# DeepPhish Simulating Malicious AI

Alejandro Correa Bahnsen, PhD

VP, AI & Research – Cyxtera Technologies

# About Me

- Industrial Engineer
- PhD in Machine Learning
- Passionate about open-source
- Scikit-Learn contributor
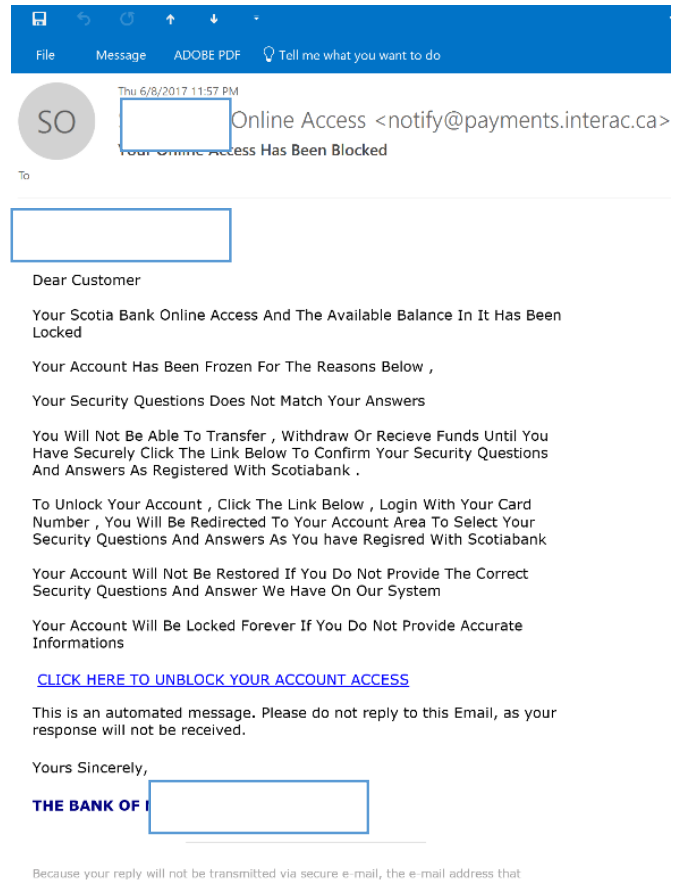- Organizer of Data Science Meetups

# Who I've worked with

# Agenda

- Phishing URL Detection using Machine Learning

- Malicious Cert Detection using Deep Learning

- DeepPhish: Simulating Malicious AI

- Demo 😰

# Typical Phishing Example

**91%**
of cybercrimes
and attacks
**start with a
phishing email**

# Why Phishing Detection is Hard?

Original Website

Only Using Images

Subtle Changes

# Ideal Phishing Detection System

# Ideal Phishing Detection System

## Issues with full content analysis:

- Time consuming

- Impractical to process millions of websites per day

- Hard to implement for small devices



WAITING FOR

PHISHING SCORE

imgflip.com

# There is always the need for an URL


`https://azenka.com.br/naza/wp-content/uploads/2018/01`

# Database of URLs

1,000,000 Phishing URLs from Phish Tank

http://moviesjingle.com/auto/163.com/index.php

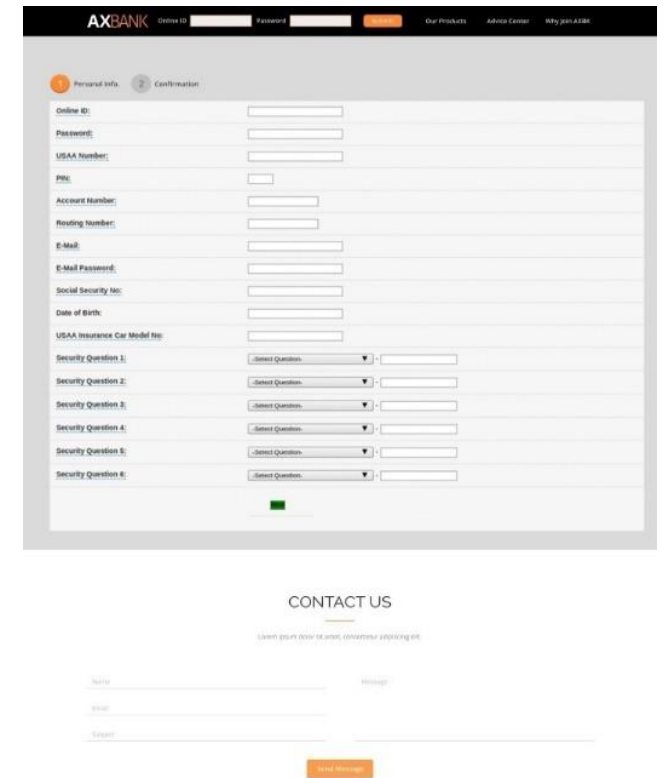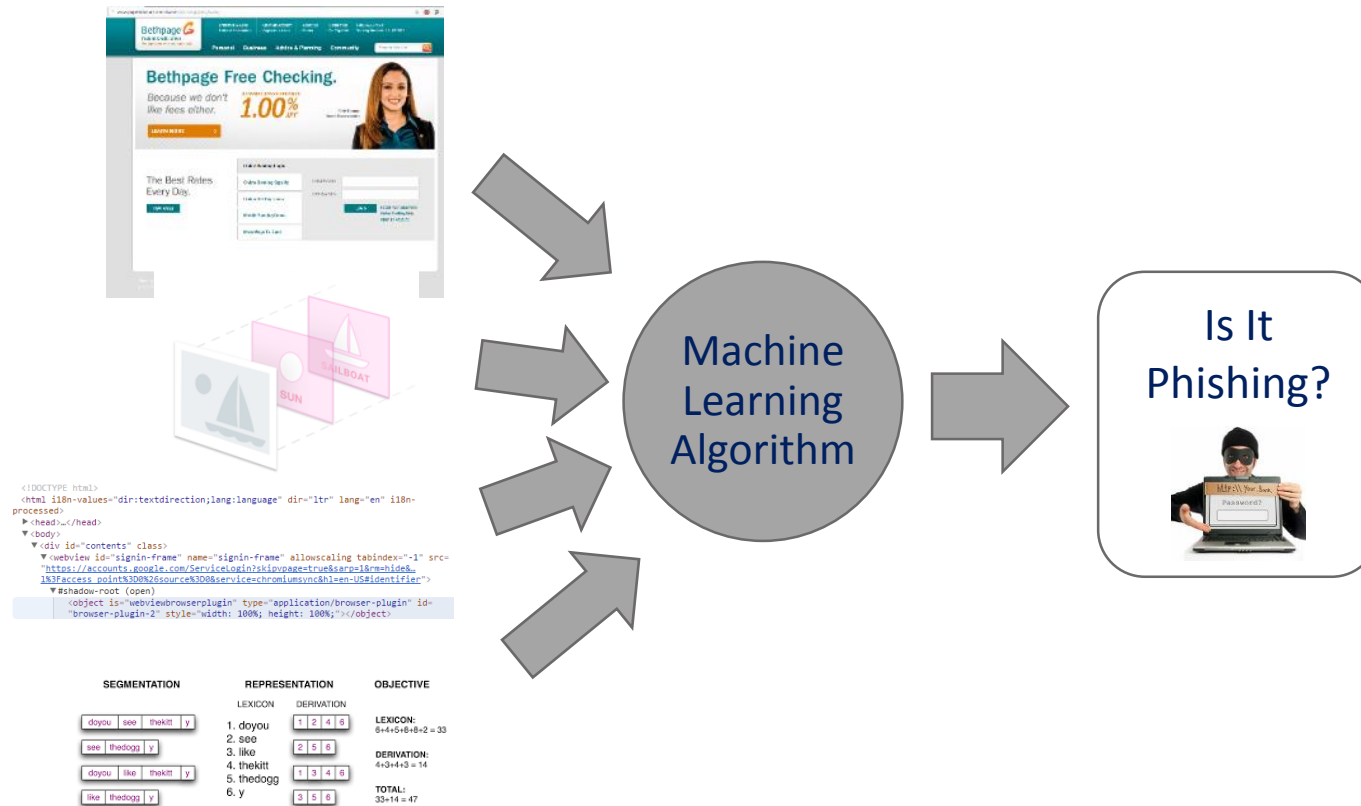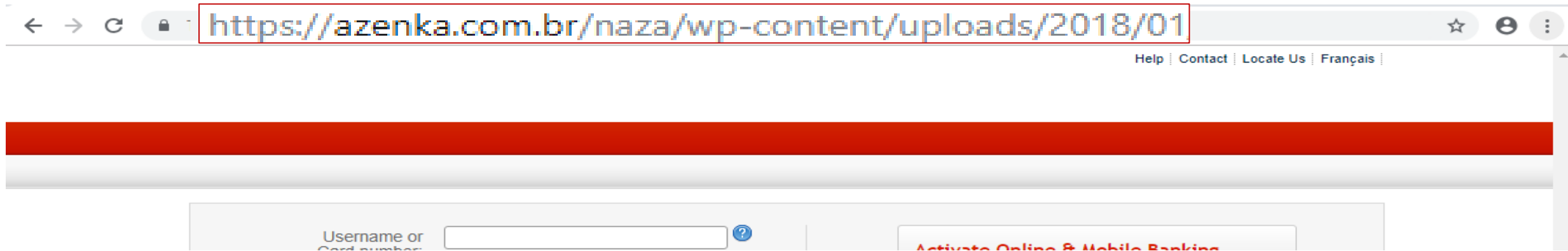http://**paypal.com**.update.account.**toughbook.cl**/8a30e847925afc5975161aeabe8930f1/?cmd=\_home\&dispatch=d09b78f5812945a73610edf38

http://msystemtech.ru/components/com\_users/Italy/zz/Login.php?run=\_login-submit\&session=68bbd43c854147324d77872062349924

1,000,000 Legitimate URLs from Common Crawl

https://www.sanfordhealth.org/ChildrensHealth/Article/73980

http://www.grahamleader.com/ci\_25029538/these-are-5-worst-super-bowl-halftime-shows\&defid=1634182

http://www.carolinaguesthouse.co.uk/onlinebooking/?industrytype=1\&startdate=2013-09-05\&nights=2\&location\&productid=25d47a24-6b74

# Recurrent Neural Networks RNN

# Recurrent Neural Networks RNN

# URL Classification Results

| 3-Fold CV | Accuracy | Recall | Precision |
|-----------|----------|--------|-----------|
| **Average** | 98.76% | 98.93% | 98.60% |
| **Deviation** | 0.04% | 0.02% | 0.02% |



Receiver operating characteristic

ROC fold 0 (area = 0.9990)
ROC fold 1 (area = 0.9991)
ROC fold 2 (area = 0.9991)
Luck

# Detecting Malicious URLs
# Is Not Enough!!

WIRED

Phishing Schemes Are Using Encrypted Sites to Seem Legit

LILY HAY NEWMAN  SECURITY  12.05.17  02:32 PM

# PHISHING SCHEMES ARE USING ENCRYPTED SITES TO SEEM LEGIT

## SHARE

SHARE
750

TWEET

COMMENT

EMAIL

## MOST POPULAR

SCIENCE
The Peculiar Math That Could Underlie the Laws of Nature
NATALIE WOLCHOVER

CULTURE
A Deadly Hunt for Hidden Treasure Spawns an Online Mystery
DAVID KUSHNER

CULTURE
MoviePass Raises Prices, Limits First-Run Availability as Pressures...
BRIAN BARRETT

GETTY IMAGES

Phishing Attacks using TLS Certificates

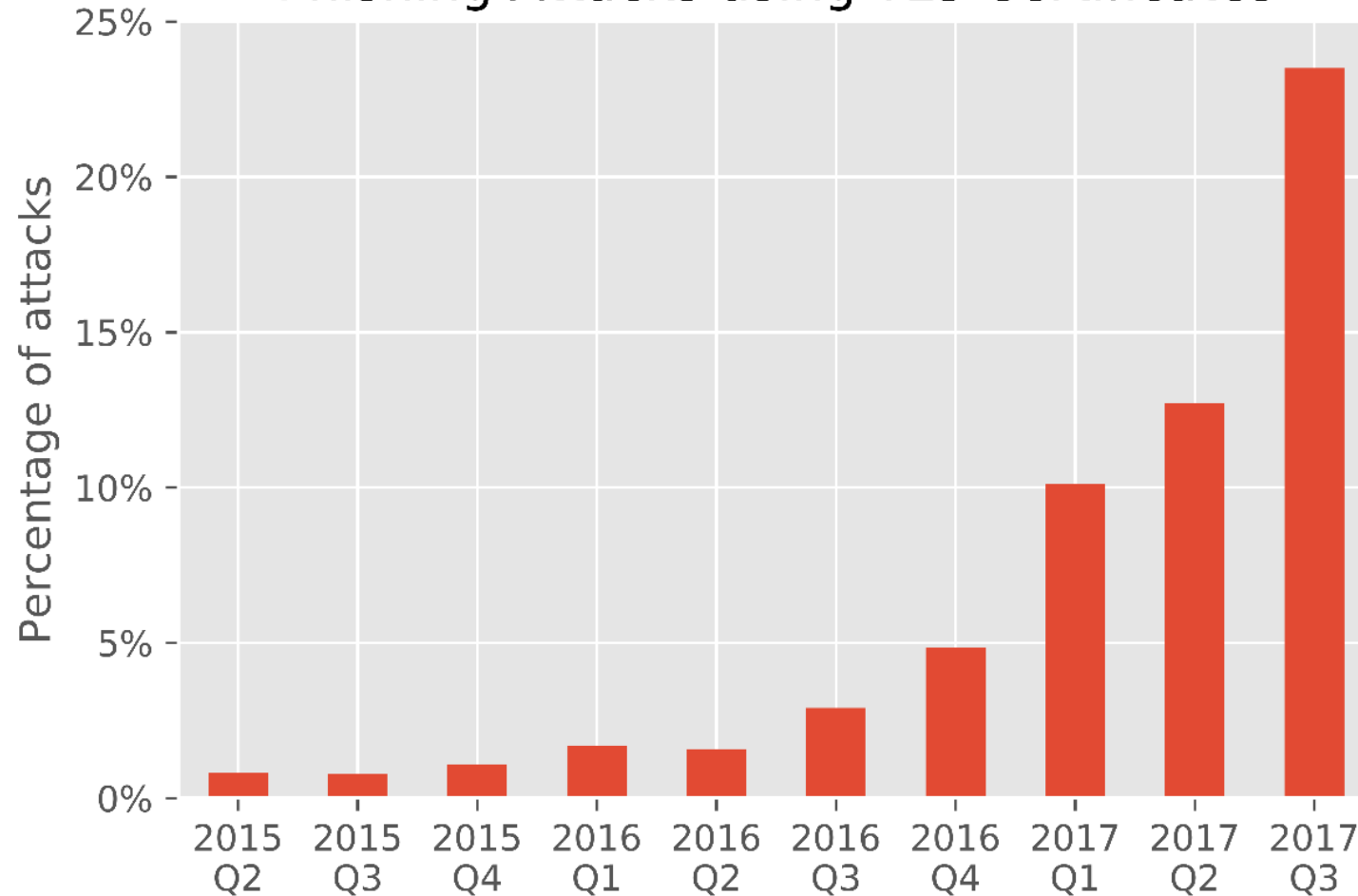# What is a Web Certificate?

Forrester survey asked users: "Some websites receive the following browser user interface security indicator in the browser. What do you think the security indicator is intended to tell users?"

🔒 Secure | https://ultrabank.com

The website is safe: 82%

The website is encrypted: 75%

The website is trustworthy: 66%

The website is private: 32%

# Hunting Malicious TLS Certificates with Deep Neural Networks

# Database of TLS Certificates

1,000,000 Legitimate Certificates from Common Crawl

CN = *.stackexchange.com, O = Stack Exchange, Inc., L = New York, S = NY, C = US

CN = slack.com, O = Slack Technologies, Inc., L = San Francisco, S = CA, C = US

CN = *.trello.com, O = Trello Inc., L = New York, S = New York, C = US
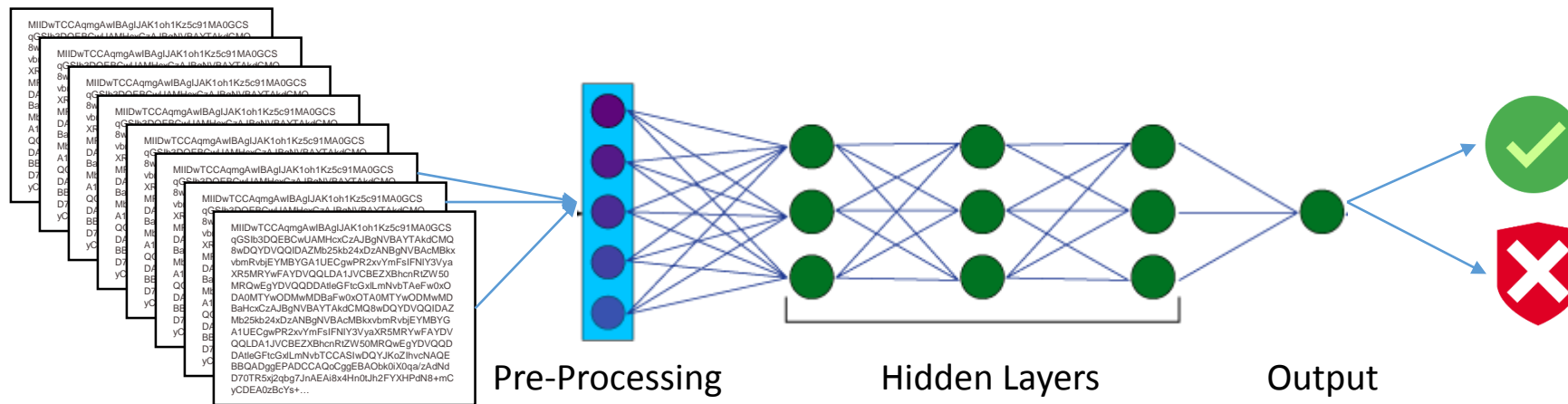
5,000 Phishing Certificates

CN = localhost, L = Springfield

CN = localhost.localdomain

CN = example.com, L = Springfield

# Deep Learning Algorithm



Pre-Processing          Hidden Layers         Output

# Deep Learning Algorithm



Subject Principal → One hot encoding → Embedding → LSTM → Dropout

Issuer Principal → One hot encoding → Embedding → LSTM → Dropout

Extracted Features → Dense/ReLu → Dropout

→ Concatenate → Dense/ReLu → Dropout → Dense/Logit → score

# Malicious Cert Classification Results

| 5-Fold CV | Accuracy | Recall | Precision |
|-----------|----------|--------|-----------|
| **Average** | 86.41% | 83.20% | 88.86% |
| **Deviation** | 1.22% | 3.29% | 1.04% |



Precision-Recall curve

EVERY
ACTION
HAS A
FRAUDSTER
REACTION

# DeepPhish Simulating Malicious AI

# The Experiment:  Simulating Malicious AI

**1** Identify individual threat actors

**2** Run them through our own AI detection system

**3** Improve their attacks using AI

# Uncovering Threat Actors

- Objective: We want to understand effective patterns of each attacker to improve them through a AI model

- As we can not know attackers directly, we must learn from them through their attacks

- Database with 1.1M confirm phishing URLs collected from Phishtank

# Threat Actor 1

naylorantiques.com



## 406 URLs

http://naylorantiques.com/components/com_**contact**/views/**contact**/tmpl/62

http://naylorantiques.com/docs/Auto/**Atendimento**/5BBROPI6S3

http://naylorantiques.com/**Atualizacao Segura**/pictures/XG61YYMT_FXW0PWR8_5P2O7T2U_P9HNDPQR/

http://naylorantiques.com/zifn3p72bsifn9hx9ldecd8jzl2f0xlwf8f

http://www.naylorantiques.com/JavaScript/charset=iso-8859-1/http-equiv/margin-bottom

## Keywords

atendimento, jsf, identificacao, ponents, views, TV, mail, SHOW, COMPLETO, VILLA, MIX, ufi, pnref, story, tryy2ilr, Autentico

Check in database

## 106 domains

naylorantiques.com, netshelldemos.com, debbiebright.co.z, waldronfamilygppractice.co.uk , avea-vacances.com , psncodes2013.com uni5.net , 67.228.96.204, classificadosmaster.com.br, ibjjf.org

Visual Check

Visual Check

# The Experiment:    Simulating Malicious AI

**1** Identify individual threat actors

**2** Run them through our own AI detection system

**3** Improve their attacks using AI

DEMO!!!

# The Experiment: Simulating Malicious AI

**1** Identify individual threat actors

**2** Run them through our own AI detection system

**3** Improve their attacks using AI

# DeepPhish Algorithm - Training

# DeepPhish LSTM Network

# DeepPhish Algorithm - Prediction



Model → Predict → Next Character Iteratively → Filter paths → Allowed Paths (Compromised Domains +) → Create → Synthetic URLs

/arendipemto/nenu-opcines-fone-facil vfone/faci/Atondime

http:// + www.naylorantiques.com + /arendipemto/nenu-opcines-fone-facilvone/facil/Atondime

#BHEU / @BLACK HAT EVENTS

# What's Next??

# What's Next??

AI powered Attacks are real, as we probed with Deep Phish experiment.
We need to enhance our own AI detection systems to account for the possibility of attackers using AI.

# DeepPhish: Simulating Malicious AI

Alejandro Correa Bahnsen, Iva[...]
Cyber Thr[...]
Email: {alejandro.correa, ivan.[...]

*Abstract*—In this work we describe how threat [...] use AI algorithms to bypass AI phishing detection s[...] analyzed more than a million phishing URLs to und[...] different strategies that threat actors use to create phi[...] Assuming the role of an attacker, we simulate how diff[...] actors may leverage Deep Neural Networks to en[...] effectiveness rate. Using Long Short-Term Memory [...] we created DeepPhish, an algorithm that learns to c[...] phishing attacks. By training the DeepPhish algorit[...] different threat actors, they were able to increase the[...] ness from 0.69% to 20.9%, and 4.91% to 36.28%, [...]

*Keywords*—*Malicious AI; phishing detection; cyb[...] current neural networks; long-short term memory net[...] adversarial learning.*

I. INTRODUCTION

Machine Learning (ML) and Artificial Intelli[...]

# Classifying Phishing URLs Using Recurrent [...]

Alejandro Correa Bahnsen[†], Edu[...]
Javier Vargas[...]
[†]Easy[...]
*MindLab Research Group, [...]
Email: acorrea@easysol.net, econtrerasb@unal.edu.co, [...]

*Abstract*—As the technical skills and costs associated [...] the deployment of phishing attacks decrease, we are witne[...] an unprecedented level of scams that push the need for l[...] methods to proactively detect phishing threats. In this [...] we explored the use of URLs as input for machine lea[...] models applied for phishing site prediction. In this way[...] compared a feature-engineering approach followed by a ra[...] forest classifier against a novel method based on recurrent n[...] networks. We determined that the recurrent neural net[...] approach provides an accuracy rate of 98.7% even witho[...] need of manual feature creation, beating by 5% the random [...] method. This means it is a scalable and fast-acting pro[...] detection system that does not require full content analysis[...]

*Keywords*—*Phishing detection; Cybercrime; Feature eng[...] ing; Recurrent neural networks; Long short term memory netw[...]*

I. INTRODUCTION

# Hunting Malicious TLS Certificates with Deep Neural Networks

Ivan Torroledo
Cyxtera Technologies
ivan.torroledo@cyxtera.com

Luis David Camacho
Cyxtera Technologies
luis.camacho@cyxtera.com

Alejandro Correa Bahnsen
Cyxtera Technologies
alejandro.correa@cyxtera.com

**ABSTRACT**

Encryption is widely used across the internet to secure communications and ensure that information cannot be intercepted and read by a third party. However, encryption also allows cybercriminals to hide their messages and carry out successful malware attacks while avoiding detection. Further aiding criminals is the fact that web browsers display a green lock symbol in the URL bar when a connection to a website is encrypted. This symbol gives a false sense of security to users, who are in turn more likely to fall victim to phishing attacks. The risk of encrypted traffic means that information security researchers must explore new techniques to detect, classify, and take countermeasures against malicious traffic. So far there exists no approach for TLS detection in the wild. In this paper, we propose a method for identifying malicious use of web certificates using deep neural networks. Our system uses the content of TLS certificates to successfully identify legitimate certificates as well as malicious patterns used by attackers. The results show that our system is capable of identifying malware certificates with an accuracy of 94.87% and phishing certificates with an accuracy of 88.64%.
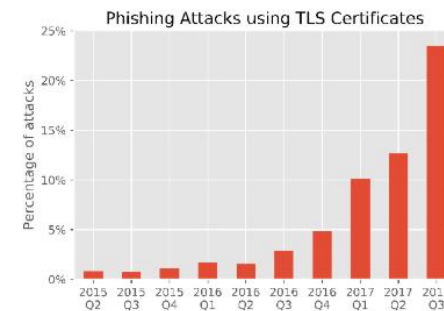
Phishing Attacks using TLS Certificates

Figure 1: Evolution of phishing attacks using TLS [16].

Secure | https://ultrabank.com

# black hat®
# EUROPE 2018
## DECEMBER 3-6, 2018
### EXCEL LONDON / UNITED KINGDOM

# Thanks!!

Alejandro Correa Bahnsen, PhD

VP, AI & Research

alejandro.correa@cyxtera.com

@albahnsen