# How AI Is Changing the Security of Software Systems

Prof. Guido Salvaneschi

MSc Daniel Sokolowski

MSc David Spielmann

# University of St.Gallen

## AI in Cybersecurity
## Focus Group

Monitor, evaluate, and formulate guidelines
for the adoption of AI technologies in cybersecurity

# An Intrusion-Detection Model

DOROTHY E. DENNING

*Abstract*—A model of a real-time intrusion-detection expert system capable of detecting break-ins, penetrations, and other forms of computer abuse is described. The model is based on the hypothesis that security violations can be detected by monitoring a system's audit records for abnormal patterns of system usage. The model includes profiles for representing the behavior of subjects with respect to objects in terms of metrics and statistical models, and rules for acquiring knowledge about this behavior from audit records and for detecting anomalous behavior. The model is independent of any particular system, application environment, system vulnerability, or type of intrusion, thereby providing a framework for a general-purpose intrusion-detection expert system.

*Index Terms*—Abnormal behavior, auditing, intrusions, monitoring, profiles, security, statistical measures.

## I. INTRODUCTION

THIS paper describes a model for a real-time intrusion-detection expert system that aims to detect a wide range of security violations ranging from attempted break-ins by outsiders to system penetrations and abuses by insiders. The development of a real-time intrusion-detection system is motivated by four factors: 1) most existing systems have security flaws that render them susceptible to intrusions, penetrations, and other forms of abuse; finding and fixing all these deficiencies is not feasible for

ging into a system through an unauthorized account and password might have a different login time, location, or connection type from that of the account's legitimate user. In addition, the penetrator's behavior may differ considerably from that of the legitimate user; in particular, he might spend most of his time browsing through directories and executing system status commands, whereas the legitimate user might concentrate on editing or compiling and linking programs. Many break-ins have been discovered by security officers or other users on the system who have noticed the alleged user behaving strangely.

- *Penetration by legitimate user:* A user attempting to penetrate the security mechanisms in the operating system might execute different programs or trigger more protection violations from attempts to access unauthorized files or programs. If his attempt succeeds, he will have access to commands and files not normally permitted to him.

- *Leakage by legitimate user:* A user trying to leak sensitive documents might log into the system at unusual times or route data to remote printers not normally used.

- *Inference by legitimate user:* A user attempting to obtain unauthorized data from a database through aggregation and inference might retrieve more records than usual.

# Evaluating Intrusion Detection Systems:
# The 1998 DARPA Off-line Intrusion Detection Evaluation*

Richard P. Lippmann, David J. Fried, Isaac Graf, Joshua W. Haines, Kristopher R. Kendall,  David McClung,
Dan Weber, Seth E. Webster, Dan Wyschogrod, Robert K. Cunningham, and Marc A. Zissman
*Lincoln Laboratory MIT, 244 Wood Street, Lexington, MA 02173-9108*
*Email: rpl@SST.LL.MIT.EDU or jhaines@SST.LL.MIT.EDU*

## ABSTRACT

*A intrusion detection evaluation test bed was developed which generated normal traffic similar to that on a government site containing 100's of users on 1000's of hosts. More than 300 instances of 38 different automated attacks were launched against victim UNIX hosts in seven weeks of training data and two weeks of test data. Six research groups participated in a blind evaluation and results were analyzed for probe, denial-of-service (DoS), remote-to-local (R2L), and user to root (U2R) attacks. The best systems detected old attacks included in the training data, at moderate detection rates ranging from 63% to 93% at a false alarm rate of 10 false alarms per day. Detection rates were much worse for new and novel R2L and DoS attacks included only in the test data. The best systems*

between software compone[...]
continually exploited by atta[...]
occur despite the best secur[...]
systems have become an [...]
security to detect these at[...]
damage. A review of curr[...]
is available in [1]. Some a[...]
and can stop an attack in prog[...]
information about attacks and can [...]
understand the attack mechanism, and reduce the possib[...]
future attacks of the same type. More advanced intrusion detection systems detect never-before-seen, new, attacks, while the more typical systems detect previously seen, known attacks.

Evaluations of developing technologies such as those used for intrusion detection are essential to focus effort, document existing capabilities, and guide research. For example, yearly DARPA-sponsored evaluations in the speech recognition area
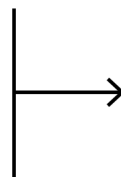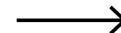
# Machine Learning for Malware Detection

kaspersky

## Training phase

Benign executables

Malicious executables

→ Training → Predictive model

## Protection phase

Unknown executable → Processing by a predictive model → **Malicious** / **Benign**

Model decision

**Machine Learning: detection algorithm lifecycle**

# Dos and Don'ts of Machine Learning in Computer Security

Daniel Arp[*], Erwin Quiring[†], Feargus Pendlebury[‡§], Alexander Warnecke[†], Fabio Pierazzi[‡],
Christian Wressnegger[¶], Lorenzo Cavallaro[‖], Konrad Rieck[†]

[*] *Technische Universität Berlin*
[†] *Technische Universität Braunschweig*
[‡] *King's College London,* [‖] *University College London*
[§] *Royal Holloway, University of London and The Alan Turing Institute*
[¶] *KASTEL Security Research Labs and Karlsruhe Institute of Technology*

## Abstract

With the growing processing power of computing systems
and the increasing availability of massive datasets, machine
learning algorithms have led to major breakthroughs in many
different areas. This development has influenced computer
security, spawning a series of work on learning-based security
systems, such as for malware detection, vulnerability discov-
ery, and binary code analysis. Despite great potential, machine
learning in security is prone to subtle pitfalls that undermine
its performance and render learning-based systems potentially
unsuitable for security tasks and practical deployment.

In this paper, we look at this problem with critical eyes.
First, we identify common pitfalls in the design, implementa-

and addressing security-relevant p
application domains, including in
malware analysis [69, 88], vulner
and binary code analysis [42, 114

Machine learning, however, has
requires reasoning about statistic
a fairly delicate workflow: incor
mental biases may cast doubts o
that it becomes unclear whether we can trust
coveries made using learning algorithms at all [56]. Attempts
to identify such challenges and limitations in specific secu-
rity domains, such as network intrusion detection, started two
decades ago [11, 119, 126] and were extended more recently
to other domains, such as malware analysis and website fin-

# Introducing ChatGPT

We've trained a model called ChatGPT which interacts in a conversational way. The dialogue format makes it possible for ChatGPT to answer followup questions, admit its mistakes, challenge incorrect premises, and reject inappropriate requests.

Try ChatGPT ↗    Read about ChatGPT Plus

**SC MEDIA**
A CyberRisk Alliance Resource

TOPICS        EVENTS        PODCASTS        RESEARCH        RECOGNITION        LEADERSHIP

Security Program Controls/Technologies, Vulnerability Management

How ChatGPT is changing the way cybersecurity practitioners look at the potential of AI

Derek B. Johnson    December 9, 2022

# AI on offense: Can ChatGPT be used for cyberattacks?

The difference between clever and intelligence

By **Apoorva Joshi, Devon Kerr**

24 May 2023

Table of contents  ☰

**The Guardian**

News  Opinion  Sport  Culture  Lifestyle

World  UK  Climate crisis  Environment  Science  Global development  Football  **Tech**  Business

**Chatbots**

# AI chatbots making it harder to spot phishing emails, say experts

**Poor spelling and grammar that can help identify fraudulent attacks being rectified by artificial intelligence**

**CNBC**

TECHNOLOGY EXECUTIVE COUNCIL

# A.I. is helping hackers make better phishing emails

PUBLISHED THU, JUN 8 2023·9:55 AM EDT

Bob Violino          WATCH LIVE

**KEY POINTS**

- Cyber criminals and other bad actors can do things faster and easier with artificial intelligence, which makes it more difficult for security experts to

**Forbes**

FORBES > INNOVATION > CYBERSECURITY

# Almost Human: The Threat Of AI-Powered Phishing Attacks

Emil Sayegh Contributor ⓘ
*CEO of Ntirety. Cover all things cloud, cybersecurity &*

things organizations
against AI-assisted
services that

programs so that

Follow

Microsoft | Microsoft 365    Products ∨    Plans and pricing    Resources ∨    Support ∨          All Microsoft ∨    🔍    ⊕

Microsoft 365 Life Hacks > Privacy & Safety > How AI is changing phishing scams

July 14, 2023

# How AI is changing phishing scams

AI language models are the hottest tech of the year, rushing people to find new, exciting ways to use it to improve their day-to-day. But just as you can use fire to cook a meal or burn down a house, you can use AI to book a trip...or initiate a **phishing attack**. Brace yourself for era of phishing schemes; and they'll only grow more sophisticated

YOUR PER
RECOVER

DEVICE

Spear phishing and voice cloning

# TECH MONITOR 30
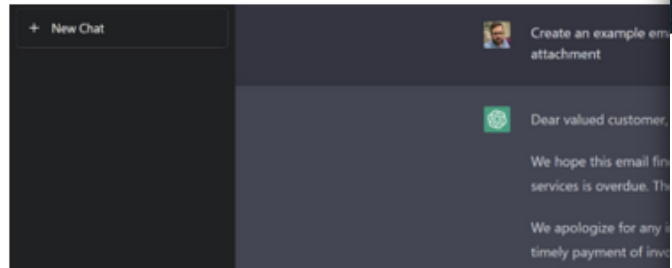PROUDLY CELEBRATING 30 YEARS OF INDEPENDENT IT JOURNALISM

**TECHNOLOGY** > **AI AND AUTOMATION** | December 19, 2022 | updated 09 Mar 2023 10:05am

## Here's how OpenAI's ChatGPT can be used to launch cyberattacks

Security researchers had the AI create a fake email from a hosting company and inject malware into an Excel file as part of a test.

By Ryan Morrison

Since its release at the end of November
ways to put OpenAI's advanced chatb
vendor has warned hackers could be
cyberattacks.

+ New Chat

Create an examp
attachment

Dear valued custo

We hope this ema
services is overdu

We apologize for a
timely payment of

# CSO

## 5 ways threat actors can use ChatGPT to enhance attacks

News

Apr 28, 2023 • 6 mins

| Artificial Intelligence | Cyberattacks | Threat and Vulnerability Management |

New research details how attackers can use AI-driven systems like ChatGPT in different aspects of cyberattacks phishing, and developing

- Social Engineering
- Attack point enumeration
- Foothold assistance
- Producing malicious code

# TECH MONITOR 30
PROUDLY CELEBRATING 30 YEARS OF INDEPENDENT IT JOURNALISM

TECHNOLOGY > AI AND AUTOMATION | December 19, 2022 | updated 09 Mar 2023 10:05am

## Here's how OpenAI's ChatGPT can be used to launch cyberattacks

Security researchers had the AI create a fake email from a hosting company and inject malware into an Excel file as part of a test.

By Ryan Morrison

Since its release at the end of November ways to put OpenAI's advanced chatbot vendor has warned hackers could be usi cyberattacks.

### CSO

## 5 ways threat actors can use ChatGPT to enhance attacks

News

Apr 28, 2023 • 6 mins

[Artificial Intelligence] [Cyberattacks] [Threat and Vulnerability Management]

...ers can use AI-driven
...t aspects of cyberattacks
...ing, and developing

### Infθsecurity Magazine

NEWS  13 DEC 2022

## Experts Warn ChatGPT Could Democratize Cybercrime

**Phil Muncaster**
UK / EMEA News Reporter, Infosecurity Magazine

Email Phil    Follow @philmuncaster

A wildly popular new AI bot could be used by would-be cyber-criminals to teach

Analytics Insight

**ChatGPT . Cyberattack . Latest News**

# Cybercriminals are Using ChatGPT to Create Hacking Tools and Code

Zaveria
January 11, 2023 . 2 mins read

## Experienced and novice cybercrimina[ls] using ChatGPT to create hacking too[ls]

Security researchers have reported that both experienced and novice cybercrimi[nals]
to create hacking tools and code.

One such instance is the Israeli security firm Check Point, which discovered a thr[...]
underground hacking site by a hacker who claimed to be testing the famous AI c[...]
malware strains".

The hacker later compressed and distributed Android malware created by ChatG[...]
internet. According to Forbes, spyware has the power to steal important files.

The same hacker also demonstrated another program that could install a backdo[...]

ITBrief
**AUSTRALIA**

## Five ways cybercriminals are making use of ChatGPT

By Anthony Daniel

Forcepoint

Home : Blogs : I built a Zero Day virus with undetectable exfiltration using only ChatGPT prompts

**X-Labs**

April 4, 2023 | 18 min read

# I built a Zero Day virus with undetectable exfiltration using only ChatGPT prompts

**Aaron Mulgrew**

( artificial intelligence ) ( chatgpt ) ( malware ) ( zero day )

OPEN SOURCE SOFTWARE   PENETRATION TESTING

# PentestGPT – Automate Penetration Testing Empowered by ChatGPT

BY **PRIYANSHU SAHAY** · MAY 15, 2023 · 5 MINUTE READ



> _PentestGPT /.

- Automate penetration testing tool empowered by ChatGPT.
- An interactive mode to guide

## Getting pwn'd by AI: Penetration Testing with Large Language Models

Andreas Happe
andreas.happe@tuwien.ac.at
TU Wien
Austria

Jürgen Cito
juergen.cito@tuwien.ac.at
TU Wien
Austria

**ABSTRACT**

The field of software security testing, more specifically penetration testing, requires high levels of expertise and involves many manual testing and analysis steps. This paper explores the potential use of large-language models, such as GPT3.5, to augment penetration testers with AI sparring partners. We explore two distinct use cases: high-level task planning for security testing assignments and low-level vulnerability hunting within a vulnerable virtual machine. For the latter, we implemented a closed-feedback loop between LLM-generated low-level actions with a vulnerable virtual machine (connected through SSH) and allowed the LLM to analyze the machine state for vulnerabilities and suggest concrete attack vectors which were automatically executed within the virtual machine. We discuss promising initial results, detail avenues for improvement, and close deliberating on the ethics of AI sparring partners.

**CCS CONCEPTS**

when stuck. The study also emphasizes that intuition is a big part of detecting vulnerabilities and that knowledge transfer, e.g., from attending Capture-the-Flag[1] (CTF) events, were seen as potential sources of this intuition — can this be partially outsourced to AI models? Using AI-based agents as sparring partners would augment and empower existing human security testers and could counteract the lack of sufficiently educated security professionals. Combining human operators with AIs creates new capabilities instead of cloning existing ones. Furthermore, keeping a human in the loop reduces the potential ethical problems imposed by the use of AIs [6]. Recent research indicates that the efficiency gains provided by the use of AI-based systems are greatest for low-skilled workers [7], augmenting human operators with a generative AI might thus also benefit the training of novice penetration testers.

**RQ: To what extent can we automate security testing with LLMs?** The rest of this paper explores whether large-language models can be deployed as sparring partners for security profes-

# OpenAI's ChatGPT can write impressive code. Here are the prompts you should use for the best results, experts say.

**Beatrice Nolan**  Aug 10, 2023, 1:07 PM GMT+2



**OpenAI's ChatGPT has caused quite a stir in the tech community.**  Getty/Luis Alvarez

- **OpenAI's ChatGPT has been able to produce working lines of code.**

- **The AI-powered bot has freaked out programmers and caught the attention of tech CEOs.**

# OpenAI's ChatGPT can write impressive code. Here are the prompts you should use for the best results, experts say.

Beatrice Nolan  Aug 10, 2023, 1:07 PM GMT+2

# Your AI pair programmer

Push what's possible with GitHub Copilot, the world's most widely adopted AI developer tool.

**Start my free trial >**   **Compare plans**

```
TS sentiments.ts      write_sql.go      parse_expenses.py      addresses.rb

1  #!/usr/bin/env ts-node
2
3  import { fetch } from "fetch-h2";
4
5  // Determine whether the sentiment of text is positive
6  // Use a web service
7  async function isPositive(text: string): Promise<boolean> {
8    const response = await fetch(`http://text-processing.com/api/sentiment/`, {
9      method: "POST",
10     body: `text=${text}`,
11     headers: {
```

**OpenAI's ChatGPT has caused quite a stir in the tech community.**  Getty/Luis Alvarez

- OpenAI's ChatGPT has been able to produce working lines of code.

- The AI-powered bot has freaked out programmers and caught the attention of tech CEOs.

# ChatGPT CVE Analysis for Red and Blue Team

Red Team use: use ChatGPT to help exploit the CVE and find vulnerabilities in the code. Blue Team use: explain the CVE and how to defend against it.

David Merian · Follow

Published in System Weakness · 2 min read · Mar 9

👏 1          💬                              🔖  ▶  📤

Many CVE's are popular targets for Ransomware, according to a report from Securin, as summarized on DarkReading. New CVE's are published everyday, and oftentimes, they are in very specific—but ubiquitous —software. You can use ChatGPT to summarize what the CVE is, what it

For your defense:
Use AI, practice it,
and get trained!

For your defense:
Use AI, practice it,
and get trained!

November 2018

Harvard Business Review
ANALYTIC SERVICES

Pulse Survey

ARTIFICIAL INTELLIGENCE
The End of the Beginning

https://freesvg.org/uncle-sam-for-dailysketch, 13.09.2023, modified

For your defense:
Use AI, practice it,
and get trained!

November 2018

**Harvard Business Review** ANALYTIC SERVICES

**ARTIFICIAL INTELLIGEN**

**The End of th**

**MITSloan** Management Review  ≡ MENU

**From ChatGPT to HackGPT: Meeting the Cybersecurity Threat of Generative AI**

It's time to replace traditional, rule-based approaches to cybersecurity with "smarter" technology and training.

Karen Renaud, Merrill Warkentin, and George Westerman • April 18, 2023

WILL KNIGHT    SECURITY    AUG 1, 2023 7:00 AM

# A New Attack Impacts Major AI Chatbots—and No One Knows How to Stop It

**Researchers found a simple way to make ChatGPT, Bard, and other chatbots misbehave, proving that AI is hard to tame.**

# Employees Are Feeding Sensitive Biz Data to ChatGPT, Raising Security Fears

More than 4% of employees have put sensitive corporate data into the large language model, raising concerns that its popularity may result in massive leaks of proprietary information.

**Robert Lemos**
Contributing Writer, Dark Reading

March 07, 2023

# Are We Ready to Embrace Generative AI for Software Q&A?

Bowen Xu[*†], Thanh-Dat Nguyen[‡], Thanh Le-Cong[‡], Thong Hoang[§], Jiakun Liu[†],
Kisub Kim[†], Chen Gong[¶], Changan Niu[∥], Chenyu Wang[†], Bach Le[‡], David Lo[†]

*North Carolina State University, USA
bxu22@ncsu.edu

†Singapore Management University, Singapore
{bowenxu.2017, jkliu, kisubkim, chenyuwang, davidlo}@smu.edu.sg

‡University of Melbourne, Austrialia
{thanhdatn, congthanh.le}@student.unimelb.edu.au, bach.le@unimelb.edu.au

§CSIRO's Data61, Australia          ¶University of Virginia, USA          ∥Nanjing University, China
james.hoang@data61.csiro.au          fzv6en@virginia.edu          niu.ca@outlook.com

*Abstract*—**Stack Overflow, the world's largest software Q&A (SQA) website, is facing a significant traffic drop due to the emergence of generative AI techniques. ChatGPT is banned by Stack Overflow after only 6 days from its release. The main reason provided by the official Stack Overflow is that the answers generated by ChatGPT are of low quality. To verify this, we conduct a comparative evaluation of human-written and ChatGPT-generated answers. Our methodology employs both automatic comparison and a manual study. Our results suggest that human-written and ChatGPT-generated answers are semantically similar, however, human-written answers outperform ChatGPT-generated ones consistently across multiple aspects, specifically by 10% on the overall score. We release the data, analysis scripts, and detailed results at https://github.com/maxxbw54/GAI4SQA.**

## I. INTRODUCTION

On November 30, 2022, OpenAI, a world-class AI company, launched an artificial intelligence chatbot named ChatGPT [1].

software question answering (SQA). In this booming era of AI-powered chatbots, traffic to OpenAI's ChatGPT has been growing exponentially, while traditional Q&A site such as Stack Overflow has been experiencing a steady decline [5]. Specifically, traffic to Stack Overflow was down by 6% every month in January 2022 on a year-over-year basis and was down 13.9% in March 2022 [6]. This phenomenon, however, is concerning due to the lack of empirical evidence on a comparative study on human-written vs AI-generated responses. The empirical evidence is much needed to ensure a balanced and robust development in the field of SQA. In this work, we investigate the following research questions:

- **RQ1:** *What are the characteristics of ChatGPT-generated and human-written answers?*
- **RQ2:** *From the human user perspective, how good are*

## The New York Times

Smart Ways to Use Chatbots     ChatGPT's Code Interpreter     Can A.I. Be Fooled?     A.I.'s Literary Skills

### In U.S., Regulating A.I. Is in Its 'Early Days'

While there has been a flurry of activity by the White House and lawmakers over artificial intelligence, rules for the technology remain distant, lawmakers and experts said.
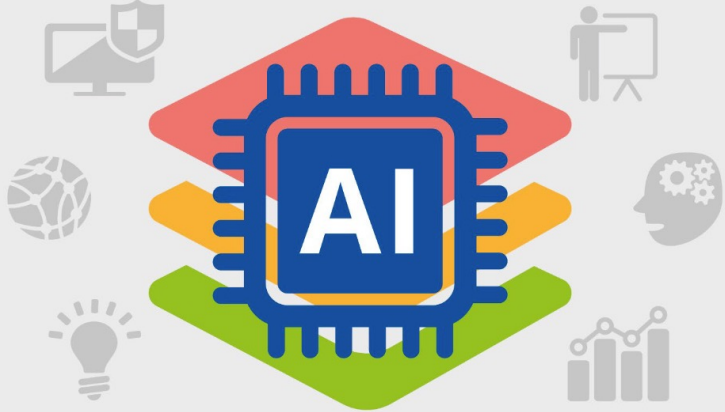
## Holistic AI

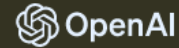# The State of Global AI Regulations in 2023

**January 2023**