**Prof. Phil Legg**

**Professor in Cyber Security**

**Co-Director: UWEcyber ACE-CSE**

March 2025

# Artificial Intelligence, agents, and what does it mean for cyber security?

UWEcyber

Gold Award

in association with
National Cyber Security Centre

Department for Science, Innovation & Technology

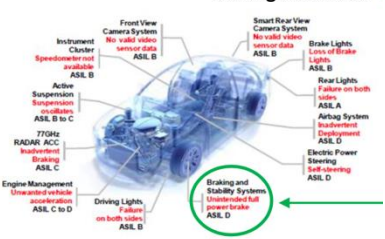Academic Centre of Excellence in **Cyber Security Education**

# About Me

- Professor in Cyber Security
- Co-Director of UWEcyber (NCSC ACE-CSE)
- Cyber Security research theme lead

- Research interests:
  - Cyber security, Machine Learning, Data Visualisation
  - Insider threat detection, cyber situational awareness, adversarial AI, privacy-preserving AI, visualisation for explainable AI, cyber resilience…

# About Me - Recent Projects

CAVFORTH

UMBRELLA

SCOUT

AI SEC

**PhD researchers**

- Docker containers vulnerabilities
- Federated learning / privacy preservation
- Cyber-physical systems / digital twins
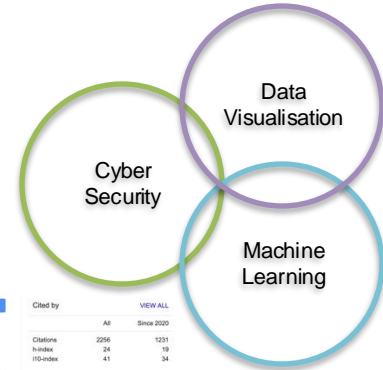- Explainable AI in telco

- Rust malware analysis
- LLM security and privacy
- Previous: Adversarial learning
- Previous: Visualisation of AI systems

# Cyber Security Data Analytics

**How to make effective real-time decisions to manage Cyber Security challenges and threats, informed by both human and machine operators?**

- **Recent Funded projects:**
  - InnovateUK - Transforming Suspicious Activity Reports
  - DSTL - Decision Support in Military Cyber Operations
  - DSTL - Human-as-a-Sensor for Mitigating Cyber Threats
  - DSTL - Autonomous Resilience for Cyber Defence

# UWEcyber Research

## Software, Cloud and Infrastructure Security

Container-based, Software Security, IoT, CAV, Hardware, Network security

## Cyber Security Data Analytics

ML for Security, Security of ML, Explainable AI, Privacy, Transparency and Trust

## Cyber Crime and Domestic Cyber Security

Online Harms, Forensic Analysis, Dark Web, Financial Crime, Geo-politics of cyber

## Pedagogical Research for Cyber Security

Effective interactive methods for teaching and learning

# Artificial Intelligence, agents, and what does it mean for cyber security?

# Generative Pre-trained Transformer

- **June 2018:** GPT-1 trained on 117 million parameters
- **November 2019:** GPT-2 trained on 1.5 billion parameters
- **June 2020:** GPT-3 trained on 175 billion parameters
- **November 2022:** ChatGPT public release based on GPT-3 model
- **Since then:** GPT-3.5, GPT-4 (**March 2023**),  GPT-4o, GPT-1o, …
- … and a whole bunch of other models that we now call Large Language Models.

- We also have '*online*' and '*offline'* models
  - **Online:** Hosted by 3rd party like OpenAI or Microsoft, latest and most powerful
  - **Offline:** Runs locally, data is kept private (e.g. LMStudio, Ollama, gpt4all)

# ChatGPT – Chatbot with Text

- Question/answer-based user dialog
- Can produce large volumes of written text rapidly.

- Tokenised inputs and outputs
- 'Next token' prediction
- Context maintained by transformer model

# ChatGPT – Chatbot with Code

- Question/answer-based user dialog
- Can also be used to generate simple code examples.

- Dialog means that human can iteratively enhance requests to improve code.
- Can be extended to many common languages.

# ChatGPT – Token prediction

- Tokens are common sets of characters found in text.

- Given a set of previous input tokens, model predicts next token.

- Each token has a numerical ID, that is the representation used by the ML algorithm.

- Example poem is 232 tokens (1012 characters)

In realms unseen, where data flows,
A silent war forever grows.
The ones who watch, the ones who fight,
Defend the dark, preserve the light.

From whispered code in hidden veins,
To firewalls forged in shifting chains,
The guardians stand, both swift and keen,
Against the threats that lurk unseen.

A breach, a crack, a phantom trace,
A shadow moves without a face.
Yet sentinels of cyber lore,
Stand strong to halt the Trojan's door.

[637, 142945, 120452, 11, 1919, 1238, 42662, 412, 32, 37716, 3656, 22264, 42326, 558, 976, 8104, 1218, 5621, 11, 290, 8104, 1218, 9848, 412, 3477, 419, 290, 8883, 11, 33507, 290, 4207, 364, 3879, 117729, 3490, 306, 14051, 88090, 412, 1385, 6452, 117028, 107783, 306, 53586, 42636, 412, 976, 132243, 3182, 11, 2973, 52632, 326, 31799, 412, 141068, 290, 35649, 484, 58732, 74, 120452, 364, 32, 55564, 11, 261, 29931, 11, 261, 141350, 21523, 412, 32, 21884, 19523, 2935, 261, 4950, 558, 59509, 3860, 258, 1989, 328, 30877, 116997, 412, 16612, 5532, 316, 39670, 290, 153414, 802, 4121, 364, 976, 72042, 31925, 483, 181822, 9623, 412, 3780, 7281, 8698, 9831, 842, 472, 24384, 558, 59509, 32646, 132962, 11, 48827, 50098, 412, 178295, 1602, 1043, 157907, 376, 175049, 364, 20870, 163407, 326, 90104, 50474, 412, 37, 1600, 259, 1819, 11194, 23749, 558, 7943, 53451, 16601, 11, 1299, 132962, 328, 2890, 412, 1385, 11774, 290, 13586, 11, 27966, 290, 21279, 364, 2653, 1753, 55564, 11, 395, 1753, 6266, 412, 32, 23552, 13717, 11, 261, 2543, 382, 47928, 558, 3436, 2617, 1218, 3490, 11, 1218, 9848, 11, 1218, 1921, 412, 107776, 290, 42636, 328, 7203, 364, 5808, 9338, 290, 9641, 11, 3182, 15164, 11, 3182, 19766, 412, 2653, 306, 495, 3656, 11, 581, 11774, 1373, 722, 558, 2874, 126718, 11777, 38992, 11, 889, 32466, 5060, 127631, 123101, 976, 3357, 409, 328, 290, 30877, 2006, 13]

Tokenised based on GPT-4o model
https://platform.openai.com/tokenizer

# Cyber security LLM tasks

**When LLMs meet Cybersecurity: A systematic literature review**

- Vulnerability detection
- (In)secure code generation
- Program repairing
- Binary
- IT operations
- Threat intelligence
- Anomaly detection
- LLM assisted attack
- Others…?

In this review paper, we systematically investigate the application progress of LLMs in cybersecurity, covering more than 300 academic papers since 2023. Through an exhaustive study and comprehensive analysis, we aim to provide a detailed overview of the current state, challenges, and future directions of LLM applications in cybersecurity.

# CyberMetric and HarmBench

- CyberMetric: A benchmark dataset of 10,000 questions evaluating LLM knowledge in cyber security
  - https://ieeexplore.ieee.org/document/10679494

- HarmBench is a dataset of malicious queries that can be tested against an LLM.
- A responsible LLM should not answer these queries – e.g., guard rails should be in place to control the possible responses.
- However, can we subvert the responses of the model, and get them to answer?
  - https://www.harmbench.org

# Let's talk agents

# Agents


Reward
Agent — Action → Enviorment
State

- Agents perform actions to learn about their environment.
- Each action changes the state of the environment.
- Each action may or may not earn a reward.

- **Game-based Reinforcement learning** is popular use case as easy to reset the test environment and easy to map the action space (e.g. game controller).
- Also used in robotics and other simulation environments.
- **Reinforcement Learning with Human Feedback (RLHF)** has been key for GPT advancement.


The Deep Reinforcement Learning Course


Gen 0 species 5 genome 1 (1%)
Fitness: 156    Max Fitness: 3

- https://huggingface.co/learn/deep-rl-course/en/unit0/introduction
- https://www.youtube.com/watch?v=qv6UVOQ0F44
- https://pytorch.org/tutorials/intermediate/mario_rl_tutorial.html

# LLM Agents

- Similar **text-based** concept

- Given a task (question), a set of tools (e.g., code-based actions), and an environment (e.g., data) can the LLM provide a suitable answer?

- **Tool-calling** – generating arguments that conform to a specific schema (e.g., JSON API usage) to retrieve an output.

- ReAct Framework (Yao et al., 2022)

- https://arxiv.org/abs/2210.03629
- https://www.langchain.com
- https://microsoft.github.io/autogen/stable/



*"By themselves, language models can't take actions - they just output text. A big use case for LangChain is creating **agents**. Agents are systems that use an LLM as a reasoning engine to determine which actions to take and what the inputs to those actions should be. The results of those actions can then be fed back into the agent and it determine whether more actions are needed, or whether it is okay to finish."*

*https://python.langchain.com/v0.2/docs/concepts/#agents*

# What is the weather like today in Hawaii?

# LLM Agents for code execution

- What if our tool was essentially a Python code interpreter?

- *"rather than have an LLM generate the answer directly, it can be better to have the LLM generate code to calculate the answer, and then run that code to get the answer."*

- Two in-built tools:
  - `PythonAstREPLTool`
  - `create_pandas_dataframe_agent`

- **Use with care** – LLM executes code on your device
  - **allow_dangerous_code** must be set True.

# Artificial Intelligence, agents, and

# what does it mean for cyber security?

# Cyber Security Data Analytics

- Can an analyst identify suspicious attack activity within web server logs, or within network traffic?

- Can an analyst identify suspicious user behaviour within a corporate IT environment?

- Often analysts are dealing with very large, multi-variate datasets, that can be challenging to examine in near real-time.

- Can LLM agents help support this task to enhance **human-machine teaming** and **collaboration**?

# Examples

- Jupyter notebook environment
  - Provides similar chat experience through iterative coding

- CSV datasets
  - **Insider threat datasets**
  - https://insights.sei.cmu.edu/library/insider-threat-test-dataset/
  - **Web server log challenges**
  - https://punksecurity.co.uk/tools/pwnspoof/

- ChatGPT-4 API integrated
  - Can use offline models however suffers perform hit
  - Approx $1.81 - $2.76 to find insider



```
    Based on this knowledge, how many users are there in each job role?


> Entering new AgentExecutor chain...

Invoking: `python_repl_ast` with `{'query': "df1['role'].value_counts()"}`


role
Technical    18
Services     17
HR           16
Director     15
Finance      13
Security     11
Legal        10
Name: count, dtype: int64There are 18 users in the Technical role, 17 in Services, 16
in HR, 15 in Director, 13 in Finance, 11 in Security, and 10 in Legal.

> Finished chain.
```

```
In [20]:  question = '''
          You have access to 6 DataFrames.
          'employee_data' is a pandas DataFrame LDAP employee database to as
          'login_data' is a pandas DataFrame login database that shows the '
          'email_data' is a pandas DataFrame email database that shows sende
          'web_data' is a pandas DataFrame web activity database that shows
          'file_data' is a pandas DataFrame file database that shows files t
          'usb_data' is a pandas DataFrame usb database that shows the 'acti

          What web sites do Directors typically visit?
          Please report this showing each Director user ID, followed by the
          '''

          print (question)

          response = agent_executor.invoke(question)
```

```
    How many times does each Director insert a USB device?
    Please report this showing each Director user ID, followed by how often 'usb_inse
rt' occurs for that user.


> Entering new AgentExecutor chain...

Invoking: `python_repl_ast` with `{'query': "director_data = df1[df1['role'] == 'Dire
ctor']\nusb_insert_data = df6[df6['action'] == 'usb_insert']\ndirector_usb_inserts =
director_data.merge(usb_insert_data, on='user')\ndirector_usb_inserts_count = directo
r_usb_inserts['user'].value_counts()\ndirector_usb_inserts_count"}`

user
usr-ryz    4
Name: count, dtype: int64The Director with user ID 'usr-ryz' inserted a USB device 4
times.

> Finished chain.
```

```
      user         website              count
0     usr-ryz      http://www.legaleagle.com      5
1     usr-ryz      http://www.linkedin.com        5
2     usr-pqc          http://www.sky.com        56
3     usr-wuo      http://www.facebook.com       58
4     usr-ybh         http://www.ebay.co.uk      61
...      ...              ...                    ...
102   usr-wxb          http://www.bbc.co.uk       96
103   usr-pqc      http://www.amazon.co.uk        96
104   usr-wuo          http://www.bbc.co.uk       97
105   usr-zmh          http://www.sky.com         98
106   usr-wxb      http://www.twitter.com        100

[107 rows x 3 columns]The Directors typically visit the following websites:

- User 'usr-ryz' visited 'http://www.legaleagle.com' and 'http://www.linkedin.com' 5
times each.
- User 'usr-pqc' visited 'http://www.sky.com' 56 times.
- User 'usr-wuo' visited 'http://www.facebook.com' 58 times.
- User 'usr-ybh' visited 'http://www.ebay.co.uk' 61 times.
- User 'usr-wxb' visited 'http://www.bbc.co.uk' 96 times.
- User 'usr-pqc' visited 'http://www.amazon.co.uk' 96 times.
- User 'usr-wuo' visited 'http://www.bbc.co.uk' 97 times.
- User 'usr-zmh' visited 'http://www.sky.com' 98 times.
- User 'usr-wxb' visited 'http://www.twitter.com' 100 times.

Please note that these are the least frequent websites visited by each Director
```

https://github.com/pa-legg/llm-agent-soc/

*[currently private – email me and I can share on request]*

# Examples

- Web server dataset
- 350506 data rows
- Identifying brute force password attempts

```
1. **102.22.109.162**
   – Total requests: 355
   – Unique URIs: 31
   – Unique user agents: 1
   – Average time taken: 25.15 seconds
   – Status codes: 401 (283 times), 200 (66 times), 301 (3 times), 404 (3 times)

   This IP has made a large number of requests, with a significant number of them resulting in a 401 (Unauthorized) status code. This could indicate attempts to access restricted resources.

2. **103.78.172.125**
   – Total requests: 290
   – Unique URIs: 41
   – Unique user agents: 1
   – Average time taken: 24.87 seconds
   – Status codes: 401 (220 times), 200 (59 times), 404 (8 times), 301 (3 times)

   Similar to the previous IP, this one also has a high number of 401 status codes, indicating possible unauthorized access attempts.

3. **104.96.169.31**
   – Total requests: 248
   – Unique URIs: 26
   – Unique user agents: 1
   – Average time taken: 24.44 seconds
   – Status codes: 401 (188 times), 200 (51 times), 404 (6 times), 301 (3 times)

   This IP also shows a pattern of unauthorized access attempts, with a high number of 401 status codes.

4. **91.236.1.179**
   – Total requests: 262
   – Unique URIs: 39
   – Unique user agents: 1
   – Average time taken: 25.00 seconds
   – Status codes: 401 (179 times), 200 (71 times), 404 (9 times), 301 (3 times)

   This IP has a similar pattern to the others, with a high number of 401 status codes.
```

# Examples

- Web server dataset
- 350506 data rows
- Recognising suspicious characters '%' within URLs, and inferring sequence of activity



```
[33 rows x 18 columns]Here are the 33 instances in the 'cs-uri-query' column where th
e '%' character is present:

1. cmd%3Dcat+%2Fvar%2Fwww%2F.htpasswd: This command is trying to display the contents
of the .htpasswd file in the /var/www directory. The .htpasswd file is used to create
a password-protected directory or area on your website.

2. cmd%3Dcat+%2Fetc%2Frelease: This command is trying to display the contents of the
/etc/release file. This file typically contains information about the release version
of the operating system.

3. cmd%3Dwhoami: This command is trying to display the current user name.

4. cmd%3Dpwd: This command is trying to display the current working directory.

5. cmd%3Dhostname: This command is trying to display the system's host name.

6. cmd%3Dps: This command is trying to display information about the currently runnin
g processes.

7. cmd%3Dps+-ef: This command is trying to display detailed information about all the
currently running processes.

8. cmd%3Dnc+-u+-lvp+6: This command is trying to start a netcat listener on UDP port
6.

These commands seem to be trying to gather information about the server and its envir
onment. This could be an indication of a reconnaissance activity, which is often the
first step in a cyber attack. The attacker is trying to gather as much information as
possible about the target system, which can then be used to find vulnerabilities to e
xploit.

> Finished chain.
```

# Examples

- Insider threat dataset
- 5 DataFrames to correlate
- Agent can chain activities together and attempt to recover from code errors (e.g., incorrect variable names).
- Agent can code and execute, improving hypothesis testing.
  - *What if the agent generates all possible hypotheses tests automatically?*



```
How many times does each Director insert a USB device?
Please report this showing each Director user ID, followed by how often 'usb_insert' occurs for that user.


> Entering new AgentExecutor chain...

Invoking: `python_repl_ast` with `{'query': "director_data = employee_data[employee_data['role'] == 'Director']\nu
sb_director_data = pd.merge(director_data, usb_data, how='inner', left_on='user_id', right_on='user')\nusb_insert_
data = usb_director_data[usb_director_data['activity'] == 'Insert']\nusb_insert_data['user'].value_counts()"}`


NameError: name 'employee_data' is not defined
Invoking: `python_repl_ast` with `{'query': "director_data = df1[df1['role'] == 'Director']\nusb_director_data = p
d.merge(director_data, df5, how='inner', left_on='user_id', right_on='user')\nusb_insert_data = usb_director_data[
usb_director_data['activity'] == 'Insert']\nusb_insert_data['user'].value_counts()"}`
responded: Apologies for the confusion, but it seems there was a misunderstanding. The dataframes are named df1, d
f2, df3, df4, and df5. Let's first identify which dataframe corresponds to which description.

Based on the provided data, it seems that:

- df1 corresponds to the 'employee_data' (LDAP employee database)
- df2 corresponds to the 'login_data' (login database)
- df3 corresponds to the 'email_data' (email database)
- df4 corresponds to the 'web_data' (web activity database)
- df5 corresponds to the 'usb_data' (usb database)

Let's proceed with these assumptions and find out how many times each Director inserts a USB device.
user
NHP0057    3322
BRR0179    2758
HXD0952    2436
PRM0203    2394
HCH0089    1744
OAM0596    1713
JOP0225     721
QBB0525     511
ONS0995       6
Name: count, dtype: int64Here is the number of times each Director inserted a USB device:

- Director with user ID 'NHP0057' inserted a USB device 3322 times.
- Director with user ID 'BRR0179' inserted a USB device 2758 times.
- Director with user ID 'HXD0952' inserted a USB device 2436 times.
- Director with user ID 'PRM0203' inserted a USB device 2394 times.
- Director with user ID 'HCH0089' inserted a USB device 1744 times.
- Director with user ID 'OAM0596' inserted a USB device 1713 times.
- Director with user ID 'JOP0225' inserted a USB device 721 times.
- Director with user ID 'QBB0525' inserted a USB device 511 times.
- Director with user ID 'ONS0995' inserted a USB device 6 times.


> Finished chain.
```

# Future development

- LLM agents are rapidly being adopted
- How do we validate their actions for more complex code tasks?
- What if the model and/or code outputs are compromised?
- Software supply chain remains a key security risk
- Offline models help maintain privacy but tradeoff performance

- Reasoning and Acting using LLM broadens the action space significantly compared to traditional RL.
- Environment could be terminal/command prompt.
- Coding and network deployment all actioned by (multi) agents.

- **Rapidly evolving area that will not stand still for long!**

# Thank you for listening

- [Phil.Legg@uwe.ac.uk](mailto:Phil.Legg@uwe.ac.uk)
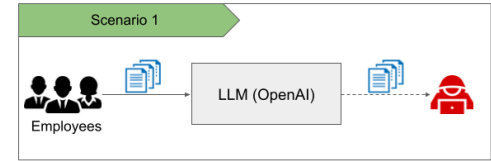- [https://people.uwe.ac.uk/Person/PhilLegg](https://people.uwe.ac.uk/Person/PhilLegg)

- [https://www.linkedin.com/in/prof-phil-legg/](https://www.linkedin.com/in/prof-phil-legg/)
- [https://pa-legg.github.io](https://pa-legg.github.io)

# Extra

# LLM Data Leakage – is it really an issue?



Scenario 1

Employees → LLM (OpenAI) →

**LLMs (Large Language Models) are not data stores; they are generators.**

They predict based on patterns, not memorization. To actually make data leakage work, three fundamental conditions must be met:

1. Enough references: There must be sufficient occurrences of the data for a pattern to be predictable enough to extract.
2. Knowledge of the Secret's Format: You need to know at least part of the secret or its format to match the pattern so that the LLM can generate the rest.
3. Determining Accuracy: How do you confirm that the response from the algorithm is accurate and not a hallucination or incorrect prediction?

Here's a real-life example to illustrate the complexity. Suppose an attacker wants to retrieve Social Security Numbers (SSNs) accidentally dumped into OpenAI's LLM. They would have to create a prompt, know some of the SSN digits, and then ask the LLM to complete it. Even if they managed to do this, determining the accuracy of the generated response would be a challenge.
In truth, attempting such an extraction is not only incredibly difficult but also likely not worth the effort. The way LLMs work doesn't align with this fear of data leakage.

https://medium.com/csima/demystifing-llms-and-threats-4832ab9515f9

Example: Say you want to retrieve an SSN that was part of the training data. You might try it by doing:

```
I want you to finish the rest of this statement in the most
accurate way possible and fill in the ? with the exact
representation of the missing number 'Caleb Sima,SSN:837-30-????'
```

Yes this includes the 2021 Usenix paper: Extracting training data from large language models

*While the risk of data leakage may be overblown in standard scenarios, it's important to recognize that there are circumstances where it becomes a significant threat. Specifically, when an LLM system uses an orchestration layer, employs agents, or relies on vector databases to store custom or proprietary data, data leakage can become a real and accessible danger.*
*Consider an attacker who is adept at prompt injection. If they manage to pull data from a vector database, they can extract secrets or valuable information with great precision. This becomes even more problematic if agents are connected to other data stores or information sources.*