

**KLASIFIKASI MALICIOUS URL PADA FILE BERBASIS
HOST-BASED FEATURE EXTRACTION MENGGUNAKAN
METODE LSTM**

SKRIPSI

**Diajukan Untuk Melengkapi Salah Satu Syarat
Memperoleh Gelar Sarjana Komputer**



OLEH :

**MUHAMMAD REALDI
09011281823070**

**JURUSAN SISTEM KOMPUTER
FAKULTAS ILMU KOMPUTER
UNIVERSITAS SRIWIJAYA
2025**

HALAMAN PENGESAHAN

SKRIPSI

Klasifikasi Malicious URL Pada File Berbasis Host-Based Feature Extraction Menggunakan Metode LSTM

Sebagai salah satu syarat untuk penyelesaian studi di

Program Studi S1 Sistem Komputer

Oleh:

MUHAMMAD REALDI

09011281823070

Pembimbing 1

: Dr. Ir. Ahmad Heryanto, M.T.

NIP. 198701222015041002

Mengetahui

Ketua Jurusan Sistem Komputer



Dr. Ir. Sukemi, M.T
196612032006041001

HALAMAN PERSETUJUAN

Telah diuji dan lulus pada

Hari : Jum'at

Tanggal : 25 Juli 2025

Tim penguji :

1. Ketua : Dr. Ahmad Zarkasi, M.T.

2. Penguji : Aditya Putra Perdana Prasetyo, M.T.

3. Pembimbing I : Dr. Ir. Ahmad Heryanto, M.T.



HALAMAN PERNYATAAN

Yang bertanda tangan di bawah ini :

Nama : Muhammad Realdi

NIM : 09011281823070

Judul : Klasifikasi Malicious URL Pada File Berbasis Host-Based Feature Extraction Menggunakan Metode LSTM

Hasil Pengecekan Software *iThenticate/Turnitin* : 4%

Menyatakan bahwa laporan tugas akhir saya merupakan hasil karya sendiri dan bukan hasil penjiplakan atau plagiat. Apabila ditemukan unsur penjiplakan atau plagiat dalam laporan tugas akhir ini, maka saya bersedia menerima sanksi akademik dari universitas sriwijaya

Demikian, pernyatan ini saya buat dalam keadaan sadar dan tanpa paksaan dari siapapun.



Indralaya, Juli 2025

Muhammad Realdi

NIM. 09011281823070

KATA PENGANTAR

Puji dan syukur Alhamdulillah penulis panjatkan atas kehadirat Allah SWT yang telah memberikan karunia dan rahmat-Nya, sehingga penulis dapat menyelesaikan penulisan Tugas Akhir ini yang berjudul “**Klasifikasi Malicious URL Pada File Berbasis Host-Based Feature Extraction Menggunakan Metode LSTM**”.

Pada penyusunan tugas akhir ini, penulis mengucapkan terima kasih kepada berbagai pihak yang telah membantu dalam menyelesaikan penulisan Tugas Akhir ini. Oleh karena itu, pada kesempatan ini penulis mengucapkan rasa syukur dan terima kasih kepada :

1. Allah Subhanahu Wa ta’ala yang memberikan rahmat dan hidayah-Nya serta nikmat yang penulis dapatkan dengan baik dan kelancaran dalam menyelesaikan penulisan Proposal Tugas Akhir ini.
2. Kedua orang tua penulis dan saudara yang telah mendukung serta doa dan motivasi yang diberikan kepada penulis.
3. Bapak Prof. Dr. Erwin, S.Si., M.Si. selaku Dekan Fakultas Ilmu Komputer Universitas Sriwijaya.
4. Bapak Dr. Ir. H. Sukemi, M.T. selaku Ketua Jurusan Sistem Komputer Fakultas Ilmu Komputer Universitas Sriwijaya.
5. Bapak Dr. Ir. Ahmad Heryanto, M.T. selaku Dosen Pembimbing Akademik serta dosen Pembimbing Tugas Akhir penulis yang telah berkenan meluangkan waktunya untuk membimbing, memberikan saran dan motivasi serta bimbingan terbaik untuk penulisan dalam menyelesaikan Tugas Akhir ini.
6. Muhammad Imam Rafi, M. Taufik, Muhammad Andiko Putra, dan Muhammad Aldi Pangestu selaku rekan yang membantu Menyusun dan menyelesaikan penulisan Tugas Akhir ini.

7. Prazna Paramitha Avi dan Rani Octaviani yang telah membantu dalam penulisan Tugas Akhir ini.
8. Kak Angga selaku admin Jurusan Sistem Komputer yang telah membantu mengurus seluruh berkas
9. Teman-teman Sistem Komputer Angkatan 2018 Indralaya.

Penulis menyadari bahwa Tugas Akhir ini masih jauh dari kata sempurna. Oleh karena itu kritik dan saran yang dapat membangun sangat dibutuhkan penulis. Akhir kata, semoga Tugas Akhir ini dapat menjadi bermanfaat dan berguna kedepannya.

Indralaya, Juli 2025
Penulis,

Muhammad Realdi
NIM. 09011281823070

KLASIFIKASI MALICIOUS URL PADA FILE BERBASIS HOST-BASED FEATURE EXTRACTION MENGGUNAKAN METODE LSTM

Muhammad Realdi (09011281823070)

Jurusan Sistem Komputer, Fakultas Ilmu Komputer, Universitas Sriwijaya
Palembang, Indonesia

Email : real27.mr@gmail.com

ABSTRAK

Bahaya di internet yang dapat membahayakan penggunanya oleh pelaku jahat dengan menyerang secara tersembunyi. Serangan – serangan tersebut adalah serangan berupa phising, malware, spyware, dan ransomware. Salah satunya sarana serangan siber yang sangat efektif dilakukan oleh penyerang yaitu dengan menggunakan URL. URL (Uniform Resource Locator) adalah sebuah alamat yang digunakan untuk menemukan lokasi dari sebuah file yang berada di Internet. Hal ini membuat URL digunakan sebagai salah satu metode untuk melakukan serangan siber disebut sebagai Malicious URL. Malicious URL atau situs berbahaya di internet memuat banyak konten berupa spam, phising, yang digunakan untuk memulai serangan. Pada penelitian ini, menghasilkan sebuah dataset URL dengan fitur URL berupa DNS Record dari URL yang akan digunakan sebagai data dalam melakukan LSTM. Dan menghasilkan sebuah visualisasi dari data hasil LSTM dengan menggunakan epoch 100 yaitu berupa benign URL dan malicious URL. Dan melakukan analisis terhadap hasil visualisasi dari LSTM dengan menggunakan uji validasi didapatkan dengan hasil pada penelitian ini, menghasilkan validasi model dengan melakukan training dataset URL pada machine learning dan menerapkan tuning Hyperparameter sehingga hasil performa setiap rasio uji yaitu benign (0) presisi 85%, Recall 97%, F1-Score 91%, dan kluster malicious (1) presisi 96%, Recall 92%, F1-score 94%, dan hasil akurasi dari model yang digunakan yaitu dengan nilai 94.6%.

Kata Kunci : URL, Uniform Resource Locator, Malicious URL Host – Based Feature Extraction, VirusTotal.

CLASSIFICATION MALICIOUS URL ON FILE BASED ON HOST-BASED FEATURE EXTRACTION USING LSTM METHOD

Muhammad Realdi (09011281823070)

Department of Computer System, Faculty of Computer Science, University of Sriwijaya

Palembang, Indonesia

Email : real27.mr@gmail.com

ABSTRACT

Dangers on the internet can endanger users by malicious actors who attack them covertly. These attacks include phishing, malware, spyware, and ransomware. One very effective cyberattack tool carried out by attackers is using URLs. A URL (Uniform Resource Locator) is an address used to find the location of a file on the internet. This makes URLs used as one method for carrying out cyberattacks called Malicious URLs. Malicious URLs or dangerous sites on the internet contain a lot of content in the form of spam and phishing, which are used to launch attacks. In this study, a URL dataset was generated with URL features in the form of DNS records from URLs that will be used as data in conducting LSTM. And produced a visualization of the LSTM results data using epoch 100, namely benign URLs and malicious URLs. And conducting an analysis of the visualization results of LSTM using the validation test obtained with the results in this study, resulting in model validation by training the URL dataset on machine learning and applying Hyperparameter tuning so that the performance results of each test ratio are benign (0) precision 85%, Recall 97%, F1-Score 91%, and malicious cluster (1) precision 96%, Recall 92%, F1-score 94%, and the accuracy results of the model used are with a value of 94.6%.

Keyword : URL, Uniform Resource Locator, Malicious URL, Host – Based Feature Extraction, VirusTotal.

DAFTAR ISI

HALAMAN PENGESAHAN	ii
HALAMAN PERSETUJUAN.....	iii
HALAMAN PERNYATAAN.....	iv
KATA PENGANTAR.....	v
ABSTRAK	vii
ABSTRACT	viii
DAFTAR ISI.....	ix
DAFTAR GAMBAR.....	xi
DAFTAR TABEL.....	xiii
BAB I PENDAHULUAN.....	1
1.1. Latar Belakang	1
1.2. Perumusan Masalah.....	4
1.3. Batasan Masalah.....	4
1.4. Tujuan	5
1.5. Manfaat	5
1.6. Metodologi Penelitian.....	5
1.7. Sistematika Penulisan	6
BAB II Tinjauan Pustaka.....	8
2.1. Pendahuluan	8
2.2. Uniform Resource Locator (URL)	12
2.3. Malicious URL	13
2.4. Host-Based Feature Extraction	14
2.5. Artificial Intelligence.....	17
2.6. Machine Learning	17
2.7. Deep Learning.....	18
2.8. Confusion Matrix.....	18
2.9. Recurrent Neural Network.....	21
2.10. Long Short-Term Memory	24
BAB III METODOLOGI PENELITIAN	27
3.1. Pendahuluan	27

3.2.	Kerangka Kerja Penelitian.....	27
3.3.	Kerangka Kerja Metodologi Penelitian.....	29
3.4.	Kebutuhan Perangkat	30
3.5.	Skenario Eksperimen	30
3.6.	Skenario Riset	31
3.7.	Pengolahan Dataset	32
3.7.1.	Persiapan Dataset	33
3.7.2.	Ekstraksi Fitur URL.....	33
3.8.	Pre-processing	36
3.8.1.	Synthetic Minorit Oversampling Technique (SMOTE)	36
3.9.	Klasifikasi LSTM.....	39
3.10.	Tahapan LSTM Dalam Mengolah Data	40
3.11.	Validasi Hasil.....	41
	BAB IV HASIL DAN ANALISA	42
4.1.	Pendahuluan	42
4.2.	Hasil Ekstraksi File PDF.....	42
4.3.	Hasil Ekstraksi Fitur Host-Based Pada URL	43
4.4.	Hasil SMOTE.....	58
4.5.	Hyperparameter LSTM	60
4.6.	Hasil Klasifikasi	61
4.7.	Validasi Hasil Klasifikasi	62
4.7.1.	Validasi Hasil Rasio Data 50:50	62
4.7.2.	Validasi Hasil Rasio Data 60:40	66
4.7.3.	Validasi Hasil Rasio Data 70:30	69
4.7.4.	Validasi Hasil Rasio Data 80:20	72
4.7.5.	Validasi Hasil Rasio Data 90:10	75
4.8.	Hasil BACC dan MCC	78
4.9.	Validasi Hasil BACC dan MCC	79
	BAB V KESIMPULAN DAN SARAN	80
5.1.	Kesimpulan	80
5.2.	Saran.....	80
	DAFTAR PUSTAKA	81

DAFTAR GAMBAR

Gambar 2.1 Confusion Matrix	19
Gambar 2.2 Visualisasi Struktur RNN.....	22
Gambar 2.3 Model RNN[28]	22
Gambar 2.4 Forward pass dan Backward pass RNN [28]	23
Gambar 2.5 Ilustrasi blok LSTM dan memory cell units.....	25
Gambar 3.1 Kerangka Kerja Penelitian	28
Gambar 3.2 Kerangka Kerja Metode Penelitian	29
Gambar 3.3 Skenario Eksperimen Malicious URL.....	31
Gambar 3.4 Skenario Riset Malicious URL	32
Gambar 3.5 Diagram Pengolahan Dataset URL	36
Gambar 3.6 Diagram Alur Proses <i>Oversampling</i>	37
Gambar 3.7 Pseudocode Proses Oversampling.....	39
Gambar 3.8 Flowchart Klasifikasi LSTM.....	39
Gambar 3.9 Flowchart Tahapan LSTM	40
Gambar 4.1 Hasil Ekstraksi PDF	43
Gambar 4.2 Fitur Obj	44
Gambar 4.3 Fitur Age	44
Gambar 4.4 Fitur Host.....	45
Gambar 4.5 Fitur TTL.....	46
Gambar 4.6 Fitur Connection_Speed.....	46
Gambar 4.7 Fitur Is_Life.....	47
Gambar 4.8 Fitur Total_Updates.....	48
Gambar 4.9 Fitur Intended_Life_Span	48
Gambar 4.10 Fitur Life_Remaining	49
Gambar 4.11 Fitur Avg_Update_Days	50
Gambar 4.12 Fitur Reg_Country.....	50
Gambar 4.13 Fitur Days_Since_Last_Seen	51
Gambar 4.14 Fitur Days_Since_First_Seen	52
Gambar 4.15 Fitur num_open_ports, isp, num_subdomains, open_ports.....	52
Gambar 4.16 Fitur Registration_Date	53
Gambar 4.17 Fitur Expiration_Date.....	54
Gambar 4.18 Fitur days_since_first_seen	54
Gambar 4.19 Fitur days_since_last_seen.....	55
Gambar 4.20 Fitur last_updates_datesb	56
Gambar 4.21 Fitur First_Seen	56
Gambar 4.22 Fitur Last_Seen	57
Gambar 4.23 Hasil Malicious URL Dengan Pengecekan VirusTotal.....	58
Gambar 4.24 Hasil Benign URL Dengan Pengecekan VirusTotal	58
Gambar 4.25 Data Imbalance.....	59
Gambar 4.26 Data Balance	60

Gambar 4.27 Hasil Persentase Klasifikasi Berbagai Rasio Data	61
Gambar 4.28 Grafik Loss Rasio Data 50:50	63
Gambar 4.29 Grafik Akurasi Rasio Data 50:50	63
Gambar 4.30 Matriks Konfusi Rasio Data 50:50.....	64
Gambar 4.31 Average Precision Recall Rasio Data 50:50	65
Gambar 4.32 Grafik Loss Rasio Data 60:40	66
Gambar 4.33 Grafik Akurasi Rasio Data 60:40	66
Gambar 4.34 Matriks Konfusi Rasio Data 60:40.....	67
Gambar 4.35 Average Precision Recall Rasio Data 60:40	68
Gambar 4.36 Grafik Loss Rasio Data 70:30	69
Gambar 4.37 Grafik Akurasi Rasio Data 70:30	69
Gambar 4.38 Matriks Konfusi Rasio Data 70:30.....	70
Gambar 4.39 Average Precision Recall Rasio Data 70:30	71
Gambar 4.40 Grafik Loss Rasio Data 80:20	72
Gambar 4.41 Grafik Akurasi Rasio Data 80:20	72
Gambar 4.42 Matriks Konfusi Rasio Data 80:20.....	73
Gambar 4.43 Average Precision Recall Rasio Data 80:20	74
Gambar 4.44 Grafik Loss Rasio Data 90:10	75
Gambar 4.45 Grafik Akurasi Rasio Data 90:10	75
Gambar 4.46 Matriks Konvusi Rasio Data 90:10	76
Gambar 4.47 Average Precision Recall Rasio Data 90:10	77
Gambar 4.48 Persentasi Hasil BACC dan MCC Berbagai Rasio Data.....	79

DAFTAR TABEL

Tabel 2.1 Rujukan Penelitian Terdahulu.....	8
Tabel 2.2 Host-Based Features	16
Tabel 2.3 Confusion Matrix	19
Tabel 3.1 Spesifikasi Perangkat Keras.....	30
Tabel 3.2 Spesifikasi Perangkat Lunak.....	30
Tabel 3.3 Host-Based Feature Extraction	34
Tabel 3.4 Spesifikasi Parameter SMOTE	38
Tabel 4.1 Jumlah Dataset URL yang berhasil diekstraksi dari File PDF.....	43
Tabel 4.2 Detail Jumlah Data Imbalance	59
Tabel 4.3 Hyperparameter Yang Digunakan	60
Tabel 4.4 Hasil Performa Klasifikasi Rasio Data 50:50	64
Tabel 4.5 Hasil Performa Klasifikasi Rasio Data 60:40	68
Tabel 4.6 Hasil Klasifikasi Performa Rasio Data 70:30	71
Tabel 4.7 Hasil Klasifikasi Performa Rasio Data 80:20	74
Tabel 4.8 Hasil Klasifikasi Performa Rasio Data 90:10	77
Tabel 4.9 Hasil Validasi BACC dan MCC	78

BAB I

PENDAHULUAN

1.1. Latar Belakang

Semakin maju perkembangan teknologi, dan juga Perkembangan Internet yang makin pesat dengan bertambahnya pengguna internet di dunia. Meliputi dari kemajuan teknologi internet di bidang software dan hardware. Hal ini juga membawa perkembangan dalam dunia penyebaran informasi yang begitu besar mendorong setiap media untuk menggunakan internet dalam membagikan informasi kepada pengguna internet melalui layanan berupa website yang juga ikut berkembang secara cepat. Dengan munculnya layanan penyedia website hosting ini, muncul juga celah untuk kejahatan cyber yang menyerang pengguna melalui URL atau halaman web yang telah dirancang khusus yang berumur pendek, untuk menjebak pengguna internet. URL berbahaya ini, dapat disebarluaskan melalui email, pesan, facebook, twitter, whatsapp, iklan pada website, file berbentuk docx, .pdf, dan .txt dan media lainnya [1]. Pada URL website tersebut terdapat elemen berbahaya berupa malware, virus (ransomware), dan juga keylogger yang digunakan untuk membuka pintu untuk serangan ke perangkat pengguna dan kemudian mencuri data dari penggunanya.

Uniform Resource Locator (URL) merupakan istilah alamat pada halaman internet. URL terdiri dari dua komponen dasar yaitu sebagai pengidentifikasi protokol yang menunjukkan protokol yang sedang digunakan dan IP atau domain dari website tersebut berada pada dalam resource name [2].

Malicious URL merupakan situs web berbahaya yang berfungsi sebagai jebakan dari pelaku cyber. Pelaku cyber ini akan menyisipkan segala jenis elemen berbahaya seperti spam, iklan tidak pantas, malware, spoofing dan lain lain. Semua ini digunakan untuk memperoleh data korban agar menjadi sasaran kedepannya seperti korban penipuan, pengambilan informasi pribadi, dan kerugian dalam hal finansial. Mengakibatkan semakin banyaknya muncul kasus phishin, malware dan spamming maka diperlukan solusi untuk mengidentifikasi dan mengklasifikasikan URL berbahaya [3].

Penelitian yang menggunakan blacklisting [4] dengan judul *Shades of Grey: On the effectiveness of reputation-based blacklist*, mendapatkan hasil dengan tingkat false negatif yang tinggi dibandingkan dengan tingkat ekspetasi dari tingkat false positif dan dengan error minimal 5% untuk mendeteksi URL. Tetapi terdapat celah apabila URL tersebut tidak termasuk pada blacklist yang diatur, Dan tidak bisa mendeteksi Malicious URL secara banyak. Dengan sedikit perubahan pada komponen string URL dan URL yang tidak ada dalam blacklist, maka sistem pendekripsi akan tertipu sebagai URL normal[5].

Menggunakan algoritma *machine learning* dan *Deep Learning* kekurangan dari penelitian sebelumnya dapat diatasi. Beberapa metode *Machine Learning* dan *Deep Learning* yang dapat digunakan untuk melakukan deteksi, klasifikasi, dan clustering, yaitu *Support Vector Machine*, *Decission Tree*, *K-nearest Neighbor*, *K-Means*, dan *Long Short-Term Memory*[6],[7]. Metode-metode memiliki fungsi prediksi yang digunakan untuk mengelompokkan URL sebagai malicious dan benign. Dengan menggunakan machine learning dan deep learning, informasi URL dapat dianalisis dengan melakukan ekstraksi fitur seperti lexical based feature extraction, host based-feature extraction, dan content-based feature extraction yang selanjutnya muncul fitur yang akan digunakan dalam klasifikasi, identifikasi, analisis, dan deteksi pada malicious URL dan benign URL, Seperti URL length, domain name length, IP address, host-name URL. Berikut ini merupakan hasil penelitian menggunakan machine learning dan deep learning dengan feature extraction yang digunakan untuk mengidentifikasi malicious URL yang menjadi acuan utama dalam penelitian ini.

Metode Host-Based Feature Extraction adalah fitur ekstraksi pada situs web berbahaya pada perangkat yang bukan dari host web konvensional. Host-Based Feature dapat mengetahui informasi “siapa” pemiliknya, menggambarkan “dimana” asal situs hostingnya, dan “bagaimana” mereka dikelola yang data Properti didapat dari URL yang diidentifikasi oleh bagian hostname dari URL [8].

Pada penelitian[6] dengan judul *Empirical study on malicious URL detection using machine learning*, penelitian ini menggunakan fitur ekstraksi yaitu lexical feature dan host-based feature extraction. Dan dua metode machine learning yaitu SVM (Support Vector Machine) dan Random Forest. Penelitian tersebut dilakukan dengan menggunakan 10 iterasi dan menggunkana tiga pembagian rasio 60:40, 70:30, dan 80:20. Yang menghasilkan tingkat akurasi pada Random Forest lebih Besar dibandingkan SVM dan hasil plot dari Random Forest dengan variasi model yang dipakai 82 – 90%.

Pada penelitian[7] dengan judul URLdeepDetect: A Deep Learning Approach for Detecting Malicious URLs Using Semantic Vector Models menggunakan metode hybrid deep-learning yang dinamakan URLdeepDetect yang dapat melakukan analisis dan klasifikasi untuk medeteksi malicious URL. URLdeepdetect menggunakan mekanisme supervised dan unsupervised yaitu LSTM (Long Short-Term Memory) dan K-Means Clustering untuk mengklasifikasi URL. Pada penelitian ini mendapatkan hasil yang baik yaitu pada LSTM 98.3% dan pada Kmeans Clustering 99.7. Dalam penelitian ini, hasil yang didapatkan baik, tetapi pada metode K-means Clustering presisi dan recall yang masih rendah dibandingkan metode LSTM dapat ditingkatkan.

Metode LSTM digunakan untuk mendeteksi dan mempelajari perilaku informasi yang diolah pada sistem yang berbahaya. Pada penelitian ini penulis akan membahas mengenai analisis terhadap data URL yang telah di parser dari PDF Malware kedalam 2 kelas yaitu URL malicious, URL benign. Penggerjaan dari tugas akhir ini akan berfokus pada penggunaan fitur-fitur ekstraksi yaitu Host-Based Feature Extraction (DNS Record) yang kemudian dilakukan analisis karakteristik dari URL yang didapat dari fitur ekstraksi yang kemudian URL tersebut akan di kelompokkan berdasarkan cluster-cluster dan fitur-fitur yang digunakan.

Dari pembahasan sebelumnya dan metode yang digunakan pada penelitian terdahulu, maka penulis mengangkat judul *Klasifikasi Malicious URL Pada File Berbasis Host-Based Feature Extraction Menggunakan Metode LSTM*

dengan menggunakan dataset pada file pdf garuda yang telah diparser untuk mengambil URL yang ada didalam file pdf.

1.2. Perumusan Masalah

Dari latar belakang yang sudah dipaparkan diatas, muncul perumusan masalah pada penelitian yaitu: Dokumen PDF atau Word kadang tedapat alamat URL belum diidentifikasi berbahaya atau tidak bagi kita sebagai pengguna[9], [10]. Dan alamat URL asing tersebut bisa saja mengandung sebuah malware yang dapat menginfeksi perangkat yang dipakai. Terjadinya kejadian cyber berhubungan pada data pribadi yang ada di device yang terinfeksi tanpa diketahui oleh pengguna dan informasi pribadi dapat tersebar. Oleh karena itu, penulis menggunakan parserpdf dan virustotal untuk mengetahui url pada file pdf mengidentifikasi yang mana memiliki malware. Pada penelitian ini akan membahas cara mengekstrak fitur dari URL menggunakan *Host – Based Feature Extraction* kemudian melakukan klasifikasi menggunakan metode LSTM dan akan melakukan validasi terhadap nilai akurasi, presisi, recall, spesifitas.

1.3. Batasan Masalah

Berikut batasan masalah dari tugas akhir ini, yaitu :

1. Penelitian ini hanya sebatas melakukan klasifikasi dengan bantuan visualisasi dengan metode LSTM berdasarkan metode Fitur Ekstraksi menggunakan Host-Based Feature Extraction.
2. Penelitian ini hanya sebatas menggunakan program dengan bahasa pemrograman *Python*.
3. Penelitian ini hanya berdasarkan file PDF dari dataset GARUDA PDF.
4. Output yang dihasilkan dari penelitian ini hanya berupa klasifikasi dari hasil ekstraksi fitur dari URL yang divisualkan

dengan menggunakan bantuan machine learning yang digunakan sebagai tolak ukur untuk melakukan klasifikasi pada URL.

1.4. Tujuan

Tujuan dari penulisan tugas akhir ini antara lain:

1. Mengolah dan menggunakan Host-Based Feature Extraction pada file PDF Garuda.
2. Melakukan Visualisasi Data dari dataset Malicious URL dengan menggunakan metode LSTM.
3. Melakukan Analisis terhadap dataset Malicious URL berdasarkan Host-based Feature Extraction dan Klasifikasi dari dataset Malicious URL dengan metode LSTM.

1.5. Manfaat

Manfaat dari penulisan tugas akhir ini, yaitu :

1. Dapat menerapkan ekstraksi fitur pada URL dengan Host-Based Feature Extraction.
2. Dapat menerapkan algoritma LSTM yang digunakan dalam klasifikasi URL berdasarkan Host-Based Feature Extraction.
3. Dapat mengetahui karakteristik dari Benign URL, dan Malicious URL dengan menganalisis terhadap hasil dari klasifikasi URL dengan LSTM berdasarkan Host-Based Feature Extraction.

1.6. Metodologi Penelitian

Pada tugas akhir ini menggunakan metodologi sebagai berikut :

1. Persiapan Data

Pada tahap ini penulis melakukan pengumpulan dataset yang akan digunakan dalam pembelajaran dan pemahaman terhadap data yang akan diolah sehingga kebutuhan untuk topik penelitian dapat terpenuhi.

2. Studi Pustaka dan Literatur

Pada metode ini mencari dan mengumpulkan informasi sebagai referensi yang berupa literatur yang terdapat pada buku dan internet mengenai Analisis URL berdasarkan ekstraksi fitur URL yang kemudian dapat digunakan untuk menjadi topik bahasan pada penulisan tugas akhir ini.

3. Metode Pengujian

Pada metode ini membangung rancangan pengujian terhadap model yang telah dibuat, apakah simulasi tersebut dapat menghasilkan data akurasi yang diinginkan.

4. Analisa dan Kesimpulan

Hasil dari pengujian pada tugas akhir ini akan dianalisis kekurangannya, sehingga dapat digunakan untuk penelitian selanjutnya.

1.7. Sistematika Penulisan

Adapun Sistematika penulisan pada skripsi ini untuk menjelaskan isi dari setiap sub bab antara lain :

BAB I. PENDAHULUAN

Dalam bab I , menjelaskan tentang latar belakang, perumusan masalah, batasan masalah, tujuan dan manfaat serta sistematika penulisan dari pembahasan topik skripsi ini yaitu Klasifikasi serangan Malicious URL dengan metode LSTM.

BAB II. TINJAUAN PUSTAKA

Dalam bab II, menampilkan literature review yang berhubungan tentang pembahasan teori serangan Malicious URL, metode LSTM dan teori-teori lainnya yang berkaitan dengan skripsi ini.

BAB III. METODOLOGI

Dalam bab III, menjelaskan tahapan proses penelitian yang dilakukan secara terstruktur dengan menampilkan tahapan-tahapan pada persiapan dataset Malicious URL, kemudian penerapan metode LSTM guna memenuhi tujuan dari pembuatan skripsi ini.

BAB IV, HASIL DAN ANALISIS

Dalam bab IV, menampilkan hasil yang telah didapatkan dari tahapan yang telah dilakukan, serta untuk melihat performa sistem kemudian melakukan analisis dari metode LSTM.

DAFTAR PUSTAKA

- [1] G. Palaniappan, S. Sangeetha, B. Rajendran, Sanjay, S. Goyal, and B. S. Bindhumadhava, “Malicious Domain Detection Using Machine Learning on Domain Name Features, Host-Based Features and Web-Based Features,” *Procedia Comput. Sci.*, vol. 171, no. 2019, pp. 654–661, 2020, doi: 10.1016/j.procs.2020.04.071.
- [2] C. Do Xuan, H. D. Nguyen, and T. V. Nikolaevich, “Malicious URL detection based on machine learning,” *Int. J. Adv. Comput. Sci. Appl.*, vol. 11, no. 1, pp. 148–153, 2020, doi: 10.14569/ijacsa.2020.0110119.
- [3] B. Cui, S. He, P. Shi, and X. Yao, “Malicious URL detection with feature extraction based on machine learning,” *Int. J. High Perform. Comput. Netw.*, vol. 12, no. 2, pp. 166–178, 2018, doi: 10.1504/ijhpcn.2018.094367.
- [4] S. Sinha, M. Bailey, and F. Jahanian, “Shades of Grey: On the effectiveness of reputation-based blacklists,” *3rd Int. Conf. Malicious Unwanted Software, MALWARE 2008*, pp. 57–64, 2008, doi: 10.1109/MALWARE.2008.4690858.
- [5] J. Ispahany and R. Islam, “Detecting malicious COVID-19 URLs using machine learning techniques,” *2021 IEEE Int. Conf. Pervasive Comput. Commun. Work. other Affil. Events, PerCom Work. 2021*, no. January, pp. 718–723, 2021, doi: 10.1109/PerComWorkshops51409.2021.9431064.
- [6] R. Patgiri, H. Katari, R. Kumar, and D. Sharma, *Empirical study on malicious URL detection using machine learning*, vol. 11319 LNCS, no. January. Springer International Publishing, 2019.
- [7] S. Afzal, M. Asim, A. R. Javed, M. O. Beg, and T. Baker, “URLdeepDetect: A Deep Learning Approach for Detecting Malicious URLs Using Semantic Vector Models,” *J. Netw. Syst. Manag.*, vol. 29, no. 3, 2021, doi: 10.1007/s10922-021-09587-8.
- [8] J. Ma, L. K. Saul, S. Savage, and G. M. Voelker, “Beyond blacklists: Learning to detect malicious web sites from suspicious URLs,” *Proc. ACM SIGKDD Int. Conf. Knowl. Discov. Data Min.*, pp. 1245–1253, 2009, doi:

- 10.1145/1557019.1557153.
- [9] Y. Li, Y. Wang, Y. Wang, L. Ke, and Y. an Tan, “A feature-vector generative adversarial network for evading PDF malware classifiers,” *Inf. Sci. (Ny.)*, vol. 523, pp. 38–48, 2020, doi: 10.1016/j.ins.2020.02.075.
 - [10] G. Meng, M. Patrick, Y. Xue, Y. Liu, and J. Zhang, “Securing Android App Markets via Modeling and Predicting Malware Spread between Markets,” *IEEE Trans. Inf. Forensics Secur.*, vol. 14, no. 7, pp. 1944–1959, 2019, doi: 10.1109/TIFS.2018.2889924.
 - [11] S. Aarthi, N. V. Kishan, V. Surya Teja, N. V Harsha, and V. Gupta, “Classification of Phishing Website Based on URL Features,” *Int. J. Emerg. Technol. Eng. Res.*, vol. 7, no. 5, pp. 9–11, 2019, [Online]. Available: www.ijeter.everscience.org.
 - [12] S. Kumi, C. Lim, and S. G. Lee, “Malicious url detection based on associative classification,” *Entropy*, vol. 23, no. 2, pp. 1–12, 2021, doi: 10.3390/e23020182.
 - [13] M. Aljabri *et al.*, “An Assessment of Lexical, Network, and Content-Based Features for Detecting Malicious URLs Using Machine Learning and Deep Learning Models,” *Comput. Intell. Neurosci.*, vol. 2022, 2022, doi: 10.1155/2022/3241216.
 - [14] J. S. Ambata, J. Gaurana, D. Jacinto, and J. De Goma, “Malicious URL Classification Using Extracted Features, Feature Selection Algorithm, and Machine Learning Techniques,” *Proc. Int. Conf. Ind. Eng. Oper. Manag.*, no. 2016, pp. 2421–2429, 2021.
 - [15] Q. Wang, L. Li, B. J. B, Z. Lu, and J. Liu, *Based on K-means and SMOTE*, vol. 1. Springer International Publishing, 2020.
 - [16] S. Nayak, D. Nadig, and B. Ramamurthy, “Analyzing Malicious URLs using a Threat Intelligence System,” *Int. Symp. Adv. Networks Telecommun. Syst. ANTS*, vol. 2019-Decem, pp. 6–9, 2019, doi: 10.1109/ANTS47819.2019.9118051.
 - [17] F. O. Catak, K. Sahinbas, and V. Dörktardeş, “Malicious URL detection

- using machine learning,” *Artif. Intell. Paradig. Smart Cyber-Physical Syst.*, no. November, pp. 160–180, 2020, doi: 10.4018/978-1-7998-5101-1.ch008.
- [18] A. B. Sayamber and A. M. Dixit, “Malicious URL Detection and Identification,” *Int. J. Comput. Appl.*, vol. 99, no. 17, pp. 17–23, 2014, doi: 10.5120/17464-8247.
 - [19] D. Chiba, K. Tobe, T. Mori, and S. Goto, “Detecting malicious websites by learning IP address features,” *Proc. - 2012 IEEE/IPSJ 12th Int. Symp. Appl. Internet, SAINT 2012*, pp. 29–39, 2012, doi: 10.1109/SAINT.2012.14.
 - [20] M. Chakraborty and V. E. Balas, *The “Essence” of Network Security : An Panorama*. .
 - [21] T. Shibahara *et al.*, “Malicious URL sequence detection using event denoising convolutional neural network,” *IEEE Int. Conf. Commun.*, 2017, doi: 10.1109/ICC.2017.7996831.
 - [22] D. Sahoo, C. Liu, and S. C. H. Hoi, “Malicious URL Detection using Machine Learning: A Survey,” vol. 1, no. 1, pp. 1–37, 2017, [Online]. Available: <http://arxiv.org/abs/1701.07179>.
 - [23] A. K. Jain and B. B. Gupta, “A machine learning based approach for phishing detection using hyperlinks information,” *J. Ambient Intell. Humaniz. Comput.*, vol. 10, no. 5, pp. 2015–2028, 2019, doi: 10.1007/s12652-018-0798-z.
 - [24] E. Aghaei and G. Serpen, “Host-based anomaly detection using Eigentraces feature extraction and one-class classification on system call trace data,” no. December, 2019, [Online]. Available: <http://arxiv.org/abs/1911.11284>.
 - [25] D. Jia *et al.*, “An Electrocardiogram Delineator via Deep Segmentation Network,” *Proc. Annu. Int. Conf. IEEE Eng. Med. Biol. Soc. EMBS*, pp. 1913–1916, 2019, doi: 10.1109/EMBC.2019.8856987.
 - [26] P. Ongsulee, “Artificial intelligence, machine learning and deep learning,” *Int. Conf. ICT Knowl. Eng.*, pp. 1–6, 2018, doi: 10.1109/ICTKE.2017.8259629.
 - [27] A. Malali, S. Hiriyannaiah, G. M. Siddesh, K. G. Srinivasa, and N. T.

- Sanjay, “Supervised ECG wave segmentation using convolutional LSTM,” *ICT Express*, vol. 6, no. 3, pp. 166–169, 2020, doi: 10.1016/j.icte.2020.04.004.
- [28] P. Zhao and X. Yang, “Opportunistic routing for bandwidth-sensitive traffic in wireless networks with lossy links,” *J. Commun. Networks*, vol. 18, no. 5, pp. 806–817, 2016, doi: 10.1109/JCN.2016.000109.
- [29] A. Sherstinsky, “Fundamentals of Recurrent Neural Network (RNN) and Long Short-Term Memory (LSTM) network,” *Phys. D Nonlinear Phenom.*, vol. 404, p. 132306, 2020, doi: 10.1016/j.physd.2019.132306.
- [30] N. Gupta, V. Jindal, and P. Bedi, “LIO-IDS: Handling class imbalance using LSTM and improved one-vs-one technique in intrusion detection system,” *Comput. Networks*, vol. 192, no. February, p. 108076, 2021, doi: 10.1016/j.comnet.2021.108076.
- [31] A. Peimankar and S. Puthusserypady, “DENS-ECG: A deep learning approach for ECG signal delineation,” *Expert Syst. Appl.*, vol. 165, no. September 2020, p. 113911, 2021, doi: 10.1016/j.eswa.2020.113911.