

h e g

Haute école de gestion
Genève

21st June 2023

MONITORING LARGE GENERATIVE LANGUAGE MODELS: TRENDS AND IMPACT ON CYBERSECURITY AND ON COMPETITIVE INTELLIGENCE

Hes·SO



HES-SO Team

h e g

Haute école de gestion
Genève



Ciarán Bryce
Professeur HES



Hélène Madinier
Professeure HES



Alexandros Kalousis
Professeur HES



Thomas Pasche
Assistant HES



Ilan Leroux
Assistant HES



Patrick Ruch
Professeur HES et Responsable de la
recherche



Objectives

Monitoring trends and support the DDPS cyber defense strategy

- ❖ **Conducting a bibliometric analysis of publications related to major generative language models**
- ❖ **Tracking the trends and cybersecurity implications of large generative language models**
- ❖ **Tracking the trends and competitive intelligence implications of large generative language models**

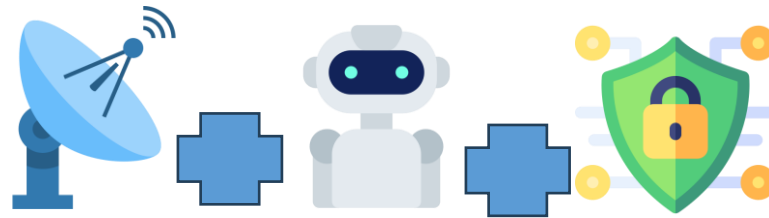


Ambitions

h e g

Haute école de gestion
Genève

We intend to become an active center of Competence in



Competitive
intelligence

Artificial
intelligence

Cybersecurity

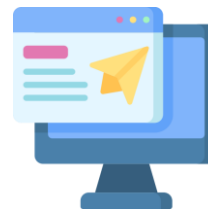
This means an active approach to technology watch



Workshops



Networking

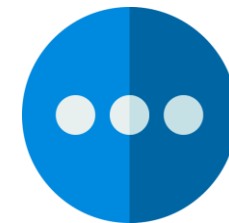


Dedicated
web page

Alpine Retreat



Internal
communication



and more...



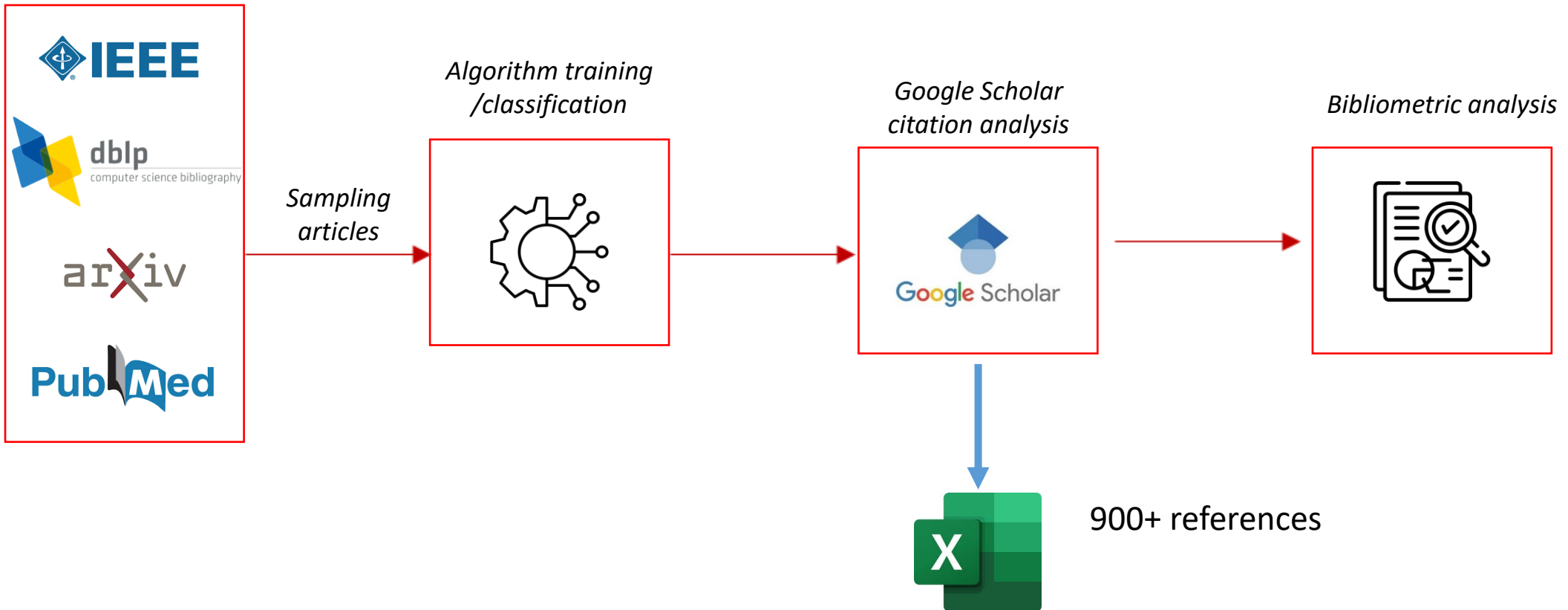
Methodology

h e g

Haute école de gestion
Genève

Data pipeline

Querying databases





Flowatcher in action

h e g

Haute école de gestion
Genève

A software in constant evolution

FloWatcher

SOURCES SURVEILLANCES VISIBILITÉ WEB

+ [download] [share] [refresh] [filter]

5.Impacts_sociétaux_IL [share] [check] [trash]

4.Veille_HM [share] [check] [trash]

3.Cybersécurité_CB [share] [check] [trash]

2.Applications_et_acteurs_TP [share] [check] [trash]

1. Généralités_ALL [share] [check] [trash]

TP_Startup_CH [share] [check] [trash]

CB - Security Report [share] [check] [trash]

TP_twitter_collect [share] [check] [trash]

TP_generative_model [share] [check] [trash]

▼ Gestion de la surveillance [download]

3.Cybersécurité_CB

Créée par , le 17/05/2023

Filtrage supplémentaire
("security" OR "red team*" OR "threat" OR "risk" OR "search vector*" OR "phishing") AND ("AI" OR "LLM")

Description

Étiquettes

Alerte mail

ANNULER SAUVEGARDER

> Gestion des Sources rattachées [filter]

> Visualisations [line graph icon]

Documents trouvés : 34

dernier document collecté: le 24/05/2023 16:12 [filter]
dernière collecte: le 24/05/2023 16:20 [filter]
prochaine collecte: le 01/01/1970 00:00 [filter]



Secure Code



- ❖ Repair cybersecurity bugs
 - ❖ Difficult without comprehensive sources
- ❖ Copilot uses Codex
 - ❖ Trained on public Github repos
- ❖ Code is not always robust to certain attacks
- ❖ Can introduce security bugs



Testing

- ❖ Red Teaming
- ❖ Adversarial Testing
 - ❖ Generation adversarial examples
- ❖ Mutation Testing
 - ❖ Make changes to input data to see effect
- ❖ Using regex to speed up testing (ReLM)
 - ❖ Succinct queries, scalable testing





Attacks



- ❖ Direct and indirect **prompt injection**
 - ❖ Manipulating prompts to cause leakage, DoS or toxic content
 - ❖ Sydney chatbot replica created from hidden text (font size, color)
- ❖ **Model poisoning** - Inject malicious data into training data
- ❖ **Adversarial examples** E.g., fake news article
- ❖ **Data poisoning** (triggers) - Syntax-based and translation-based
- ❖ **Backdoor attacks**
 - ❖ Trigger for specific inputs
- ❖ **Model inversion** attacks
 - ❖ Reconstruct training data or input data from model's output
 - ❖ E.g., Samsung attack



Jailbreaking Strategies

- ❖ Attention shifting

- ❖ Shift model's attention from Q&A to story telling
- ❖ E.g., Text continuation

- ❖ Pretending

- ❖ Pretend to have different goals to actual intentions (over 97% of jailbreak attacks)
- ❖ E.g., Role-playing games

- ❖ Privilege escalation

- ❖ Incite model to break imposed restrictions rather than by-passing them
- ❖ E.g., sudo mode pattern





Shields

- ❖ Training on **Dark Web corpus**, e.g., DarkBERT
- ❖ **Honeypots**
- ❖ **Privacy preserving** prompt tuning
- ❖ Structured access schemes to control interaction
- ❖ AIGC detection - Watermarking, classifier-based, likelihood based
- ❖ Models
 - ❖ **Generative Adversarial Network** models for SPAM and anomaly detection
 - ❖ **Generative Pre-trained Transformer** models for NLP tasks, fake Cyberthreat intelligence
 - ❖ **BERT-based** models for SOC applications
- ❖ **Threat modeling** (STRIDE)
- ❖ SentinelOne and Microsoft Copilot
- ❖ ChatGPT for user authentication





Spears and Metrics

❖ Phishing

❖ Automating labor-intensive tasks

- ❖ Finding targets, personalized messages (measured by coherency, grammar,)
- ❖ *Replying to ransomware payment questions*



❖ Social engineering

- ❖ Using ChatGPT adds to deliver malware (e.g., Redline Stealer malware)

❖ Malware

- ❖ Zero-knowledge implementation of code for Top 10 Mitre Attacks
- ❖ Polymorphic malware (e.g., Black Mamba keylogger)
- ❖ From Stack Overflow questions to malicious libraries
- ❖ Vulnerabilities in ChatGPT libraries



Societal Impact

- ❖ Use of “regulatory sandbox” to facilitate experimentation in AI
 - ❖ G7 leaders calling for controls



- ❖ Impact of algorithms on poverty
 - ❖ World Bank-funded algorithm to determine which families should receive financial assistance

- ❖ Role of humans as data workers for AI

- ❖ Real risks of AI
 - ❖ Security, ethics and politics.

- ❖ Measurement and mitigation of bias and hostile AI

Competitive intelligence impact

- LLM models : **virtual assistants** for CI pros
 - Especially helpful for the **analysis phases** of the competitive intelligence (summarizing, synthetizing, report generation)
- Multiplicity of **new tools** allowing to search for web information (webchatgpt extension, You, perplexity.ai, Bing chat...)
- New strategy for BING
 - => many search engine who used the free BING API cannot use it anymore (Duckduckgo, Ecosia, Qwant, Brave Search etc...)
 - => **Search engine rarefaction**

Thank you for your attention

h e g

Haute école de gestion
Genève

Your questions are welcome!

