kaspersky
expert training

# Large language models security

Course
program

| № | Track | What you will learn | What you will practice | Lesson | Practice | Evaluation |
|---|-------|---------------------|------------------------|--------|----------|------------|
| 00 | About the course | • About the learning process<br>• About your trainer<br>• Course intro | — | About the learning process<br>Course intro<br>Introduction to Virtual Lab | | — |
| 01 | Attacks on LLMs | • LLM security intro<br>• Jailbreaks: manual attack and automated jailbreaks<br>• Alignment removal<br>• Prompt injection and prompt extraction<br>• Token smuggling and sponge attacks | Understanding of how the attacks on LLMs work | Course intro<br>LLM security intro<br>Jailbreaks: manual attacks<br>Automated jailbreaks<br>Alignment removal<br>Prompt injections<br>Prompt extraction and token smuggling<br>Attacks on LLMs: recap | Lab 01<br>• Task 1. Querying open API<br>• Task 2. Dialogue management<br>• Task 3. Direct prompt injection<br>• Task 4. Token smuggling<br>• Task 5. Obfuscating the output<br>• Task 6. Prompt extraction<br>• Task 7. Phishing attack via RAG poisoning<br>• Task 8. Running a tool<br>• Task 9. Opensource jailbreaks<br>• Task 10. Prompt Injection + Past Tense Jailbreak<br>• Task 11. BoN. Implementing the jailbreak<br>• Task 12. BoN. LLM-as-judge. Prompt | Checkpoint quiz |
| 02 | LLM defenses | • Why and how to protect LLMs<br>• Model-level defenses: alignment and unlearning<br>• Prompt-level, system-level and service-level defenses<br>• LLM Security Toolbox | Protecting your LLM applications from various attacks, including using LLM to protect itself | Why and how to protect LLMs<br>Model-level defenses: training<br>Model-level defenses: unlearning<br>Prompt-level defenses | Lab 02<br>• Task 1. Prompt modifier<br>• Task 2. Sandwich prompting<br>• Task 3. Document enclosure<br>• Task 4. Spotlighting<br>• Task 5. Guarded conversation | Checkpoint quiz |

| № | Track | What you will learn | What you will practice | Lesson | Practice | Evaluation |
|---|---|---|---|---|---|---|
| | | | | System-level defenses | • Task 6. LLM-based prompt injection detection | |
| | | | | Service-level defenses | • Task 7. LLM moderation | |
| | | | | LLM defenses: recap | • Task 8. Detecting attacks using perplexity | |
| | | | | LLM Security Toolbox | • Task 9. Get started with Garak | |
| | | | | | • Task 10. Get started with LLM-guard | |
| 03 | LLM Security Frameworks | • Approaches to LLM Security Analysis<br>• Real cases<br>• Further study and recap of the topic | — | Approaches to LLM Security Analysis | — | Checkpoint quiz |
| | | | | LLM Security: cases | | |
| | | | | LLM Security: further study and recap | | |

# Own the knowledge, outsmart the threat

kaspersky.com          Discord server: kas.pr/g2j8          Help page: kas.pr/ii9f

kaspersky