



BREAKING CYBERSECURITY NEWS: NEARLY 12,000 API KEYS AND PASSWORDS FOUND IN AI TRAINING DATASET

Vairav Cyber Security News Report

Date: 2025-03-03

Vairav Cyber Threat Intelligence Team

Vairav Technology Security Pvt. Ltd.

Thirbam Sadak 148

Baluwatar, Kathmandu

Phone: +977 4541540

Mobile: +977-9820105900

Email: sales@vairavtech.com

EXECUTIVE SUMMARY

A recent cybersecurity investigation uncovered that nearly 12,000 valid API keys and passwords were exposed within the Common Crawl dataset, a resource often utilized for training large language models (LLMs). These credentials, hardcoded by developers, included sensitive keys for services such as Amazon Web Services (AWS), MailChimp, and WalkScore. The exposure of these secrets underscores the critical need for secure coding practices and vigilant data sanitization during AI model training.

DETAILS OF THE INCIDENT

Description of the Cyber Threat: Analysis of 400 terabytes of data from 2.67 billion web pages in the Common Crawl December 2024 archive, identifying 11,908 valid secrets hardcoded into front-end HTML and JavaScript. This practice exposes sensitive information, making it accessible to unauthorized parties.

Identification: The exposure was identified during Truffle Security's comprehensive scan of the Common Crawl dataset, aiming to detect sensitive data inadvertently included in publicly accessible web archives.

Affected Entities/Industries: The exposed credentials pertain to various services, notably:

- **Amazon Web Services (AWS):** Including root keys found in HTML forms.
- **MailChimp:** Approximately 1,500 unique API keys hardcoded in front-end code.
- **WalkScore:** A single API key appeared 57,029 times across 1,871 subdomains.

Organizations utilizing these services are at heightened risk.

Potential Impact: The risks associated with this exposure include unauthorized access, data exfiltration, service abuse and reputational damage.

Exploitation Methods: Attackers could exploit these hardcoded credentials by directly accessing services and data using the exposed keys, utilizing services like MailChimp to

distribute malicious content or exploiting credentials to impersonate legitimate services or users.

RECOMMENDED ACTIONS

Immediate Mitigation Steps

- **Revoke Exposed Keys:** Immediately invalidate any exposed API keys and generate new ones.
- **Audit Codebases:** Search for hardcoded credentials in code repositories and refactor to use secure methods.
- **Monitor for Abuse:** Implement monitoring to detect unauthorized use of services.

Security Best Practices

- **Environment Variables:** Store sensitive credentials in environment variables, not in code.
- **Access Controls:** Apply the principle of least privilege to API keys and services.
- **Regular Audits:** Conduct periodic security reviews of code and configurations.

For Advanced Security Teams

- **Automated Scanning:** Utilize tools like TruffleHog to detect secrets in codebases.
- **Data Sanitization:** Ensure comprehensive data cleaning processes before using datasets for AI training.
- **Threat Intelligence Integration:** Incorporate threat intelligence feeds to stay informed about exposed credentials.

ADDITIONAL RESOURCES AND OFFICIAL STATEMENTS

- <https://www.bleepingcomputer.com/news/security/nearly-12-000-api-keys-and-passwords-found-in-ai-training-dataset/>
- <https://trufflesecurity.com/blog/research-finds-12-000-live-api-keys-and-passwords-in-deepseek-s-training-data>

CONTACT US

Vairav Technology Security Pvt. Ltd.

Cyber Defender from the land of Gurkha

Thirbam Sadak 148, Baluwatar

Kathmandu, Nepal

Phone: +977-01-4541540

Mobile: +977-9820105900

Email: sales@vairavtech.com

Website: <https://vairavtech.com>