

## VIFDD Data Card



VIFDD	VIFDD (Visual Intrusion and Fraud Detection Dataset) is a dataset designed for the task of detecting visual intrusions and fraudulent websites. The dataset consists of 240,000 images of webpage screenshots, evenly divided into two categories: SCAM and NORMAL.
DATASET LINK <a href="#">Dataset</a>	DATA CARD AUTHOR(S)  Will be added after acceptance

Dataset Owners		
TEAM(S)	CONTACT DETAIL(S)	AUTHOR(S)
VIFDD Team	<b>Dataset Owner(s):</b> will be revealed on acceptance  <b>Affiliation:</b> will be revealed on acceptance  <b>Contact:</b> will be revealed on acceptance	Will be revealed on acceptance

Dataset Overview																
DATA SUBJECT(S)	DATASET SNAPSHOT	CONTENT DESCRIPTION														
Data about systems or products and their behaviors	<table><tr><td>Size of Dataset</td><td>50 GB</td></tr><tr><td>Number of Instances</td><td>240,000</td></tr><tr><td>Labeled Classes</td><td>2</td></tr><tr><td>Number of Labels</td><td>2</td></tr><tr><td>Average Labels Per Instance</td><td>1</td></tr><tr><td>Algorithmic Labels</td><td>None</td></tr><tr><td>Human Labels</td><td>Some<sup>1</sup></td></tr></table>	Size of Dataset	50 GB	Number of Instances	240,000	Labeled Classes	2	Number of Labels	2	Average Labels Per Instance	1	Algorithmic Labels	None	Human Labels	Some <sup>1</sup>	<p>Each datapoint in the VIFDD dataset represents a screenshot of a webpage. The content of a datapoint includes:</p> <ul style="list-style-type: none"><li>• <b>Image Data:</b> A 224x224 pixel image in RGB format, stored as either a PNG or JPG file.</li><li>• <b>Label:</b> A categorical label indicating whether the webpage is a SCAM or NORMAL. This label is determined based on the source of the URL:<ul style="list-style-type: none"><li>• <b>SCAM:</b> The image is labelled as SCAM if the screenshot is from a known fraudulent website.</li><li>• <b>NORMAL:</b> The image is labelled as NORMAL if the screenshot is from a legitimate website.</li></ul></li></ul>
	Size of Dataset	50 GB														
	Number of Instances	240,000														
	Labeled Classes	2														
	Number of Labels	2														
	Average Labels Per Instance	1														
	Algorithmic Labels	None														
	Human Labels	Some <sup>1</sup>														
<p><b>Above:</b> Summary of VIFDD dataset.</p> <p><sup>1</sup> Some websites needed to be cross-verified manually before being labelled as SCAM.</p>																
Sensitivity of Data																
SENSITIVITY TYPE(S)	FIELD(S) WITH SENSITIVE DATA															
None	<p><b>Intentionally Collected Sensitive Data</b></p> <p>The dataset does not intentionally collect sensitive data.</p> <p><b>Unintentionally Collected Sensitive Data</b></p> <p>Since the screenshots are taken from live webpages, there is a potential for unintentional capture of sensitive data. This could include:</p> <ul style="list-style-type: none"><li>• <b>Personal Information:</b> Names, email addresses, phone numbers, or other personal information visible on the webpage.</li><li>• <b>Pornographic Content:</b> Explicit content that might be displayed on some scam or adult websites.</li><li>• <b>Violent Content:</b> Images or text depicting violence that might be present on certain webpages.</li></ul>															
Dataset Version and Maintenance																
MAINTENANCE STATUS	VERSION DETAILS	MAINTENANCE PLAN														

<b>Limited Maintenance</b> The data will not be updated, but any technical issues will be addressed.	<b>Current Version:</b> 1.0 <b>Last Updated:</b> 05/2024 <b>Release Date:</b> N/A	Since the dataset is static, there will be no additions or modifications to the dataset itself. Any discovered issues or updates will be documented and communicated through the dataset's documentation.
---	---	---

Example of Data Points

PRIMARY DATA MODALITY	SAMPLING OF DATA POINTS	DATA FIELDS									
Image Data	Examples of data in VIFDD										
		<table><tr><th>Field Name</th><th>Type</th><th>Description</th></tr><tr><td>image_data</td><td>PNG/JPG</td><td>Pixel data of the screenshot.</td></tr><tr><td>label</td><td>Integer/String</td><td>Class of the screenshot.</td></tr></table>	Field Name	Type	Description	image_data	PNG/JPG	Pixel data of the screenshot.	label	Integer/String	Class of the screenshot.
	Field Name	Type	Description								
	image_data	PNG/JPG	Pixel data of the screenshot.								
label	Integer/String	Class of the screenshot.									
label: SCAM		<b>Above:</b> Summary of data fields in VIFDD.									
label: NORMAL											
TYPICAL DATA POINT	ATYPICAL DATA POINT										
A typical data point.	The dataset does not contain atypical data points as far as we know.										
<table><tr><th>Field</th><th>Value</th></tr><tr><td>image_data</td><td>“\xFF\xD8\xFF\xE0\x00\x10JFIF\x00\x01\x01\x00\x00\x01\x00\x01\x...”</td></tr><tr><td>Label</td><td>[“SCAM”, “NORMAL”]</td></tr></table>	Field	Value	image_data	“\xFF\xD8\xFF\xE0\x00\x10JFIF\x00\x01\x01\x00\x00\x01\x00\x01\x...”	Label	[“SCAM”, “NORMAL”]					
Field	Value										
image_data	“\xFF\xD8\xFF\xE0\x00\x10JFIF\x00\x01\x01\x00\x00\x01\x00\x01\x...”										
Label	[“SCAM”, “NORMAL”]										

Motivations & Intentions		
Motivations		
PURPOSE(S)	DOMAIN(S) OF APPLICATION	MOTIVATING FACTOR(S)
Research	Cybersecurity, Fraud Detection, Machine Learning, Computer Vision, Web Security, Artificial Intelligence, Data Science, E-commerce Security, Internet Safety, Phishing Detection, Scam Prevention	<ol style="list-style-type: none"> <li><b>Enhancing Fraud Detection.</b></li> <li><b>Improving Cybersecurity.</b></li> <li><b>Supporting E-commerce Security.</b></li> <li><b>Filling Data Gaps.</b></li> <li><b>Encouraging Academic Research.</b></li> <li><b>Promoting Internet Safety.</b></li> </ol> <p>The VIFDD-2024 dataset aims to enhance fraud detection, improve cybersecurity, advance machine learning in visual content analysis, and support e-commerce security. It fills data gaps in scam detection research, encourages academic studies, and promotes internet safety by enabling the development of tools to identify and block scam websites.</p>
Intended Use		
DATASET USE(S)	SUITABLE USE CASE(S)	UNSUITABLE USE CASE(S)
Safe for research use	<ul style="list-style-type: none"> <li><b>Fraud Detection Research:</b> Developing and evaluating machine learning models for detecting scam websites based on visual content.</li> <li><b>Cybersecurity Studies:</b> Conducting research in web security to identify patterns and features indicative of fraudulent sites.</li> <li><b>Phishing and Scam Prevention:</b> Training and testing tools that detect and prevent phishing attacks and scam websites.</li> <li><b>Computer Vision Applications:</b> Exploring computer vision techniques to analyse and classify webpage screenshots.</li> </ul>	<ul style="list-style-type: none"> <li><b>Production Environments:</b> Using the dataset in live production environments for real-time fraud detection or security applications without additional validation and testing.</li> </ul>
	RESEARCH AND PROBLEM SPACE(S)	CITATION GUIDELINES

The VIFDD dataset addresses the problem space of visual intrusion and fraud detection on the internet. The dataset supports the advancement of computer vision techniques for webpage classification and the development of machine learning models that differentiate between legitimate and fraudulent sites.

**Guidelines & Steps:**

1. Include the full citation in your references.
2. Provide a direct link to the dataset wherever possible.

**BiBTeX:**

```
`` ` @dataset{vifdd,
  title={VIFDD: Visual Intrusion
and Fraud Detection Dataset},
  author={on acceptance},
  year={2024},
  url={
https://www.kaggle.com/datasets/ae6
753dd33076e09f4803961a000e8e5adbd6a
1d6c16829195f422513720af3c }
} `` `
```

Provenance		
Collection		
METHOD(S) USED	METHODOLOGY DETAIL(S)	SOURCE DESCRIPTION(S)
Scraped or Crawled	<p><b>Collection Type:</b> Scraped or Crawled</p> <p><b>Source:</b></p> <ul style="list-style-type: none"> <li><b>Normal URLs:</b> Gathered from the Alexa Top 1 Million Websites list.</li> <li><b>Scam URLs:</b> Collected from various online sources including <a href="#">PhishTank</a>, <a href="#">Scamwatch</a>, <a href="#">Scamalytics</a>, and <a href="#">Spamhaus</a>.</li> </ul> <p><b>Is this source considered sensitive or high-risk?</b> No</p> <p><b>Dates of Collection:</b> Oct 2023 - Mar 2024</p> <p><b>Primary Modality of Collected Data:</b></p> <ul style="list-style-type: none"> <li><b>Image Data:</b> Screenshots of webpages resized to 224x224 pixels, stored in PNG and JPG formats.</li> </ul> <p><b>Update Frequency for Collected Data:</b></p> <ul style="list-style-type: none"> <li>Static</li> </ul>	<p><b>Normal URLs:</b></p> <ul style="list-style-type: none"> <li>Normal URLs were sourced from the Alexa Top 1 Million Websites list, representing a diverse range of legitimate websites across various domains and industries.</li> </ul> <p><b>Scam URLs:</b></p> <ul style="list-style-type: none"> <li>Scam URLs were collected from several reputable online sources, including PhishTank, Scamwatch, Scamalytics, and Spamhaus. These sources specialize in identifying and cataloging fraudulent websites, serving as valuable repositories for detecting online scams and threats.</li> </ul> <p><b>Additional Notes:</b> The Alexa Top 1 Million Websites list is a widely recognized resource for gathering website popularity and traffic data. The selection from this list provides a comprehensive sample of normal webpages for the dataset. The inclusion of URLs from multiple sources enhances the diversity and representativeness of scam webpages in the dataset.</p>
COLLECTION CADENCE	DATA INTEGRATION	DATA PROCESSING

Static

Data was collected once from single or multiple sources.

Included Fields

Field Name	Description
URL	This field contains the URLs of normal webpages sourced from the Alexa Top 1 Million Websites list.
Screenshot	Screenshots of the webpages in PNG and JPG formats, resized to 224x224 pixels. Each screenshot provides a visual representation of the webpage content, capturing layout, design elements, and textual information. These screenshots are the primary data used for training and evaluating models in the dataset

Excluded Fields  
None

Scraped or Crawled

**Description:** Data for the VIFDD dataset was collected through web scraping and crawling methods. Normal URLs were gathered from the Alexa Top 1 Million Websites list, while scam URLs were collected from various online sources including PhishTank, Scamwatch, Scamalytics, and Spamhaus.

**Methods employed:** Web scraping and crawling techniques were used to extract URLs from the selected sources. This involved programmatically navigating webpages, identifying relevant URLs, and extracting them for further processing.

**Tools or libraries:** Node.js-based libraries such as Puppeteer were utilized for web scraping and crawling tasks. These libraries provide functionality for parsing HTML content, navigating webpage structures, and extracting desired information.

**Additional Notes:** The collected URLs were then used as input to the custom script responsible for taking screenshots of the corresponding webpages. The screenshots were resized to a standardized format of 224x224 pixels and stored in PNG/JPG formats for further analysis.



# Transformations

## Synopsis

TRANSFORMATION(S) APPLIED	FIELD(S) TRANSFORMED	LIBRARY(IES) AND METHOD(S) USED
Others (Resizing)	<p><b>Transformation Type:</b> Resizing</p> <ul style="list-style-type: none"><li><b>Original Data:</b> Image files in PNG and JPG formats of all sizes.</li><li><b>Transformed Data:</b> Image files resized to 224x224 pixels</li><li><b>Description:</b> The image data, consisting of screenshots of webpages, was transformed by resizing each image to a standardized dimension of 224x224 pixels. This transformation ensures uniformity in the size of images across the dataset, facilitating consistent analysis and modelling.</li></ul> <p><b>Additional Notes:</b></p> <ul style="list-style-type: none"><li>The transformation of image data to a standardized size simplifies the processing and analysis pipeline, as models can expect inputs of consistent dimensions.</li><li>Resizing the images to a common size also helps in reducing computational complexity during model training and inference.</li></ul>	<p><b>Transformation Type:</b> Resizing</p> <p><b>Method:</b> The image resizing process was implemented using the OpenCV library in Python. OpenCV (Open Source Computer Vision Library) is a popular open-source computer vision and image processing library that provides a wide range of functions for manipulating images. The <b>cv2.resize()</b> function from the OpenCV library was utilized to resize each image to the desired dimensions of 224x224 pixels.</p> <p><b>Platforms, tools, or libraries:</b></p> <ul style="list-style-type: none"><li>OpenCV</li></ul> <p><b>Transformation Results:</b> The resizing process transformed the original images into a standardized size of 224x224 pixels, ensuring consistency in the dimensions of all images in the dataset.</p>

## Breakdown of Transformations

Others (Resizing)	METHOD(S) USED	COMPARATIVE SUMMARY
-------------------	----------------	---------------------

The raw image data, consisting of screenshots of webpages is of various dimensions depending on the web page content and length. The raw image data was transformed by resizing each image to a standardized dimension of 224x224 pixels.

**Platforms, tools, or libraries:**  
OpenCV

**Before transformation:**



**After transformation:**



**Above:** Before and after comparison of transformation.