logs and the news corpus. With VICTOR's intervention, a recommendation path between the fake news and the possible verification will be planned for the at-risk users to follow.**

In this case, there is little chance that users who have already been exposed to the fake news will come back again to see later attached warnings (if any) as recommenders naturally avoid duplicate recommendations, or because the story no longer exists. Moreover, it is rare that the recommender knows where the verification is, and even fewer systems are equipped with the ability to propose this to the user, not to mention the possible side effects of clumsy verification such as backfiring. In the current web environment, such exposed users are at risk of viewing more similar misinformation. However, no treatment has been proposed or implemented in the current web environment for those at risk users.

To better help these users mitigate the impacts of misinformation, we turn to effective implicit intervention. In this paper, we study misinformation intervention on top of the recommendation system. We propose VICTOR (**V**erified **I**n **C**ontinuous Reading **TO**ward **R**eal News), a novel intervention approach that gradually lays out the user's reading path to reach verified news only after she or he reads a fake news article (see Figure 1). These users are defined as potentially at-risk users and are the target audience of the proposed method. VICTOR takes known fake news articles as guidance. Once the system is made aware that a user has been exposed to any of fake news articles, the intervention mechanism is automatically activated. To smooth the intervention process, we maintain users' reading preferences by constructing the user preference graph using the result of a recommender system (RS). Each node in the graph represents an individual news article, and connected neighbors are the news articles recommended by the RS. We use reinforcement learning (RL) to uncover potential verified news that might be several steps away from the fake news in the graph, and lay out the path from the fake news towards the verified news. In this process, the RL model decides a recommendation list from the neighbors at each step. Consequently, we frame this decision process as a question-answering (QA) task, where the question is the fake news,

and the answer is the verified news [5, 35]. To evaluate the real-world feasibility of VICTOR, we conduct extensive experiments on a large quantity of news articles as well as a user simulation study. Results from both aspects show that VICTOR yields better performance compared to other baseline methods for fake news intervention. The key contributions of this paper include:

- We propose VICTOR, the first implicit fake news intervention method that avoids the negative effects of current explicit interventions methods.
- All of the experiment, the offline simulation, and the user study results illustrate that VICTOR successfully guides the reading path from fake news towards diverse verified news by intervening in the recommendation results.

## 2 RELATED WORK

In this section, we review previous studies related to VICTOR. In recent years, much progress has been made on fake news detection [17, 28, 29, 32]. However, little work has been conducted to effectively deliver the detection results to users [25]. Here, we first review common solutions to tackle fake news in the real world, and focus on techniques to mitigate the influence of fake news and related fields such as misinformation and rumor.

Some researchers tackle fake news from a proliferation aspect. They attempt to detect malicious accounts at an early stage as the source of misinformation, and once the malicious account is detected, to ban the account. Such a method combats the spread of misinformation [11, 27, 34] with the advantage of stopping the spread of fake news at an early stage.

Another idea is to "sanitize" users exposed to fake news. The intuition of this method is to provide him/her with real news as the sanitizer [29]. Nguyen et al. [20] use an independent cascade or a linear threshold model to model the diffusion process and limit the propagation of misinformation. Farajtabar et al. [7] propose a fake news intervention framework that models the spread of fake news and mitigates events via a multivariate Hawkes process [13] combined with an RL model. This framework restricts the spread of fake news and optimizes the spread of real news. Goindani and Neville [10] adopt reinforcement learning to learn a fake news intervention model to promote the spread of real news. As mentioned above, such sanitation has the power to limit the spread of fake news and optimize real news propagation.

An alternative method is crowd-sourcing, specifically crowd-sourcing feedback to label fake news. Kim et al. [15] leverage the crowd to determine whether the news should be verified. Users report news as fake news, and once it is reported as fake by a sufficient number of users, it is flagged for verification. In another study, crowd-sourcing is used to judge the quality of the media source. Pennycook and Rand [23] conduct crowd-sourcing to rate the trustworthiness of the media outlet. Their results suggest that incorporating the trust ratings of laypeople into social media ranking algorithms may prove an effective intervention against misinformation with heavy political bias. In general, crowd-sourcing has the benefit of eliminating the high cost of labeling fake news.

However, these approaches have their limitations. For example, malicious account detection is liable to false positives, e.g., regular accounts may be deleted incorrectly. Most sanitization methods

| Condition | Fake only | Real only | Fake and Real |
|---|---|---|---|
| # Users | 445 | 1275 | 65 |

**Table 1: Number of users who read only known fake news articles, read only verified (real) news, and who read both.**

| VERI | | |
|---|---|---|
| News corpus | Stories | 1,481,125 |
| Fake-news events | Events | 570 |
| | Fake articles | 1198 |
| | Real articles | 2649 |
| User browsing logs | Users | 485,522 |
| | At-risk users | 510 |
| | Risky sequences | 541 |

**Table 2: Statistics of VERI dataset.**

are point process-based methods [13], and it is difficult to set the parameters for such methods; inappropriate parameter settings easily ruin performance [26]. Likewise, crowd-sourcing involves a stark trade-off between collecting more samples for fake news and exposing more users to potential fake news.

To tackle these issues and mitigate the impact of fake news, we propose VICTOR, an implicit intervention approach which has none of these shortcomings, as there is no need to label fake news during intervention, much less show the labels. For example, VICTOR does not involve human effort; there is no risk of exposing more users to fake news as labeling is now superfluous; and the omission of labels could eliminate the implied truth effect.

## 3 TASK DEFINITION

From real-world news websites[2], we manually collected 570 fake news events—the news events, known fake news articles, and corresponding verified news—over 31 months (the collection method is introduced in the next section). We then traced back to the users' browsing logs on the news websites, and observed that most online news consumers were exposed to either fake news (445, 25%) or their verified news (1275, 71%), but very few of them skimmed over both (65, 4%) (see the statistics in Table 1). We define those 510 users (445 read fake news only, and 65 read both fake news and verifications), who ever read fake news articles, as the *at-risk users*, and their browsing logs as *risky sequences*. We further discuss these at-risk users and risky sequences below; the remaining 1275 users are not the targets in this paper. As a result, we find that although current commercial news recommenders are skilled at maintaining users' reading preferences, they fail to bring up nuanced information about a fake-news event from a different perspective. In other words, once users have consumed a piece of fake news, they could be vulnerable to more exposure to similar misinformation.

To avoid harming users' reading preferences, and reveal timely true information to mitigate users' belief in fake news, the goal of the task is to intervene implicitly and direct contaminated users to the verified news in their reading session. To begin with, we collected the VERI dataset, which consists of the news corpus, the labeled fake-real news pairs, and the users' browsing logs. To better maintain users' reading preferences and simulate a realistic reading environment, we utilize a recommender system, trained with the logs presented by VERI, and construct a user preference graph by the recommender system's results. To present the verified news to the fake news readers, we use VICTOR, a module to intervene in the recommendations within users' reading sessions. VICTOR begins to operate when users are reading either a debunked fake news story or any suspicious news detected by any mechanism.

When VICTOR is activated, it attempts to guide the users to reach possible verifications with regard to the fake news.

In this task, we focus on evaluating whether VICTOR successfully delivers verifications to the fake news readers. Detecting whether users have read a piece of fake news is out of the scope of the current work. Instead, we directly feed a fake news article as VICTOR's input, indicating that the users' exposure to a piece of fake news has been detected. We describe the details of the dataset and methodology in the following sections.

## 4 THE VERI DATASET

To tackle this task, we require a dataset including (1) a corpus of news articles with headlines and bodies; (2) event-based fake-real news pairs; (3) user browsing logs to establish a recommender system. However, to the best of our knowledge, there is no existing dataset that meets these requirements. Therefore, we constructed VERI, a new dataset which consists of three major components: a news corpus in Chinese, a set of fake-news events, and a set of user browsing logs. The dataset statistics are shown in Table 2. We explain each component in detail below.

### 4.1 News Corpus

All news articles in the news corpus were collected over a period of 31 months—from April 2018 to October 2020—from four news websites[3]. Each article included a headline, body and was numbered by a unique ID. In total, the corpus contains 1,481,125 news articles.

### 4.2 Fake-news Events

Within the period of the news corpus collection, we manually collected 570 fake-news events, each of which includes at least one fake news article and corresponding verified news articles. We strictly filtered the fake news and verified news articles by checking the governmental and credible FactCheck websites. Altogether, there are 1198 fake news articles and 2649 verified news articles.

### 4.3 User Browsing Logs

From the commercial news websites, we collected up to 485,522 users' browsing logs. With the collected fake-news events and the logs, we labeled those users who had read fake news articles as *at-risk users*. We then gathered these at-risk users' browsing logs into sequences starting from a piece of fake news they had read.
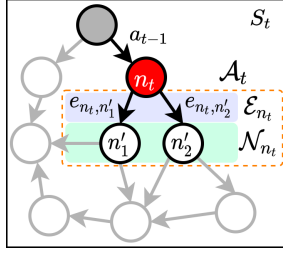
Figure 2: Illustration of state at time $t$ ($S_t$). $n_t$ denotes the currently visited node; $a_{t-1}$ indicates the previous selected action; $\mathcal{E}_{n_t}$ and $\mathcal{N}_{n_t}$ show the set of all edges and all nodes connected to $n_t$ respectively.

More than one sequence could be extracted for one at-risk user since they could read multiple fake news articles. We define these extracted sequences as *risky sequences*. In sum, we identified 510 at-risk users and 541 risky sequences in the dataset.

## 5 METHODOLOGY

### 5.1 Model Overview

VICTOR is a reinforcement learning (RL)-based model which walks a user preference graph to find paths that link fake news articles and their corresponding verifications. The starting point of VICTOR can be either a piece of debunked fake news or a piece of suspicious news by any fake news detection mechanism. When users reach the starting point, VICTOR is activated by taking the piece of fake news or the suspicious article as its input. In subsequent recommendation steps, VICTOR follows the learned RL policy based on current state information—the read news articles, the input fake news article, and all candidate articles provided by the recommender system—to rearrange the recommendations. Via such intervention, VICTOR attempts to guide the users to reach potential verified news articles. After ten steps, the intervention ends and control is handed back to the underlying recommender system.

### 5.2 User Preference Graph

In this work, we attempt to maintain the user's reading preferences. Therefore, we utilized a content-based news recommender system [3], which was well-trained on VERI's news corpus and user browsing logs to construct the user preference graph according to the results. Each node in the graph represents a news article, and each directional edge connecting two nodes indicates that the tail node is a recommendation of the head node suggested by the recommender system. There are 607,332 articles involved in the graph. We consider the top-20 recommendations of each article. That is, the out-degree of every node is equal to 20.

### 5.3 Reinforcement Learning for Intervention

We frame the intervention process as a graph-based question-answering problem. Given a user preference graph ($\mathcal{G}$), a question ($n_0$, a fake news article), and a query ($q$, "*isFakeNewsOf*"), the goal is to make a series of decisions walking on the graph to reason a path from the question to a target ($\widehat{n}_T \in \{n_{T1}, n_{T2}, ..., n_{Tm}\}$, any one
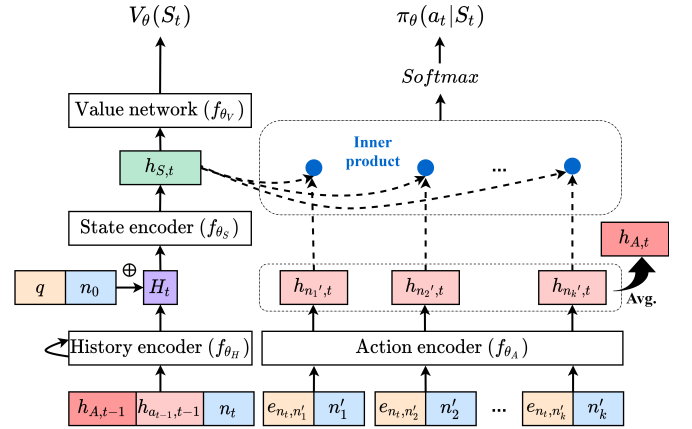


Figure 3: The action encoder encodes all action candidates, and the history encoder encodes earlier steps' decisions. The current state ($h_{S,t}$) is encoded by the state encoder with the history representation ($H_t$), the question ($n_0$) and the query ($q$) information. All candidate actions are rated by conducting inner product to current state. Finally, the value network estimates the accumulated reward value.

article of the corresponding $m$ verifications). We formulate the decision process as a Markov Decision Process (MDP), defined by the tuple ($\mathcal{S}, \mathcal{A}, \delta, R$), each of which denotes states, actions, transitions, and rewards in the procedure.

To make better decisions at each time step, we usually need not just the query, but also the entire history of traversed nodes. We follow Das et al. [5], Xian et al. [35] in adopting a history-dependent policy for searching decisions. The $S_t \in \mathcal{S}$ denotes the state at time $t$ and is defined by the following recursion:

$$S_t = S_{t-1} \cup \{a_{t-1}, n_t, \mathcal{E}_{n_t}, \mathcal{N}_{n_t}\}, \tag{1}$$

where $a_{t-1}$ denotes the previous selected action; $n_t$ denotes the currently visited node; and $\mathcal{E}_{n_t}$ and $\mathcal{N}_{n_t}$ indicate the set of all edges and the set of all nodes connected to $n_t$ respectively (see Figure 2). We specify the initial state as $S_0 \triangleq \{q, n_0, \mathcal{E}_{n_0}, \mathcal{N}_{n_0}\}$, where $q$ and $n_0$ represent the query (i.e., "*isFakeNewsOf*") and the question (i.e., the fake news article) respectively. For each state $S_t$, all possible outgoing edges $e_{n_t, n'}$ and the connected neighbors $n'$ compose a set of candidate actions $\mathcal{A}_t = \{(e_{n_t, n'}, n')\} \in \mathcal{A}$. The RL agent, in each state, decides an action outweighs others, then transits the state by equation (1) (with $t$ replaced by $t + 1$). After transition, the agent receives a corresponding reward $R_t$. If the current location $n_t$ reaches any of the targets in $\widehat{n}_T$ (i.e., a set of verified news articles), a positive reward of +1 is given; otherwise, the reward is 0. The cumulative discounted reward is $G_t = \sum_{k=0}^{T-t-1} \gamma^k R_{t+k+1}$, where $T$ is the final time step (in this work, we set $T = 10$), and $\gamma$ is a parameter, $0 \leq \gamma \leq 1$, termed the discount rate [31]. The agent's goal is to maximize the $G_t$ received in the sequential decision process.

## 5.4 The VICTOR Agent

To better capture the semantic information of the news articles, we first encode every node $n$ in the user preference graph $\mathcal{G}$ by the pre-trained language model, BERT [6]. For all edges in $\mathcal{G}$, we consider the single relation of "*REC*", which shows the recommendation relationship between two connected nodes. We randomly initialize the representations of the edge (with the relation of "*REC*") as well as the query (with the relation of "*isFakeNewsOf*") into the vectors with a dimension of 100.

The complete network architecture of VICTOR is shown in Figure 3. For the $n'$-th candidate action of $S_t$, we concatenate the representation of the edge $e_{n_t, n'}$ with the associated representation of the node $n'$, and feed them into the action encoder ($f_{\theta_A}$) to obtain the $n'$-th action representation $h_{n', t}$, where $f_{\theta_A}$ is a fully-connected neural network with the model parameter $\theta_A$. We use a vector $h_{A, t}$ for the fusion representation of all candidate actions of $S_t$ by calculating element-wise average-pooling over all candidate actions' representations. For better decision making, we adopt a history-dependent policy network that takes the traversed nodes into account via a history encoder:

$$H_t = f_{\theta_H}(H_{t-1}, [h_{A, t}; h_{a_{t-1}, t-1}; n_t]), \qquad (2)$$

$$H_0 = f_{\theta_H}(q, [0; 0; n_0]), \qquad (3)$$

where $H_t$ denotes the history representation at time $t$, $f_{\theta_H}(\cdot)$ is a LSTM neural network encoder [14] with the model parameter $\theta_H$, $h_{a_{t-1}, t-1}$ indicates the chosen action's representation at time $t-1$, and [;] denotes vector concatenation. For the initial history representation $H_0$, $n_0$ and $q$ denote the question (i.e., the fake news) and the query (i.e., "*isFakeNewsOf*") representations.

To keep the policy network focused on the problem at hand when it is making decisions, we concatenate the question ($n_0$) and the query ($q$) representations with the history representation, and feed them into a state encoder to get the state representation: $h_{S, t} = f_{\theta_S}([q; n_0; H_t])$, where $f_{\theta_S}(\cdot)$ is a fully-connected neural network with the model parameter $\theta_S$. The candidate action probabilities are determined according to

$$u_{n'} = \langle h_{S, t}, h_{n', t} \rangle, \quad n' \in \mathcal{N}_{n_t}, \qquad (4)$$

$$\pi_\theta(\cdot \mid S_t) = \text{softmax}(u'_{n_1}, u'_{n_2}, ..., u'_{n_k}), \qquad (5)$$

where $\langle \cdot, \cdot \rangle$ denotes the inner product; $\theta$ is the parameter set of the policy network; and $k$ in this work is 20, since we consider the top-20 recommendations of each article in $\mathcal{G}$.

To encourage VICTOR to find diverse verified news articles within the traversal trajectory, we feed the state representation into a value network $f_{\theta_V}$ to obtain the value of $V_\theta(S_t)$, which is used to estimate the reward value accumulated so far. We jointly train the value network and the policy network to help VICTOR better estimate which action will achieve greater rewards.

## 6 EXPERIMENTS

Here we evaluate the verified news delivery performance of recommender systems with and without the intervention. For all intervention-aided cases, in the experiments, we considered only the optimal decisions made by the RL agents, and disregarded the factor of user engagement. That is, according to the policy network, the selected articles generally outweigh other candidates at

each time step. Starting from a fake news story and ending after ten actions are taken, we assess the delivery performance of such trajectories as the optimal performance that the intervention can achieve. We further discuss situations that account for user engagement in the "User Simulation" section.

## 6.1 Experimental Settings

Regarding the publishing time of fake news articles and the time the intervention model is trained, we observe three conditions: *inside test*, *few-shot learning*, and *zero-shot learning*, from simple to difficult. These conditions encompass most of the scenarios that could happen in reality. We split the VERI dataset by three different settings to reflect the conditions. Below, we describe the scenarios of these conditions and the dataset-splitting methods respectively.

**Inside Test.** The fake news event is outdated, and no further associated fake news stories are published, where the intervention model has completely seen related news in the training process. In this scenario, we conduct an inside test for the model. All fake news articles in the testing set exist in the training set as well. We monitor the upper bound on performance that the model reaches when all the related fake news have been explored in training.

**Few-shot Learning.** The event has continued for a period of time and is still active; thus the intervention model saw some of the fake news when it was training. In this scenario, we apply few-shot learning, where the model has already seen a few fake news articles of some events and must now handle unseen fake news about these seen events during testing. Here we presume that all events in the dataset have been partially seen. For every event, we divide associated fake news articles into 60%, 10%, and 30% for training, validation, and testing respectively.

**Zero-shot Learning.** The event has just come up, and no former published fake news has been seen by the intervention model In this scenario, we apply zero-shot learning, where in the testing phase, the model tackles unseen news events. Here we divide the dataset according to the ratio by fake news events. That is, 60% of the events among a total of 570 events and their news are for training, 10% are for validation, and 30% are for testing.

## 6.2 Models

For the intervention-aid cases, we compared the results of VICTOR with two other representative RL-based models. For cases without intervention aid, on the other hand, we simply conducted random selects as the baseline.

**1. Random Select (RAND).** We conducted a random select as a naive baseline without any intervention, i.e., randomly select the next-to-read news article from the recommender system's results.

**2. REINFORCE (PG).** REINFORCE is a vanilla policy gradient (PG) algorithm which at each time step considers only the current article's representation as the state to make the decision.

**3. MINERVA (MIN).** MINERVA is a benchmark model proposed for knowledge base completion [5, 35]. Similar to VICTOR, it adopts a history-dependent policy to help the policy network remember the historical decision path when exploring in knowledge graphs. Since MINERVA was originally designed for finding an answer to the given question, the traversal can end by taking a special action

Optimized actions taken

| Model | PGA | | | | PLM-NR | |
| --- | --- | --- | --- | --- | --- | --- |
| | RAND* | PG | MIN | VIC | RAND* | VIC |
| Inside test | | | | | | |
| Reach@1 | 6.13% | 38.44% | **40.94%** | <u>39.83%</u> | 8.63% | **43.45%** |
| Reach@5 | 13.37% | 57.38% | 64.90% | **71.87%** | 16.43% | **74.65%** |
| Reach@10 | 13.37% | 60.45% | 66.30% | **73.82%** | 17.83% | **82.45%** |
| Diversity | 1.167 | 1.313 | <u>1.424</u> | **1.777** | 1.266 | **1,791** |
| Few-shot learning | | | | | | |
| Reach@1 | 6.96% | <u>27.58%</u> | 27.02% | **28.13%** | 9.19% | **31.48%** |
| Reach@5 | 14.21% | 44.29% | 45.96% | **49.58%** | 15.04% | **47.63%** |
| Reach@10 | 14.48% | 45.96% | 47.63% | **50.70%** | 18.11% | **55.15%** |
| Diversity | 1.135 | 1.455 | <u>1.468</u> | **1.511** | 1.2308 | **1.576** |
| Zero-shot learning | | | | | | |
| Reach@1 | 5.57% | 15.32% | <u>17.55%</u> | **20.61%** | 10.31% | **27.58%** |
| Reach@5 | 13.09% | 21.73% | 30.36% | **31.75%** | 17.27% | **33.98%** |
| Reach@10 | 14.21% | 24.79% | 30.92% | **35.38%** | 18.66% | **37.33%** |
| Diversity | 1.078 | 1.202 | <u>1.207</u> | **1.228** | 1.254 | **1.306** |

**Table 3: Experiment results for inside test, few-shot learning, and zero-shot learning. "*" denotes no intervention aid.**

of "NO_OP" when it reaches a possible answer. We remove this special action to proceed with the traversal.

**4. VICTOR (VIC).** This is the proposed model described in the Methodology section.

In the experiments, we evaluate aforementioned models upon two recommender systems **PGA** [3] and **PLM-NR** [33] to demonstrate VICTOR can work on different recommender systems.

## 6.3 Experimental Results and Findings

Table 3 shows the experimental results of the different models tested in different conditions. To measure different models' ability to deliver verifications, we utilize two criteria—*Reach@K* and *Diversity*—to assess the achievement rate and the fixation on proposing the verfications respectively. Reach@K is the reach rate, which indicates the ratio of paths that achieve the goal within $K$ steps. That is, it indicates the ratio of fake news in the testing sets that can access any of the associated verifications within $K$ decisions made by the model. Diversity shows the average number of different verifications reached in a successful path, that is, a path in which the goal is achieved. Paths with higher diversity imply that the model, in the best case, brings more corrections from different sources for the misinformation. We masked all traversed articles in the action space to prevent models from getting stuck at the selected articles when making decisions. That is, all picked articles in each trajectory are unique. For few-shot and zero-shot learning, we conducted 5-fold cross validation and present the macro performance in Table 3. The results are summarized in the findings below.

**Finding 1: Relying only on the recommender system rarely achieves verification.** The consistently poor performance of the random select baseline (RAND), and the fact that all intervention-aided cases outweigh RAND in all three conditions, proves conclusively that naively relying on the recommender system without intervention rarely leads readers from fake news to verified news.

**Finding 2: The history-dependent policy network helps VICTOR and MINERVA reason between fake news and verifications.** In the most optimal case where the model has already seen all of the fake news in the training process (the inside test), and with the best decisions, VICTOR successfully directs the subsequent path of 73.82% of fake news to their verifications within ten steps. It outperforms RAND (13.37%) and beats other RL-based baselines easily: PG (60.45%) and MINERVA (66.30%). As the difficulty increases (few-shot and zero-shot learning conditions), the weak baseline (PG) cannot maintain its reach rate (45.96% in few-shot learning and 24.79% in zero-shot learning). In contrast, VICTOR and the strong baseline (MINERVA) maintain strong performance, indicating that the history-dependent policy network results in better reasoning trails between fake news and potential verifications.

**Finding 3: The value network of VICTOR contributes to finding diverse verifications.** Note that although MINERVA yields an acceptable reach rate in all conditions, it fails to propose diverse verifications, as its diversity is consistently less than that of VICTOR. If we further jointly consider the proposal rate and the diversity, we find that both MINERVA and VICTOR have relatively high proposal rates. Given that the recommender does not recommend articles that have already been read by the user, the high proposal rate and the low diversity of MINERVA indicate a constant recommendation of the same verified news which was not clicked by the user, whereas the high proposal rate plus high diversity shows that VICTOR recommends a variety of verified news stories to users whether they have clicked it or not. This consistent attempt toward verification by VICTOR is indeed a merit for our task.

## 7 USER SIMULATION

To evaluate in a more realistic scenario, we constructed a virtual reading environment to simulate users clicking on recommendations. In this reading environment, we set the starting point to a piece of fake news. The top-5 recommendations were provided via different models (PG, MINERVA, and VICTOR for intervention-aided cases, and recommender system results for the non-intervention-aided cases). One of the five was automatically selected to mimic different users' decisions on the next news article to read. Each simulation path was complete after ten news articles were selected. Following the settings in the experiments, we also considered the three conditions (inside test, few-shot learning, and zero-shot learning) in the simulation. For each fake news story (the starting point) in the testing set, we ran 80 simulations to present 80 randomly selected users experiencing the virtual reading environment starting from the piece of fake news.

## 7.1 Simulation Results and Observations

Here we additionally introduce the exposure and proposal rates to probe behavior differences across virtual reading environments utilizing different models within the simulation. The *exposure rate* shows the ratio of simulation paths were proposed any of the verifications regarding the fake news within the recommendations,

80 simulations with randomly selected actions

| Model | PGA | | | | PLM-NR | |
|---|---|---|---|---|---|---|
| | RAND* | PG | MIN | VIC | RAND* | VIC |
| Inside test | | | | | | |
| Exposure rate | 34.68% | 56.40% | 57.04% | **66.02%** | 37.51% | **68.13%** |
| Proposal rate | 1.24% | 2.66% | 2.81% | **4.70%** | 1.43% | **4.88%** |
| Diversity | 1.793 | 2.358 | 2.461 | **3.560** | 1.901 | **3.583** |
| Few-shot learning | | | | | | |
| Exposure rate | 34.46% | 51.43% | 53.98% | **59.31%** | 37.02% | **62.97%** |
| Proposal rate | 1.26% | 2.49% | 2.61% | **3.59%** | 1.42% | **3.90%** |
| Diversity | 1.831 | 2.425 | 2.419 | **3.024** | 1.913 | **3.094** |
| Zero-shot learning | | | | | | |
| Exposure rate | 32.18% | 48.12% | 49.97% | **54.53%** | 37.11% | **56.4%** |
| Proposal rate | 1.18% | 2.03% | 2.12% | **2.61%** | 1.41% | **2.92%** |
| Diversity | 1.827 | 2.113 | 2.121 | **2.398** | 1.899 | **2.509** |

**Table 4: Simulation results for the inside test, few-shot learning, and zero-shot learning. "*" denotes no intervention aid.**

and the *proposal rate* indicates the total proposed verified news stories among all recommendations in all of the simulation paths. Following the experiment, we use *diversity* to demonstrate the average number of different verified news being proposed in each simulation. Table 4 shows the simulation results. As with the experiment, for few-shot and zero-shot learning, we conducted 5-fold cross validation and present the macro performance in the table.

**Finding 1.** Table 4 shows high exposure rates for VICTOR, even in the most difficult condition (zero-shot learning): the exposure rate of VICTOR (54.53%) exceeds that of MINERVA in the few-shot learning condition (53.98%). The simulation results, with a high variance of user engagements, suggest that VICTOR is the most stable model for fake news intervention, in which more VICTOR users who are exposed to fake news are then exposed to potential verified news than in other baselines.

**Finding 2.** VICTOR's higher proposal rate and verification diversity in Table 4 indicate that VICTOR more frequently proposes verifications from different sources to users. In comparison, MINERVA does achieve relatively high proposal rates (e.g., 2.81% in the inside test). However, its low diversity (2.461) implies that MINERVA keeps on proposing the same verifications to the users until they click on it. This confirmation bias behavior could potentially decrease users' interest in reading the verified news.

### 7.2 Simulation on Risky Sequences

The VERI dataset contains 541 risky sequences (see Table 2) from 510 at-risk users who rely on the recommender system in a real news media platform yet only read fake news articles without reaching verifications. Every sequence records a segment of one real user's browsing history, starting from a piece of fake news. Given VICTOR's superior performance in the simulations, here we test the possibility of a user being exposed to verifications (exposure rate) when walking through the sequences with VICTOR's intervention.
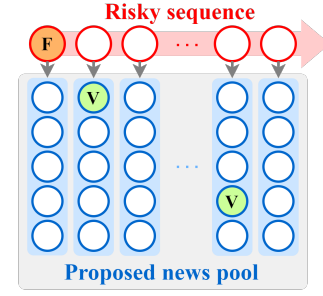


**Figure 4: VICTOR aid in risky sequences. "F" indicates a piece of fake news, and "V" indicates the verified news.**

In this experiment, we follow the user's actual footprints. That is, we move forward by taking the next step recorded in the sequence. At each step, the top-5 recommendations are provided by VICTOR. All recommendations are listed in a proposed news pool. As in the experiments, all traversals on the sequences are terminated after ten steps. We evaluate whether the at-risk users being exposed to the verified news articles by checking the existence of veifications in the proposed news pool. (see Figure 4). After doing so, the verified news articles are proposed to 347 out of 510 (68.04%) at-risk users in 377 out of 541 (69.69%) risky sequences.

## 8 USER STUDY

We have shown the potentials of VICTOR using a user simulation. Next, we present a user study to examine the effectiveness of VICTOR's intervention in reducing users' susceptibility to COVID-19 misinformation in a *field* setting. In particular, we compare the VICTOR's intervention with the other two conditions: a recommender system (*RS*, non-intervention-aid case) and a recommender system with a warning tag to the fake news articles (*RS-tag*, a common intervention method).

### 8.1 User Study Method

We posted our study to questionnaire-filling groups on social media platforms, e.g., Facebook and LINE. We sent the link to our study website to 296 participants interested in our study. The experiment took four days to complete. On average, participants spent about 30 minutes each day on the study. Each participant who completed the experiment received 20 US dollars as payment. Participants could
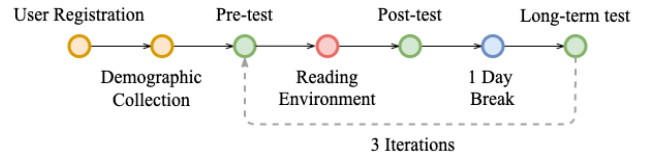


**Figure 5: The flowchart of field study. Orange dots depict pre-study phases, green dots represent news veracity evaluation, the red dot shows the reading environment experiencing, and the blue dot indicates a one-day gap. Three test iterations were conducted for each participant, corresponding to the three conditions.**

contact us via e-mail or the phone number posted on the flyer for any questions throughout the study. Our study was approved by the IRB. The participation of the study was anonymous and voluntary.

**Study Process.** Using a within-subject design, each participant evaluated all three conditions (i.e., model settings: VICTOR, RS, RS-tag) across four days. After registering on the website and completing the demographic survey, participants were randomly assigned to one of the three conditions (see Figure 5). In each condition, there were three test phases (i.e., pre-test, post-test, and long-term test). We implemented three different target (fake) events across the three model settings. For each target event, we constructed three test cases that were used across the three test phases, respectively.

In the pre-test phase, participants made veracity judgments about four pieces of fake news and three pieces of real news. The pre-test thus evaluated whether participants were able to identify fake news before any intervention. The reading environment in the user study was a replica of the offline simulation. In the reading environment, the fake news article about the target event was presented. The recommended articles were accompanying listed, each of which was clickable to read. After going through the reading environment in the assigned condition, participants evaluated the veracity of another news set (4 fake and 3 real) immediately and the next day, respectively. During the following two days, each participant went through the other two conditions in the same way as the first condition. The order of the three conditions and the dispatch order of the test cases among the conditions and test phases were presented using a Latin square design.

Since reading articles is critical to our user study, we implemented an attention check in each test phase [12]. Specifically, a real news article unrelated to any events in each test case was presented and participants were asked to answer the specified correct answer based on the instructions. Altogether, we implemented nine sets of COVID-19 related news articles. Each set included eight news articles: four fake and four real.

**Evaluation Measures.** In addition to exposure rate and proposal rate evaluated in the offline simulation, we evaluated participants' news veracity judgment using the signal detection theory [18]. We compared the participants' *sensitivity* ($d'$) and *response bias* ($c$) between the conditions of VICTOR and RS-tag. Specifically, we set the signal as the fake news, and noise as the verified news. We calculate the sensitivity by $d' = z(H) - z(FA)$[4], showing participants' ability to differentiate fake news from real news. We also calculate the response bias by $c = -0.5[z(H) + z(FA)]$, which points out the tendency of the participants to believe a piece of news article is real or fake.

## 8.2 User Study Results and Analysis

We included results of 165 valid participants in the data analysis. Most of the participants were female, had at least a bachelor's degree, and were below 60 years old (see Table 5 for the details).

**Finding 1.** Table 6 first two columns show the exposure rate and proposal rate of the target events' verifications for each condition in the user study reading environment. VICTOR outperforms both

---

| Education | High School- | Bachelor | Master+ |
|---|---|---|---|
| | 10.91% | **63.03%** | 26.06% |

| Age | 20–24 | 25–44 | 45–59 | 60+ |
|---|---|---|---|---|
| | 26.76% | **34.51%** | 27.46% | 11.27% |

| Gender | Male | Female |
|---|---|---|
| | 32.72% | **67.28%** |

**Table 5: Demographic information of the valid participants (n=165). Bold values show the majority in each category.**

| Condition | Exposure rate | Proposal rate | Click rate |
|---|---|---|---|
| VICTOR | **89.09%** | **13.75%** | **43.46%** |
| RS | 84.24% | 11.82% | 34.10% |
| RS-tag | 83.03% | 10.26% | 33.75% |

**Table 6: Statistics of model performances**

| Criterion | Method | Pre-test | Post-test | Long-term |
|---|---|---|---|---|
| Sensitivity ($d'$) | VICTOR | 0.220 | 0.524 | 0.617 |
| | RS-tag | 0.187 | 0.605 | 0.543 |
| Response bias ($c$) | VICTOR | 0.059 | -0.045 | -0.058 |
| | RS-tag | 0.036 | -0.043 | 0.004 |

**Table 7: The participants' sensitivity and response bias in different phases with different intervention methods. The larger $d'$ indicates the participants are more sensitive to the fake news; a lower $c$ implies the participants retain skeptical of the news articles' authenticity.**

RS-based conditions in delivering verified news to users. Moreover, participants revealed an increased click rate in the VICTOR condition (about 10%), suggesting that verification according to users' preference might increase their intention to read the verification.

**Finding 2.** As shown in Table 7, participants' in both VICTOR and RS-tag conditions increased their sensitivity and reduced their bias toward believing news articles are real in the post-test (i.e., the short term). Critically, in the long-term test after one day break, participants in VICTOR maintained or somewhat increased their sensitivity and vigilance to fake news. In contrast, participants in RS-tag reduced their sensitivity and showed the bias to believe news as real again. Altogether, those results indicate that VICTOR outperforms adding warning tags to the fake news (the most common intervention method) in enhancing users' ability to distinguish fake and real news in the long term.

## 9 CONCLUSION

We propose a novel model VICTOR and a new VERI dataset to expose users to verified information and encourage them to read it continuously. Outperforming model evaluation results reveal VICTOR's effectiveness in conveying verification. The field user study results demonstrate that VICTOR increases users' sensitivity

and vigilance to fake news in the long term. We believe such an implicit approach is a promising direction to defend liberty in crisis.

## ACKNOWLEDGMENTS

## REFERENCES

[1] Davey Alba and Kate Conger. 2020. Twitter moves to target fake videos and photos (Accessed Sept 5, 2020). https://www.nytimes.com/2020/02/04/technology/twitter-fake-videos-photos-disinformation.html

[2] Cristian Bravo-Lillo, Lorrie Cranor, Saranga Komanduri, Stuart Schechter, and Manya Sleeper. 2014. Harder to ignore? Revisiting pop-up fatigue and approaches to prevent it. In *10th Symposium On Usable Privacy and Security ({SOUPS} 2014)*. 105–111.

[3] Chia-Wei Chen, Sheng-Chuan Chou, Chang-You Tai, and Lun-Wei Ku. 2019. Phrase-Guided Attention Web Article Recommendation for Next Clicks and Views. In *Proceedings of the 2019 IEEE/ACM International Conference on Advances in Social Networks Analysis and Mining*. 315–324.

[4] Katherine Clayton, Spencer Blair, Jonathan A Busam, Samuel Forstner, John Glance, Guy Green, Anna Kawata, Akhila Kovvuri, Jonathan Martin, Evan Morgan, et al. 2019. Real solutions for fake news? Measuring the effectiveness of general warnings and fact-check tags in reducing belief in false stories on social media. *Political Behavior* (2019), 1–23.

[5] Rajarshi Das, Shehzaad Dhuliawala, Manzil Zaheer, Luke Vilnis, Ishan Durugkar, Akshay Krishnamurthy, Alexander J. Smola, and Andrew McCallum. 2018. Go for a Walk and Arrive at the Answer: Reasoning Over Paths in Knowledge Bases using Reinforcement Learning. *CoRR* abs/1711.05851 (2018). arXiv:1711.05851 http://arxiv.org/abs/1711.05851

[6] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2018. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. *CoRR* abs/1810.04805 (2018). arXiv:1810.04805 http://arxiv.org/abs/1810.04805

[7] Mehrdad Farajtabar, Jiachen Yang, Xiaojing Ye, Huan Xu, Rakshit Trivedi, Elias Khalil, Shuang Li, Le Song, and Hongyuan Zha. 2017. Fake news mitigation via point process based intervention. *arXiv preprint arXiv:1703.07823* (2017).

[8] Mingkun Gao, Ziang Xiao, Karrie Karahalios, and Wai-Tat Fu. 2018. To label or not to label: The effect of stance and credibility labels on readers' selection and perception of news articles. *ACM CHI* 2, CSCW (2018).

[9] R Kelly Garrett and Brian E Weeks. 2013. The promise and peril of real-time corrections to political misperceptions. In *Proceedings of the 2013 Conference on Computer Supported Cooperative Work*. 1047–1058.

[10] Mahak Goindani and Jennifer Neville. 2020. Social reinforcement learning to combat fake news spread. In *Uncertainty in Artificial Intelligence*. PMLR, 1006–1016.

[11] Qingyuan Gong, Yang Chen, Xinlei He, Zhou Zhuang, Tianyi Wang, Hong Huang, Xin Wang, and Xiaoming Fu. 2018. DeepScan: Exploiting deep learning for malicious account detection in location-based social networks. *IEEE Communications Magazine* 56, 11 (2018), 21–27.

[12] David J Hauser and Norbert Schwarz. 2016. Attentive Turkers: MTurk participants perform better on online attention checks than do subject pool participants. *Behavior research methods* 48, 1 (2016), 400–407.

[13] Alan G Hawkes. 1971. Spectra of some self-exciting and mutually exciting point processes. *Biometrika* 58, 1 (1971), 83–90.

[14] Sepp Hochreiter and Jürgen Schmidhuber. 1997. Long Short-Term Memory. *Neural Comput.* 9, 8 (Nov. 1997), 1735–1780. https://doi.org/10.1162/neco.1997.9.8.1735

[15] Jooyeon Kim, Behzad Tabibian, Alice Oh, Bernhard Schölkopf, and Manuel Gomez-Rodriguez. 2018. Leveraging the Crowd to Detect and Reduce the Spread of Fake News and Misinformation. In *Proceedings of the Eleventh ACM International Conference on Web Search and Data Mining* (Marina Del Rey, CA, USA) *(WSDM '18)*. Association for Computing Machinery, New York, NY, USA, 324–332. https://doi.org/10.1145/3159652.3159734

[16] Justin Kosslyn and Cong Yu. 2017. Fact check now available in Google search and news around the world (Accessed Sept 5, 2020). https://www.blog.google/products/search/fact-check-now-available-google-search-and-news-around-world/

[17] Yi-Ju Lu and Cheng-Te Li. 2020. GCAN: Graph-aware Co-Attention Networks for Explainable Fake News Detection on Social Media. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*. Association for Computational Linguistics, Online, 505–514. https://doi.org/10.18653/v1/2020.acl-main.48

[18] Neil A Macmillan and C Douglas Creelman. 2004. *Detection theory: A user's guide*. Psychology press.

[19] Sina Mohseni, Eric Ragan, and Xia Hu. 2019. Open Issues in Combating Fake News: Interpretability as an Opportunity. arXiv:1904.03016 [cs.SI]

[20] Nam P. Nguyen, Guanhua Yan, My T. Thai, and Stephan Eidenbenz. 2012. Containment of Misinformation Spread in Online Social Networks. In *Proceedings of the 4th Annual ACM Web Science Conference* (Evanston, Illinois) *(WebSci '12)*. Association for Computing Machinery, New York, NY, USA, 213–222. https://doi.org/10.1145/2380718.2380746

[21] Gordon Pennycook and Tyrone Cannon. 2018. Prior exposure increases perceived accuracy of fake news. *Journal of Experimental Psychology General* (06 2018). https://doi.org/10.2139/ssrn.2958246

[22] Gordon Pennycook, Jonathon McPhetres, Yunhao Zhang, Jackson G Lu, and David G Rand. 2020. Fighting COVID-19 misinformation on social media: Experimental evidence for a scalable accuracy-nudge intervention. *Psychological Science* 31, 7 (2020), 770–780.

[23] Gordon Pennycook and David G Rand. 2019. Fighting misinformation on social media using crowdsourced judgments of news source quality. *Proceedings of the National Academy of Sciences* 116, 7 (2019), 2521–2526.

[24] Guy Rosen, Katie Harbath, Nathaniel Gleicher, and Rob Leathern. 2019. Helping to protect the 2020 US elections (Accessed Sept 5, 2020). https://about.fb.com/news/2019/10/update-on-election-integrity-efforts/

[25] Haeseung Seo, Aiping Xiong, and Dongwon Lee. 2019. Trust It or Not: Effects of Machine-Learning Warnings in Helping Individuals Mitigate Misinformation. In *Proceedings of the 10th ACM Conference on Web Science* (Boston, Massachusetts, USA) *(WebSci '19)*. ACM, New York, NY, USA, 265–274. https://doi.org/10.1145/3292522.3326012

[26] Karishma Sharma, Feng Qian, He Jiang, Natali Ruchansky, Ming Zhang, and Yan Liu. 2019. Combating fake news: A survey on identification and mitigation techniques. *ACM Transactions on Intelligent Systems and Technology (TIST)* 10, 3 (2019), 1–42.

[27] Kai Shu, H Russell Bernard, and Huan Liu. 2019. Studying fake news via network analysis: detection and mitigation. In *Emerging Research Challenges and Opportunities in Computational Social Network Analysis and Mining*. Springer, 43–65.

[28] Kai Shu, Deepak Mahudeswaran, Suhang Wang, and Huan Liu. 2019. Hierarchical Propagation Networks for Fake News Detection: Investigation and Exploitation. *CoRR* abs/1903.09196 (2019). arXiv:1903.09196 http://arxiv.org/abs/1903.09196

[29] Kai Shu, Amy Sliva, Suhang Wang, Jiliang Tang, and Huan Liu. 2017. Fake News Detection on Social Media: A Data Mining Perspective. *SIGKDD Explor. Newsl.* 19, 1 (Sept. 2017), 22–36. https://doi.org/10.1145/3137597.3137600

[30] Jeff Smith, Grace Jackson, and Seetha Raj. 2017. Designing against misinformation (Accessed Sept 5, 2020). https://medium.com/facebook-design/designing-against-misinformation-e5846b3aa1e2

[31] Richard S. Sutton and Andrew G. Barto. 2018. *Reinforcement Learning: An Introduction* (second ed.). The MIT Press. http://incompleteideas.net/book/the-book-2nd.html

[32] Yaqing Wang, Weifeng Yang, Fenglong Ma, Jin Xu, Bin Zhong, Qiang Deng, and Jing Gao. 2020. Weak Supervision for Fake News Detection via Reinforcement Learning. In *The Thirty-Fourth AAAI Conference on Artificial Intelligence, AAAI 2020, The Thirty-Second Innovative Applications of Artificial Intelligence Conference, IAAI 2020, The Tenth AAAI Symposium on Educational Advances in Artificial Intelligence, EAAI 2020, New York, NY, USA, February 7-12, 2020*. AAAI Press, 516–523. https://aaai.org/ojs/index.php/AAAI/article/view/5389

[33] Chuhan Wu, Fangzhao Wu, Tao Qi, and Yongfeng Huang. 2021. Empowering News Recommendation with Pre-Trained Language Models. In *Proceedings of the 44th International ACM SIGIR Conference on Research and Development in Information Retrieval* (Virtual Event, Canada) *(SIGIR '21)*. Association for Computing Machinery, New York, NY, USA, 1652–1656. https://doi.org/10.1145/3404835.3463069

[34] Liang Wu, Fred Morstatter, Xia Hu, and Huan Liu. 2016. Mining misinformation in social media. *Big Data in Complex and Social Networks* (2016), 123–152.

[35] Yikun Xian, Zuohui Fu, S. Muthukrishnan, Gerard de Melo, and Yongfeng Zhang. 2019. Reinforcement Knowledge Graph Reasoning for Explainable Recommendation. In *Proceedings of the 42nd International ACM SIGIR Conference on Research and Development in Information Retrieval* (Paris, France) *(SIGIR'19)*. Association for Computing Machinery, New York, NY, USA, 285–294. https://doi.org/10.1145/3331184.3331203