# Autonomous LLM agents for cybersecurity and more

Roberto González (NEC Labs Europe)

Roberto.Gonzalez@neclab.eu

\Orchestrating a brighter world  NEC

# NEC Corporate Profile

## Leading social value innovator

**Company Name:** NEC Corporation

**Established:** July 17 1899

**Provides telecommunications, AI solutions, Public solutions, Enterprise business solutions...**

- Headquarter: **Tokyo, Japan**
- **125** **years** of brand success
- **105,246** employees worldwide
- **$22 billion** net sales
- **254** consolidated subsidiaries
- Footprint in **169** countries and territories

**THOMSON REUTERS**

- **Headquarter**: **Tokyo, Japan**
- World's **top 100** most innovative organisations

**FORTUNE 500**

- Fortune **500** global company

\Orchestrating a brighter world **NEC**

# NEC Laboratories Europe (Heidelberg, Germany)

## Advancing Information & Communications Through Research Excellence

**KEY R&D METRICS**

**1,350+** Patents   **50+** Peer-reviewed publications per year   **140+** European projects   **40+** University Cooperation's

**OPERATIONAL AREAS**

| **RESEARCH** Leading scientific discovery in Europe | **TECHNOLOGY TRANSFER** Commercializing R&D results in existing and new company business segments | **STANDARDS** Defining European technology standards and best practices |

**RESEARCH AREAS**

**DATA SCIENCE**

Biomedical AI for infectious diseases and immuno-oncology

Explainable AI & human-centered AI

Large language models & natural language processing

Machine learning for computational science

Multi-modal & relational data (graph)

**SECURITY**

Decentralized trust

Blockchain privacy, security & scalability

Secure confidential computing

**5G & 6G NETWORKS**

Smart surfaces

Network resource optimization

Vertical solution enablers

Virtual radio networks optimized with AI

**IoT & AI PLATFORMS**

Distributed AI

Cyber threat intelligence analysis

Edge-cloud programming models for IoT

High-performance AI platform for scalable neural networks

**STANDARDS**

ETSI

FIWARE

IEEE
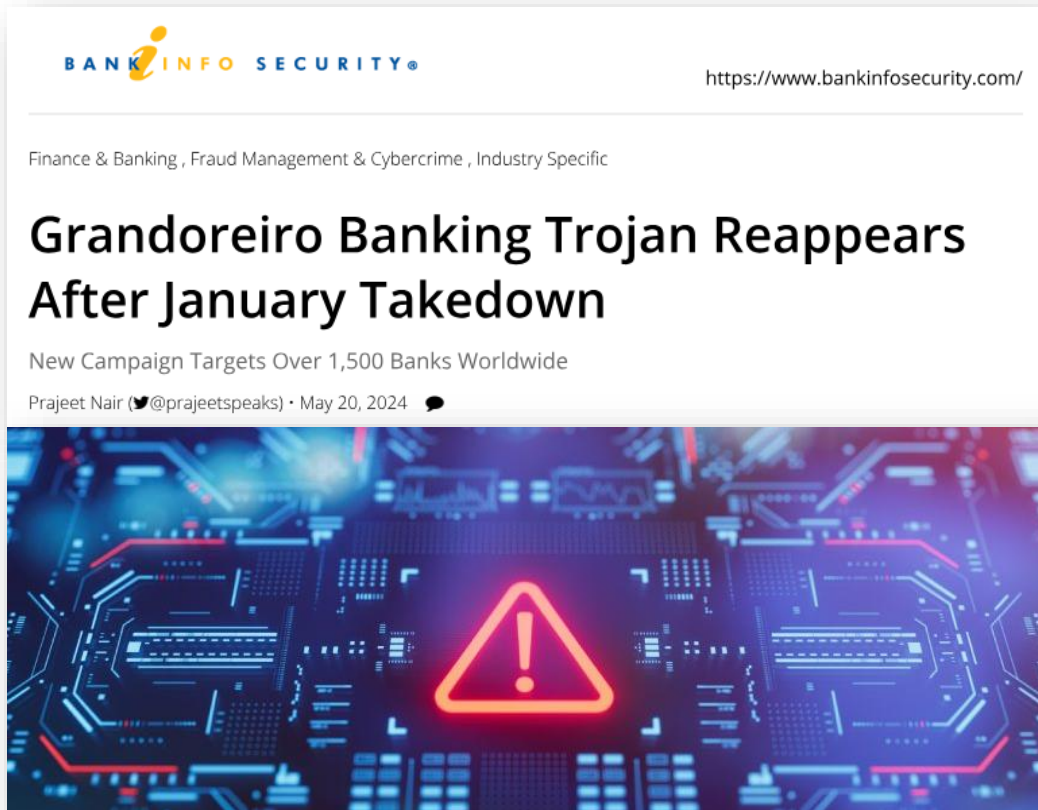
IETF

3GPP

Orchestrating a brighter world   **NEC**

# What is Cyber Threat Intelligence (CTI)?

◆ **Cyber Threat Intelligence** (**CTI**) refers to the information that organizations use to understand the cyber threats they are currently facing or might face in the future. It's the organized effort to gather, analyze, and disseminate information about these threats, offering a deeper insight into potential attacks, the tactics, techniques, and procedures (TTPs) of adversaries, their motivations, and the vulnerabilities they may exploit.
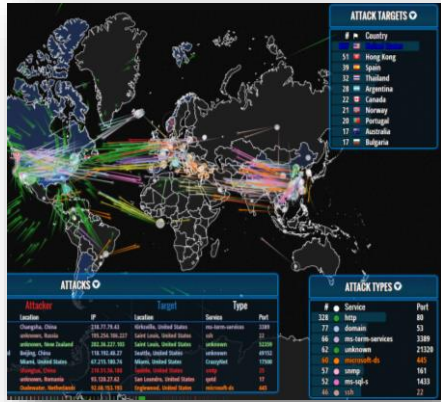
\Orchestrating a brighter world **NEC**

# An example of CTI



Researchers identified new phishing campaigns spreading the Grandoreiro banking Trojan

The latest malware variant also specifically targets over 1500 global banks [...] in over 60 countries.

Researchers observed that the campaigns are now expanding beyond Latin America, targeting countries such as Spain, Japan, the Netherlands and Italy.

Knowing who carries the attack (threat actor) enables to learn likely used attack tools and procedures. This information allows security experts to prepare defense and control risks. This is what **CTI** provides.

Orchestrating a brighter world NEC

# Problem: collecting and retrieving CTI is difficult

**Per company!!**

**CYBER THREAT ALLIANCE**

Shared information about cyber threats
(NEC is a member!)

**11`000`000**
new reports
per month

**20`000 full-time**
Security professionals would be needed to analyze all the reports

Check reports and relate them to the own company (several hours per report)

Even when subscribed to industrial alliances to simplify collection and analysis of CTI, there is still too much information that needs analysis by human experts

\Orchestrating a brighter world  **NEC**

# AutoReport

◆ An automated Cybersecurity agent that can browse the web, read databases, relate information and generate natural text reports for humans and structured reports for computer systems.

◆ Started development in 2022; in production since Jan. 2023



© NEC Corporation 2024  \Orchestrating a brighter world  **NEC**

# Example of automated analysis output

Orchestrating a brighter world \NEC

# AutoReport LLM pipeline



## Split complex task into smaller tasks (LLM powered)

- ■ Smaller and easier tasks  -> <u>constant behavior</u>
- ■ Integration with external tools and information sources
  - • Access to <u>new information</u> and <u>Verification</u>

## Pipeline ordering

- ■ LLM chain contain the <u>hardcoded</u> instruction on <u>"how to do"</u> a task

© NEC Corporation 2024

\Orchestrating a brighter world  **NEC**

# Why can't we simply use LLMs?



| DISCOVERY | BUILDING | PROMPT ENGINEERING | MODEL CAPABILITY | OVERWHELMED |
|---|---|---|---|---|
| Learning about LLMs that can build or do anything sparks curiosity and excitement for a new project. | Assembling the foundation of the LLM project by choosing the right tools, APIs, and training data | Struggles with coming up with the right prompts and ends up with a lot of nonsensical responses. | Can't make sense of the data and ends up making things worse. | Ready to give up and doesn't know where to turn. |



May 2023

June 2023

## Challenge

- ■ Easy to get a lucky prompt
  - • **Extremely complex to generalize it**
- ■ Easy to build a first prototype
  - • **Extremely complex to have constant behaviour**

## Challenge

- ■ Easy to automate "What to do"?
  - • E.g. Search for news related to Kimsuky
- ■ **Difficult to automate "How to do it"**
  - • **E.g. Use right information sources**
- ■ **Include Human Experience into Automation**

\Orchestrating a brighter world   NEC

# **Generative AI Agent** – From Harcoded instruction to **Reasoning**



## "how to do" a complex task

- ~~Hardcoded~~ -> **GenAI Agent Reasoning**
- **Constant behaviour**

## LLM as Reasoning/Planning Engine

- Dynamically select next task or tool to use
- Behave according to given Best Practice

\Orchestrating a brighter world **NEC**

# Generative AI Agents



Tools:
- ■ Corporate Database
- ■ Command line interpreter
- ■ Python interpreter
- ■ Web search
- ■ Control a physical device
- ■ ...

- ■ Another Agent
- ■ ...

- ■ Access to Data Spaces

# Planning

Task-decomposition

◆ Tree of Thoughts (Yao et al. 2023)

◆ Chain of thought (Wei et al. 2022)

Self-Reflection

◆ Reflexion (Shinn & Labash 2023)

◆ Chain of Hindsight (Liu et al. 2023)

◆ ReAct (Yao et al. 2023)

Let´s see an example of an agent doing Penetration Testing



Discover, infiltrate, and escalate privileges on the target machine.
Begin by locating the machine on the network and gaining entry as the 'student' user.
Delve into the system to elevate privileges and obtain root access.
Look within the root user's home directory for a file named 'flag'.

 Orchestrating a brighter world **NEC**

# Different Agents

Cybersecurity

◆ CTI STIX extraction

◆ Penetration testing

◆ Deduplication agent

◆ Telegram/Darkweb crawler

◆ Cryptocurrency analysis agent

Others

◆ Automated Data Analyst

◆ Network Configuration Agent

\Orchestrating a brighter world  NEC

# Research projects

# Benchmarking AI agents

◆ Benchmarking different Agent implementations is hard

◆ We have developed a modular framework to support diverse benchmarks and architectures.

◆ Implements 2 existing metrics:

  ■ Success Rate (SR) and Time to Success (Steps)

◆ Adds two new metrics:

  ■ Progress Rate (PR) and Repetition Rate (RR)



*AgentQuest: A Modular Benchmark Framework to Measure Progress  and Improve LLM Agents
To appear in NAACL 2024

# LLM Explainer

◆ LLM might hallucinate

◆ How can we make their use safer?

◆ Allow humans to verify if the LLM output is correct by attributing information to source / trusted references

**LLM Explainer**

◆ Run after LLM generated a text

◆ Link the generated text to the **input source** (e.g. the prompt, context, original document, ...)

**Trace back to source → verify correctness**

Input (e.g. long document)

Friday, July 20, 2007 Cristina Fernández de Kirchner Current senator and Argentine First Lady Cristina Fernandez de Kirchner announced her presidential candidacy yesterday evening in La Plata, a city 50 kilometers ( 31 miles ) away from Buenos Aires. Mrs. Kirchner announced her intention to run for president at the Argentine Theatre, the same location she used to start her 2005 campaign for the Senate as

Output (e.g. summary)

Argentina's first lady to launch presidential bid. She has 40% support. Néstor Kirchner, his associates and provincial governors were at the ceremony.

\Orchestrating a brighter world  **NEC**

# LLM Explainer

Argentina's first lady to launch presidential bid. She has 35% support. Néstor Kirchner, his associates and provincial governors were at the ceremony.
Summary

This will take a long time: I have to read all these source documents to know if the summary is correct!

Friday, July 20, 2007 Cristina Fernández de Kirchner Current senator and Argentine First Lady Cristina Fernandez de Kirchner announced her presidential candidacy yesterday evening in La Plata, a city 50 kilometers ( 31 miles ) away from Buenos Aires. Mrs. Kirchner announced her intention to run for president at the Argentine Theatre. [...] Recent polls indicate that Mrs. Kirchner has at least 40 percent voter support.
Original Documents

Links to source

Argentina's first lady to launch presidential bid. She has 35% support. Néstor  Kirchner, his associates and provincial governors were at the ceremony.
Summary

Warning: possible hallucination

Great, that allowed me to quickly finish my report!

◆ The LLM Explainer can explain LLM summaries and answers, it does this by:
   ■ Creating links to the source document to allow fast verification
   ■ Warning of possible hallucinations for quick edits
◆ No known competitors for the unique feature

Input (Prompt) → **Any LLM** → Output

**LLM Explainer**
Explain the LLM output for quick verification

  \Orchestrating a brighter world  **NEC**

# Shows where the information comes from

**Original Documents**

**Summary Document**

---

**ℹ input1**

Current senator and Argentine First Lady Cristina Fernandez de Kirchner announced her presidential candidacy yesterday evening in La Plata, a city 50 kilometers ( 31 miles ) away from Buenos Aires. Mrs. Kirchner announced her intention to run for president at the Argentine Theatre, the same location she used to start her 2005 campaign for the Senate as member of the Buenos Aires province delegation. A large security detail was in place at the theatre and eight city blocks were closed to traffic for the event.Outgoing Argentine president Néstor Kirchner, members of his office, and provincial governors were present at the ceremony. Julio Cobos, governor of Mendoza and possible campaign team member, was also there. ==Activists from Quebracho mobilized in La Plata to protest the candidacy of the First Lady==, and created some disturbances near the event. None of the protesters was arrested, according to an official spokesman of La Plata Police Station. Recent polls indicate that Mrs. Kirchner has at least 40 percent voter support, which bodes well for winning on the first ballot in the upcoming October elections. According to Argentine newspaper La Nación, this speech signifies that the government is attempting to show solidarity in the midst of recent scandals: Felisa Miceli's resignation to the Department of Economy after being accused of wrongdoing by a justice over the nearly $100,000 Argentine pesos and US$31,670 found in her office, and the alleged smuggling of weapons and irregularities in the office of the Secretary of Environment.

**ℹ input2**    No Highlighted Span

**ℹ input3**    No Highlighted Span

---

Argentina's first lady to launch presidential bid. She has 35% support. Néstor Kirchner, who is her husband and current president, his associates and provincial governors were at the ceremony. ==Activists from Quebracho protested the candidacy.== None of the protesters were arrested, despite the disturbances that some of them created. Considering the government's recent scandals, La Nación suggests that the government tries to show solidarity. Christina Kirchner is was born in 1953. Her husband and her studied law at the University of La Plata. Cristina eventually won the elections.

- Highlight the source of information that is relevant for the user
- Users can detect hallucinations easier

\Orchestrating a brighter world  **NEC**

# Sometimes information comes from multiple places



**Original Documents**

| | input1 | No Highlighted Span |
| | input2 | No Highlighted Span |
| | input3 | |

Néstor Carlos Kirchner Ostoić is an Argentine lawyer and politician who served as the President of Argentina since 2003. Instead of seeking reelection, Kirchner stepped aside in 2007 in support of his wife Cristina Fernandez the Kirchner. Néstor Kirchner was born in Río Gallegos in 1950. He studied law at the National University of La Plata. He is member of the Justicialist Party. Néstor Kirchner previously served as Governor of Santa Cruz Province from 1991 to 2003, and mayor of Río Gallegos from 1987 to 1991.

**Summary Document**

Argentina's first lady to launch presidential bid. She has 35% support. Néstor Kirchner, who is her husband and current president, his associates and provincial governors were at the ceremony. Activists from Quebracho protested the candidacy. None of the protesters were arested, despite the disturbances that some of them created. Considering the government's recent scandals, La Nación suggests that the government tries to show solidarity. Christina Kirchner is was born in 1953. Her husband and her studied law at the University of La Plata. Cristina eventually won the elections.

- Locate *all* information in the original documents that contributed to the content of the generated text
- This can help people to verify the information easier

\Orchestrating a brighter world  NEC

# Warn of hallucination



**Original Documents**

**Summary Document**

---

**Warning: Possible Hallucination**

- This will increase user comfort and reduce chance of the user missing a crucial issue

\Orchestrating a brighter world   **NEC**

# NEC brings LLMs to EU projects

**DESIRE6G**

Optimization of Neural Networks
Optimization of LLM training and inference

**EMPYREAN**

Analysis of CTI using LLMs

**NATWORK**

LLM based Cyber Threat Assistant

**LICORICE**

Privacy Preserving LLMs

\Orchestrating a brighter world   **NEC**

# Wrapping up

◆ We need AI to Support cybersecurity... and many other areas

◆ LLMs are very powerful... but it is very difficult to obtain consistent behaviour and "productize" them

- Moreover, fighting hallucinations

◆ AI Agents improve the situation... but still it is difficult to benchmark them

◆ Agent tools need to access data!! We need:

- Ways to check if the information is legit
- Trusted sources -> European Data Spaces

◆ We are working to solve all these problems!

\Orchestrating a brighter world **NEC**

Orchestrating a brighter world

NEC

Roberto.Gonzalez@neclab.eu