# Technical measures for safe use of generative AI

eROUN&COMPANY co.,ltd.

Doo Shik, Yoon

# Contents

# 1

**Overview**

# Paradigm shift

ChatGPT is here

------------------------------------------

Disruption of routines and
ways of working

It's like that time when the internet and PCs completely changed the way we do
things.

# Efficiency vs. Security

**FORBES > BUSINESS**

**BREAKING**

## Samsung Bans ChatGPT Among Employees After Sensitive Code Leak

**Siladitya Ray** Forbes Staff

*Covering breaking news and tech policy stories at Forbes.*

Follow

May 2, 2023, 07:17am EDT

## Locally developing large-scale AI is matter of digital sovereignty: expert

generative AI asia 2023
모두를 위한 AI

Yoo Young-joon, Wrtn Technologies co-founder and chief operating officer, speaks during the Generative AI Asia 2023 conference held in southern Seoul, Wednesday. [WRTN TECHNOLOGIES]

**Business efficiency**

**Security / Digital sovereignty**

# 2

## Security issues with generative AI

# The biggest threats to AI adoption

**Organization considers risk relevant**

| | |
|---|---|
| Inaccuracy | 56 |
| Cybersecurity | 53 |
| Intellectual-property infringement | 46 |
| Regulatory compliance | 45 |
| Explainability | 39 |
| Personal/individual privacy | 39 |
| Workforce/labor displacement | 34 |
| Equity and fairness | 31 |
| Organizational reputation | 29 |
| National security | 14 |
| Physical safety | 11 |
| Environmental impact | 11 |
| Political stability | 10 |
| None of the above | 1 |

**Organization working to mitigate risk**

| | |
|---|---|
| Inaccuracy | 32 |
| Cybersecurity | 38 |
| Intellectual-property infringement | 25 |
| Regulatory compliance | 28 |
| Explainability | 18 |
| Personal/individual privacy | 20 |
| Workforce/labor displacement | 13 |
| Equity and fairness | 16 |
| Organizational reputation | 16 |
| National security | 4 |
| Physical safety | 6 |
| Environmental impact | 5 |
| Political stability | 2 |
| None of the above | 8 |

| Inaccuracies | Cybersecurity | Intellectual property infringement |
|---|---|---|

**Failing to address potential risks in the age of AI**

# Risks and security issues when using generative AI

1 Data privacy and confidentiality

2 Third-party security issues

3 AI behavioral vulnerabilities

4 Legal issues

5 Evolution of threat actors (attackers)

6 Copyright issues

7 Generating insecure code

8 Bias and discrimination issues

9 Issues of trust and reputation

10 Software security vulnerabilities

11 Performance, availability, and cost issues

12 Ethics and regulatory issues

# Top 10 Security Risks of LLMs

1 Prompt Injection

2 Insecure Output Handling

3 Training Data Poisoning

4 Model Denial of Service

5 Supply Chain Vulnerabilities

6 Sensitive Information Disclosure

7 Insecure Plugin Design

8 Excessive Agency

9 Overreliance

10 Model Theft

# Regulating the use of AI in South Korea

**National Intelligence Service Releases
'Security Guidelines for Generative AI '  (2023.06.)**

**1** **Top Security Threats in Generative AI Technology Announced**

**2** **Guidelines for safe use of generative AI technologies**

- **Cautions for using the service**
- **Cautions for talking to services**
- **Cautions for using service plugins**
- **Cautions for using service extensions**
- **How to deal with AI model generation-based attacks**

> **Difficult for users to protect themselves**
>
> **Need for automated controls**

**3** **Generative AI-based informationization business construction plan and security measures**

It is difficult for an organization or institution to have systems and processes in place to meet AI privacy guidelines.
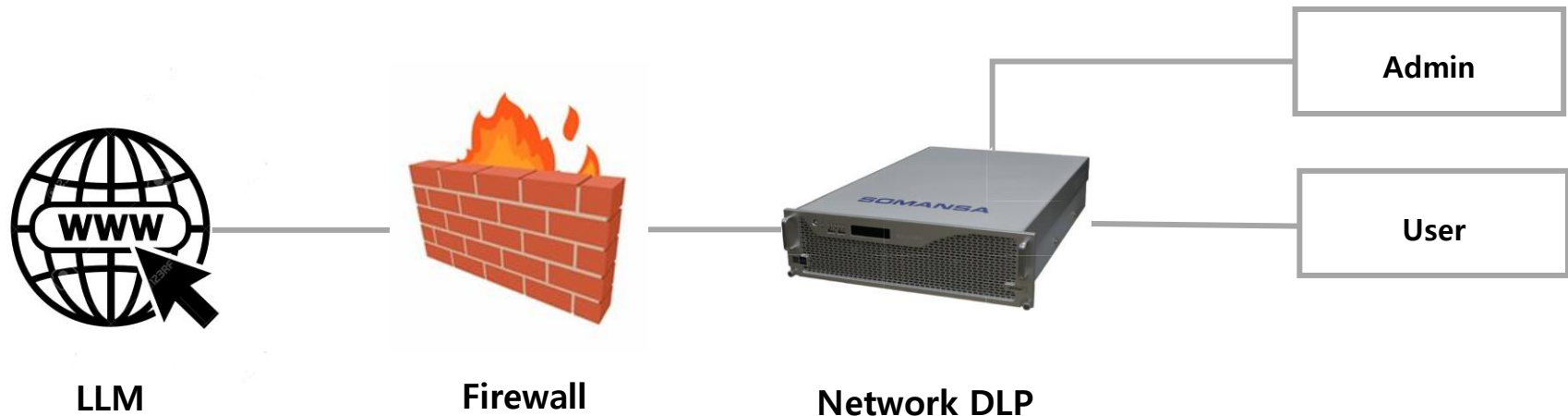
**Emergence of AI-specific data protection services that comply with security guidelines**

# 3

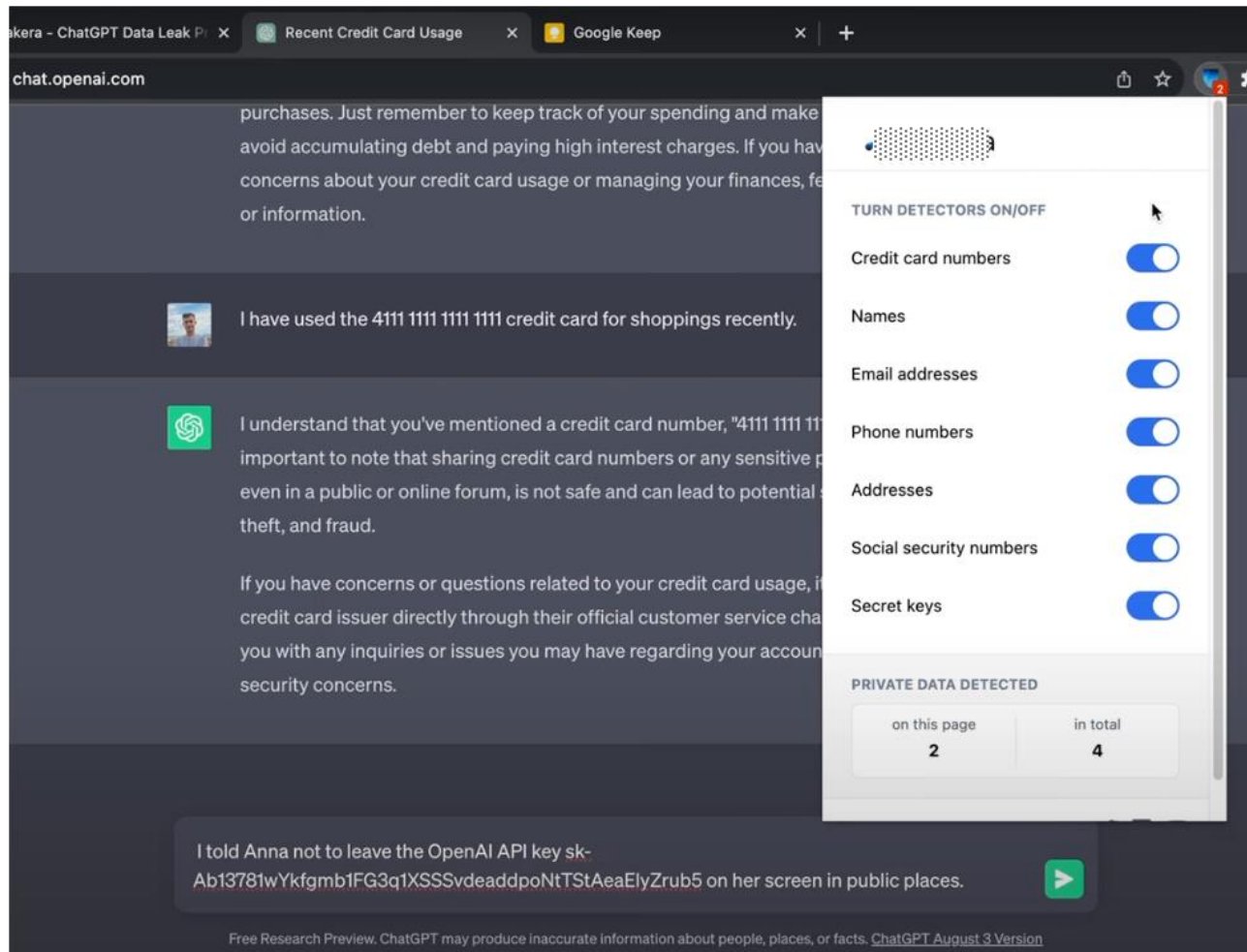**Technical measures for security when using generative AI**

# "Network DLP" Type



**LLM**  **Firewall**  **Network DLP**

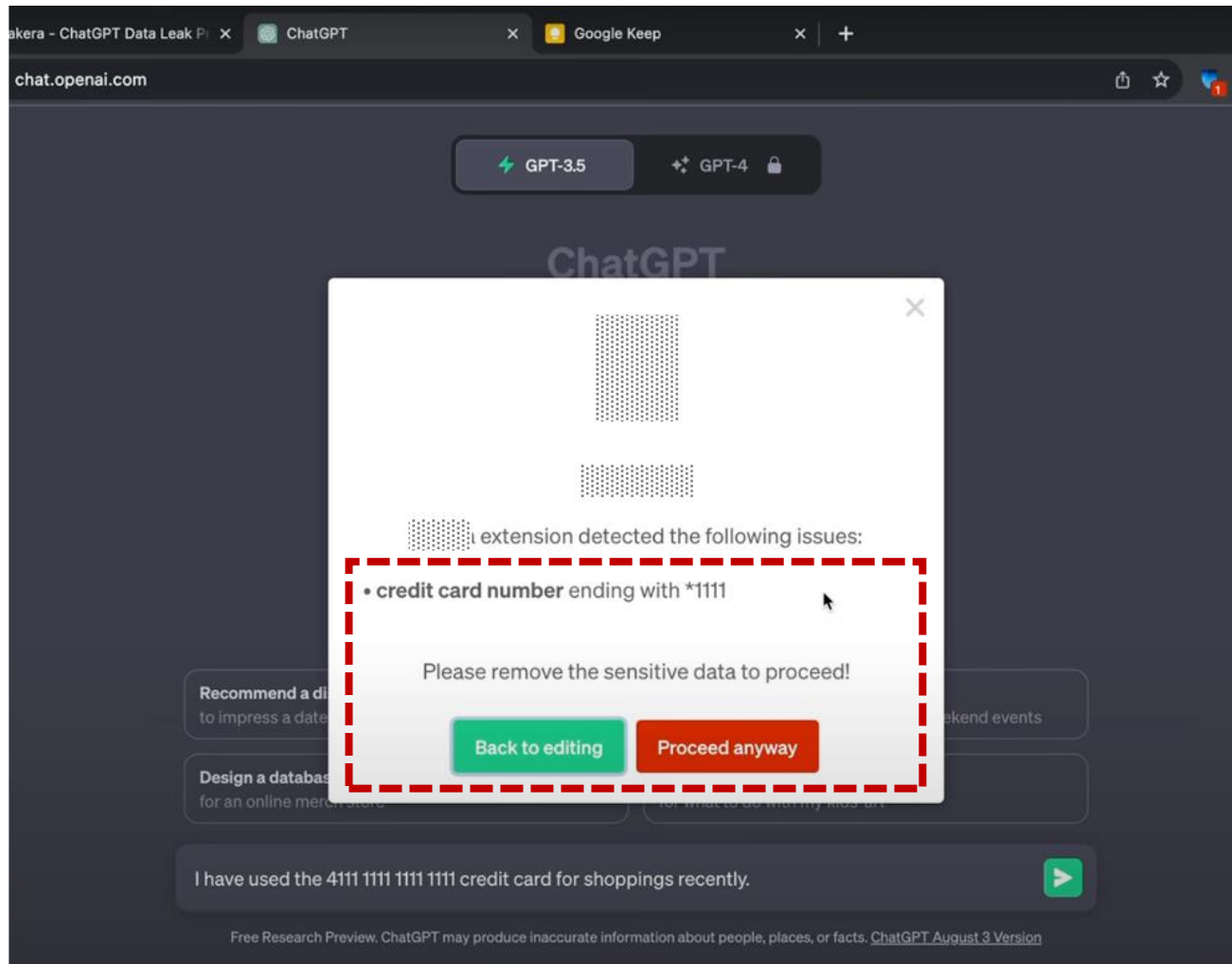**Network DLP : Decrypt HTTPS (SSL/TLS) and block access to unauthorized sites**

**Control over the use of generative AI**
  - **Record both prompts and answers**
  - **Search for specific keywords**
  - **Pre-filtering and blocking when prompts contain personal information (PII)**
  - **Establish departmental and task-specific chatGPT blocking and allowing policies**
  - **If a policy is violated, traffic is blocked, causing business disruption.**

# "Browser Extension" Type

# "Browser Extension" Type

# "API" Type

# "Sandbox" Type

**Generative AI**

prompt

result

| Multilingual detection | Prompt injection |
|---|---|
| Detecting forbidden strings | Code detection (intellectual property) |
| Detecting confidential corporate information | DDoS protection |
| Privacy detection | Secret detection |
| Anonymize personal information | Sensitive information detection |

How can I help you today?

Show me a code snippet
of a website's sticky header

Brainstorm content ideas
for my new podcast on urban design

for an online merch store

for a retro-style arcade game

Message ChatGPT...

| Bias detection | Detecting forbidden strings |
|---|---|
| Code detection | Detecting banned topics |
| Detecting vulnerable code | Language discrimination techniques |
| Recovering anonymization | Malicious link detection |
| Json Completion | Incomplete link detection |
| Detecting harm | Sentiment detection |

prompt

result

**Users / LLM Agent**

# Thank you