

# Is This URL Safe: Detection of Malicious URLs Using Global Vector for Word Representation

Rohit Bharadwaj<sup>1</sup> Ashutosh Bhatia<sup>1</sup> Laxmi Chhibbar<sup>1</sup>  
Kamlesh Tiwari<sup>1</sup> Ankit Agrawal<sup>1</sup>

<sup>1</sup>Department of Computer Science and Information System,  
Birla Institute of Technology and Science Pilani, Pilani, Jhunjhunu, Rajasthan, 333031 India

International Conference on Information Networking (ICOIN 2022)



**ICOIN 2022**

The 36th International Conference on Information Networking (ICOIN 2022)

January 12-15, 2022, Jeju Island, Korea & Virtual Conference

# Background

Users are frequently exposed to many unknown links through advertisements and emails.

---

<sup>1</sup> *Help Net Security*. <https://www.helpnetsecurity.com/2019/10/09/phishing-increase-2019/>. Accessed: 2021-09-28.

<sup>2</sup> *Kaspersky Security Bulletin 2020-2021*. <https://securelist.com/category/kaspersky-security-bulletin/>. Accessed: 2021-09-28.

Users are frequently exposed to many unknown links through advertisements and emails.

- ▶ Malicious URLs pose a prime threat to the information security of organizations.

---

<sup>1</sup> *Help Net Security*. <https://www.helpnetsecurity.com/2019/10/09/phishing-increase-2019/>. Accessed: 2021-09-28.

<sup>2</sup> *Kaspersky Security Bulletin 2020-2021*. <https://securelist.com/category/kaspersky-security-bulletin/>. Accessed: 2021-09-28.

Users are frequently exposed to many unknown links through advertisements and emails.

- ▶ Malicious URLs pose a prime threat to the information security of organizations.
- ▶ According to a report, nearly 1 in 50 URLs are infected with malicious code.<sup>1</sup>

---

<sup>1</sup> *Help Net Security*. <https://www.helpnetsecurity.com/2019/10/09/phishing-increase-2019/>. Accessed: 2021-09-28.

<sup>2</sup> *Kaspersky Security Bulletin 2020-2021*. <https://securelist.com/category/kaspersky-security-bulletin/>. Accessed: 2021-09-28.

Users are frequently exposed to many unknown links through advertisements and emails.

- ▶ Malicious URLs pose a prime threat to the information security of organizations.
- ▶ According to a report, nearly 1 in 50 URLs are infected with malicious code.<sup>1</sup>
- ▶ Trusted domain URLs are being used to host malicious code.

---

<sup>1</sup> *Help Net Security*. <https://www.helpnetsecurity.com/2019/10/09/phishing-increase-2019/>. Accessed: 2021-09-28.

<sup>2</sup> *Kaspersky Security Bulletin 2020-2021*. <https://securelist.com/category/kaspersky-security-bulletin/>. Accessed: 2021-09-28.

Users are frequently exposed to many unknown links through advertisements and emails.

- ▶ Malicious URLs pose a prime threat to the information security of organizations.
- ▶ According to a report, nearly 1 in 50 URLs are infected with malicious code.<sup>1</sup>
- ▶ Trusted domain URLs are being used to host malicious code.
- ▶ In 2019, 85% of detected web threats were malicious URLs.<sup>2</sup>

---

<sup>1</sup> *Help Net Security*. <https://www.helpnetsecurity.com/2019/10/09/phishing-increase-2019/>. Accessed: 2021-09-28.

<sup>2</sup> *Kaspersky Security Bulletin 2020-2021*. <https://securelist.com/category/kaspersky-security-bulletin/>. Accessed: 2021-09-28.

- ▶ **IDS (Intrusion Detection Systems):**

## Current Approach

- ▶ **IDS (Intrusion Detection Systems)**: uses URL blacklisting or signature blacklisting.



# Current Approach

- ▶ **IDS (Intrusion Detection Systems)**: uses URL blacklisting or signature blacklisting.
- ▶ These blacklist URLs maybe secured behind a paywall.

# Current Approach

- ▶ **IDS (Intrusion Detection Systems)**: uses URL blacklisting or signature blacklisting.
- ▶ These blacklist URLs maybe secured behind a paywall.
- ▶ Keeping updated list of blacklist URLs is challenging.

# Current Approach

- ▶ **IDS (Intrusion Detection Systems)**: uses URL blacklisting or signature blacklisting.
- ▶ These blacklist URLs maybe secured behind a paywall.
- ▶ Keeping updated list of blacklist URLs is challenging.
- ▶ Attackers have started using **DGAs** which renders blacklist method ineffective.

# Machine Learning Based Approaches

- ▶ Train a learning model based on URLs' **static** and **dynamic** features.

# Machine Learning Based Approaches

- ▶ Train a learning model based on URLs' **static** and **dynamic** features.

## Dynamic Features

- ▶ Activities performed in the target system due to the URL.
- ▶ Features associated with the content of the webpage referred to by the URL.

# Machine Learning Based Approaches

- ▶ Train a learning model based on URLs' **static** and **dynamic** features.

## Dynamic Features

- ▶ Activities performed in the target system due to the URL.
- ▶ Features associated with the content of the webpage referred to by the URL.

## Static Features

- ▶ Directly derived from the URL.
- ▶ Ex: URL length, presence of unsafe extension.

## Related Work

- ▶ Static classifiers have been modeled using lexical features of URL. (i.e. host properties, network traffic, etc.)
- ▶ The **bag-of-words model** is often used with these approaches, but not helpful in our case due to significant latency.
- ▶ The bag-of-words model leads to the loss of contextual information from the data.

## Related Work

- ▶ (James et al., 2013)<sup>3</sup> detected phishing websites using lexical features. However, their approach depends on specially crafted features that are not suitable for large datasets.
- ▶ (Verma & Dyer, 2015)<sup>4</sup> built machine learning classifiers using URL's lexical features. However, it is not robust to spelling errors in the malicious URLs.
- ▶ Our work builds upon the work by (Hai & Hwang, 2018)<sup>5</sup>, but we use different feature extraction methods, different dataset sizes and introduce **GloVe** based embedding learning instead of Word2Vec embedding used by previous authors.

---

<sup>3</sup>Joby James, L Sandhya, and Ciza Thomas. "Detection of phishing URLs using machine learning techniques". In: *2013 international conference on control communication and computing (ICCC)*. IEEE. 2013, pp. 304–309.

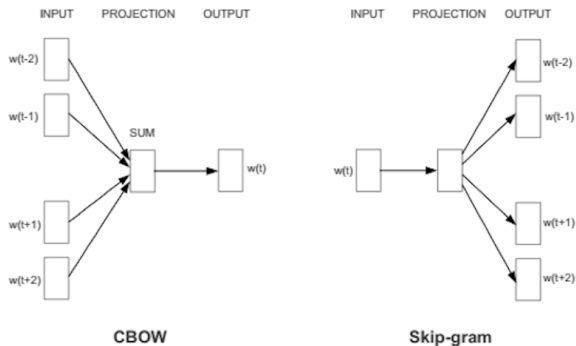
<sup>4</sup>Rakesh Verma and Keith Dyer. "On the character of phishing URLs: Accurate and robust statistical learning classifiers". In: *Proceedings of the 5th ACM Conference on Data and Application Security and Privacy*. 2015, pp. 111–122.

<sup>5</sup>Quan Tran Hai and Seong Oun Hwang. "Detection of malicious URLs based on word vector representation and ngram". In: *Journal of Intelligent & Fuzzy Systems* 35.6 (2018), pp. 5889–5900.



# Word2Vec

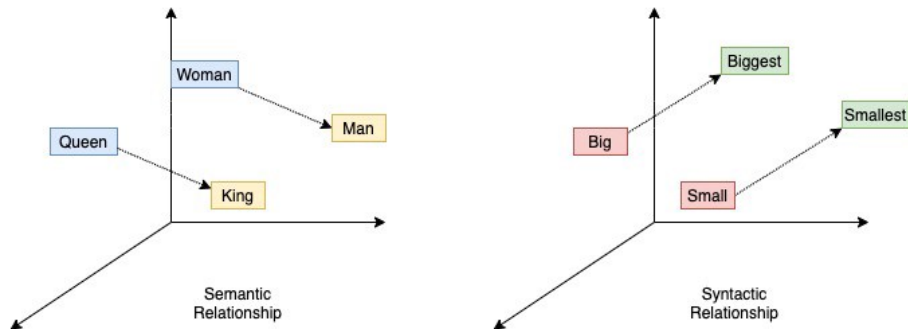
Word2Vec focuses on learning word embedding. For this, a two-layer neural network is trained to remodel the word's semantic contexts.



**Figure:** The two variants of Word2Vec model. CBOW predicts center word given the context words, while skip-gram predicts the context words given the center word.

# Word2Vec Analogy Tasks

These word embeddings are capable of identifying similarities between words. They use vector arithmetic to reconstruct the syntactic and linguistic patterns.



**Figure:** The figure shows the different relationship captured by Word2Vec models.  $\text{Embedding}(\text{"King"}) - \text{Embedding}(\text{"Man"}) + \text{Embedding}(\text{"Woman"})$  yields a vector which is closest in similarity to the vector  $\text{Embedding}(\text{"Queen"})$ .

Word2Vec model considers only local perspectives to learn the embedding. Whereas, GloVe model builds upon the work of Word2Vec by also considering the global word-word co-occurrence statistics.

GloVe based models have also been shown to give better results in many natural language processing tasks compared to Word2Vec.<sup>6,7</sup>

With this work we show that the ability to incorporate global word co-occurrence statistics through GloVe model helps to better discriminate between malicious and benign URLs.

---

<sup>6</sup>Jeffrey Pennington, Richard Socher, and Christopher D Manning. "Glove: Global vectors for word representation". In: *Proceedings of the 2014 conference on empirical methods in natural language processing (EMNLP)*. 2014, pp. 1532–1543.

<sup>7</sup>Christopher Ifeanyi Eke et al. "The significance of global vectors representation in sarcasm analysis". In: *2020 International Conference in Mathematics, Computer Engineering and Computer Science (ICMCECS)*. IEEE. 2020, pp. 1–7.

# Proposed Approach

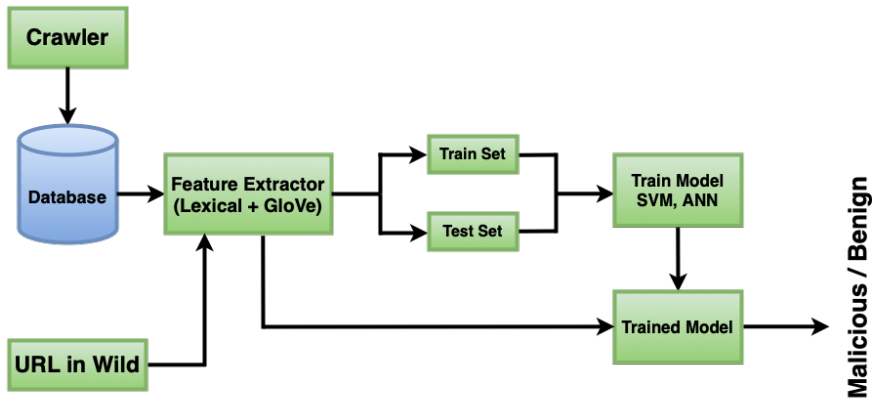
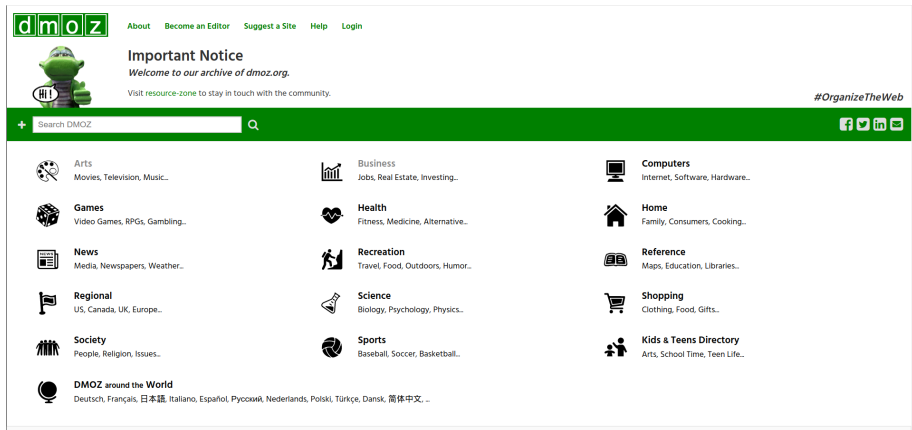


Figure: Proposed Malicious URL Detector.

# Getting the Data

The data for benign URLs were collected from DMOZ. The website contains several categories and subcategories where there are many links.



The screenshot shows the DMOZ website interface. At the top left is the 'dmoz' logo. Navigation links include 'About', 'Become an Editor', 'Suggest a Site', 'Help', and 'Login'. A green banner features a cartoon character saying 'Hi!' and an 'Important Notice' about the archive of dmoz.org. A search bar is located below the banner. The main content area displays a grid of category icons and titles: Arts, Games, News, Regional, Society, DMOZ around the World, Business, Health, Recreation, Science, Sports, Computers, Home, Reference, Shopping, and Kids & Teens Directory. A social media bar with icons for Facebook, Twitter, LinkedIn, and YouTube is at the bottom right.

**dmoz** About Become an Editor Suggest a Site Help Login

**Important Notice**  
Welcome to our archive of dmoz.org.  
Visit resource-zone to stay in touch with the community.

#OrganizeTheWeb

+ Search DMOZ

- Arts**  
Movies, Television, Music...
- Games**  
Video Games, RPGs, Gambling...
- News**  
Media, Newspapers, Weather...
- Regional**  
US, Canada, UK, Europe...
- Society**  
People, Religion, Issues...
- DMOZ around the World**  
Deutsch, Français, 日本語, Italiano, Español, Русский, Nederlands, Polski, Türkçe, Dansk, 简体中文, ...
- Business**  
Jobs, Real Estate, Investing...
- Health**  
Fitness, Medicine, Alternative...
- Recreation**  
Travel, Food, Outdoors, Humor...
- Science**  
Biology, Psychology, Physics...
- Sports**  
Baseball, Soccer, Basketball...
- Computers**  
Internet, Software, Hardware...
- Home**  
Family, Consumers, Cooking...
- Reference**  
Maps, Education, Libraries...
- Shopping**  
Clothing, Food, Gifts...
- Kids & Teens Directory**  
Arts, School Time, Teen Life...

Figure: DMOZ website with several top-level categories.

# Getting the Data

## ▼ Subcategories 29

- |   |  |   |   |
|---|--|---|---|
| <ul style="list-style-type: none"><li>☐ Anime</li><li>☐ Production<ul style="list-style-type: none"><li>➤ Animated Graphics</li><li>➤ Animation Art Galleries</li></ul></li><li>☐ Artists</li><li>☐ Audio</li><li>☐ Awards</li><li>☐ Chats and Forums</li></ul> | <ul style="list-style-type: none"><li>☐ Cartoons</li><li>☐ Collectibles</li><li>☐ Computer</li><li>☐ Contests</li><li>➤ Directors</li><li>☐ DVD</li><li>☐ Experimental</li></ul> | <ul style="list-style-type: none"><li>☐ Movies</li><li>☐ Festivals</li><li>☐ Magazines and E-zines</li><li>☐ News and Media</li><li>☐ Organizations</li><li>➤ Puppetry</li><li>➤ Shopping</li></ul> | <ul style="list-style-type: none"><li>☐ Voice Actors</li><li>☐ Stop-Motion</li><li>☐ Training</li><li>➤ Video Games</li><li>☐ Web</li><li>☐ Web Rings</li><li>☐ Writers</li></ul> |
|---|--|---|---|

## ▼ Related categories 2

- Kids and Teens > Entertainment > Animation
- Regional > Europe > United Kingdom > Arts and Entertainment > Animation

## ▼ Sites 5



-  **Animation World Network** ★  
Provides information resources to the international animation community. Features include searchable database archives, monthly magazine, web animation guide, the Animation Village, discussion forums and other useful resources.
-  **About.com: Animation Guide**  
Keep up with developments in online animation for all skill levels. Download tools, and seek inspiration from online work.

Figure: Subcategories and website links found in DMOZ.

# Getting the Data

```
try:
    url = "https://dmoz-odp.org"
    myreq=req.Request(url,headers=headers)
    resp=req.urlopen(myreq)
    respData=str(resp.read())
    patt=r'<h2 class="top-cat"><a href="[^"]+'          ###Regular Expression for match
    links=re.findall(patt,respData) #Parsing Data Using Regex
    links = links[1:-1]
    print('The total of %i links are as follows : '%len(links))

    for i in links:
        i = i.replace(r'<h2 class="top-cat"><a href="','')
        i = url + i
        getSubCat(i,True)
    print('\n')
```

**Figure:** Python code to extract categories from DMOZ in a DFS based manner.

In this work, **80,128** benign URLs were collected using the above-stated algorithm. Further, **147,781** malicious URLs were obtained and used in this work.<sup>8,9</sup>

<sup>8</sup>AK Singh. "Malicious and Benign Webpages Dataset". In: *Data in brief* 32 (2020), p. 106304.

<sup>9</sup>URL: <https://www.kaggle.com/antonyj453/urldataset>.

# Getting the Data

```
headers={}  
headers['user-agent']="Mozilla/5.0 (X11; Linux x86_64) AppleWebKit/537.36 (KHTML, like Gecko) Chrome/50.0.2661.102 Safari/537.36"  
  
visited = {}  
visited["a"] = -1  
  
def getSubCat(link, isNotFirstWorld):  
    #saveURL(link)  
  
    myreq=req.Request(link,headers=headers)  
    resp=req.urlopen(myreq)  
    respData=str(resp.read())  
    # <div class="cat-item">  
    #<a href=  
    #print(respData[0:10])  
    patt=r'<div class="cat-item">[\r\n]+([>]+)'      ###Regular Expression for matching and searching Relevant Link on the Page  
    links=re.findall(patt,respData) #Parsing Data Using Regex  
    for i in range(len(links)):  
        links[i] = str(links[i]).strip()[1:-1]  
  
    #Need to remove /World entries  
    links2 = []  
    for e in links:  
        if ( (not e.startswith('<a href="/World') ) and (isNotFirstWorld) ) or (not isNotFirstWorld) ):  
            links2.append(e)  
  
    #links2 = links  
  
    #links = links2  
    link = "https://dmoz-odp.org"  
    # print(links2)  
    for i in links2:  
        i = i.replace(r'<a href=', '')  
        i = link + i  
        print("Saving All Links in: ", i)  
        try:  
            if visited.get(i) == None:  
                saveURL(i)  
                visited[i] = 1  
                getSubCat(i, True)
```

Figure: Code to extract URL from categories.

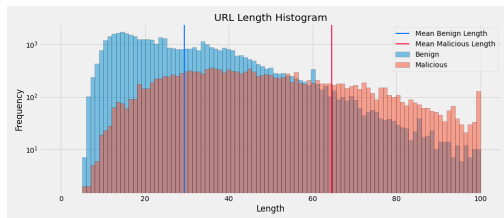


# Extracting Features

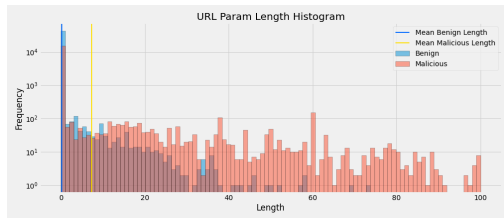
Feature Set	
Statistical Features	Type
URL Length	Continuous
Path Length	Continuous
Parameter Length	Continuous
is_exe	Categorical
is_dll	Categorical
letter-count	Continuous
digit-count	Continuous
count@	Continuous
count#	Continuous
count%	Continuous
use_of_ip	Categorical
vowByCons	Continuous
DigByLetter	Continuous
blackNGScore	Continuous
New Feature	Type
GloVe	Continuous

Table: URL Lexical Features considered for training the ML Model.

# Histogram



(a) Histogram of URL length.



(b) Histogram of URL parameter length.

Figure: Histogram of URL length and URL parameter length.

# Proposed GloVe Feature

Following from (Hai & Hwang, 2018)<sup>10</sup>, we train the GloVe model on only benign URL tokens. By doing so, the malicious URL tokens would be far from benign tokens in the vector space.

The GloVe score is calculated using Algorithm 1.

---

**Algorithm 1:** Computing GloVe scores.

---

**Input:** tokenList // Obtained from URL

**Output:** s // GloVe score for the URL

```
1 foreach token  $w \in$  tokenList do
2   if  $w \in$  embeddingsIndex then
3      $myVec \leftarrow myVec + embeddingsIndex[w]$ 
4   end
5   else
6      $myVec \leftarrow myVec + \mathbf{0}$ 
7   end
8 end
9  $myVec \leftarrow \frac{myVec}{len(tokenList)}$ 
10 return  $s \leftarrow ||myVec||$ 
```

---

<sup>10</sup>Quan Tran Hai and Seong Oun Hwang. "Detection of malicious URLs based on word vector representation and ngram". In: *Journal of Intelligent & Fuzzy Systems* 35.6 (2018), pp. 5889–5900.

# Classifiers

SVM with linear kernel and a simple ANN with architecture shown in the below figure were used as classifiers.

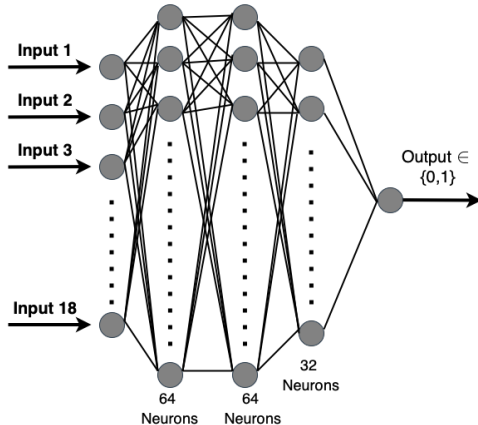


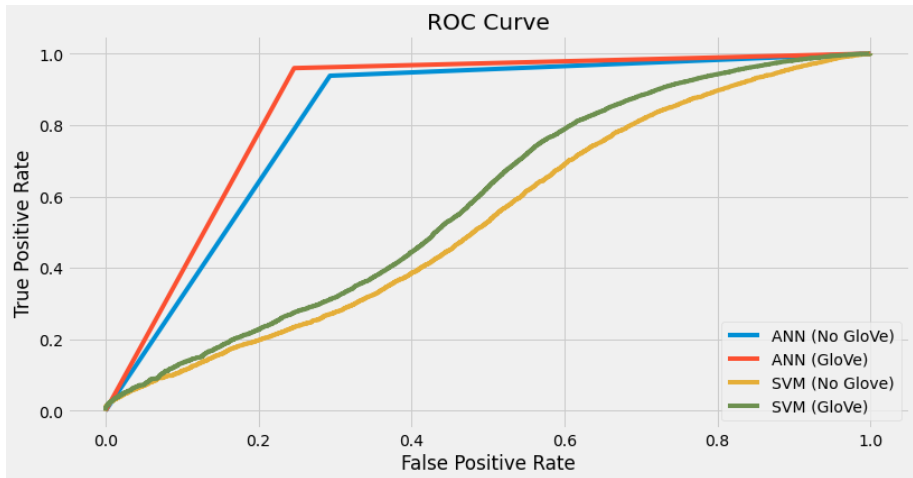
Figure: Neural Network Architecture

# Results

Classifier	Feature Type	Accuracy (%)	Precision (%)	F1
Support Vector Machine (SVM)	Only Statistical Features	69.72	0.69	0.80
Artificial Neural Network (ANN)	Only Statistical Features	86.00	0.86	0.90
Support Vector Machine (SVM)	GloVe with Statistical Feature	<b>77.26</b>	<b>0.80</b>	<b>0.83</b>
Artificial Neural Network (ANN)	GloVe with Statistical Feature	<b>89.00</b>	<b>0.88</b>	<b>0.92</b>

**Table:** Performance of the proposed system as compared with other setting.

# ROC Curve



**Figure:** Receiver Operating Curve (ROC) while using the proposed feature on SVM and ANN classifier. Higher is the area under the curve, the better is the model.

# Conclusion

We proposed using GloVe as one of the features to improve the accuracy of machine learning models for detecting malicious URLs. We observe that the proposed system has improved the overall performance of malicious URL detection, as it has reduced the error by **63.33%** when compared to the traditional approach.

The performance gain is observed in both ANN and SVM, by including the GloVe based feature. This validates our hypothesis and demonstrates the effectiveness of using GloVe based features to classify malicious URLs.