

AI Red Teaming

Introduction

Foundational Knowledge

- AI / ML Fundamentals
- Cybersecurity Principles

Prompt Hacking

Model Vulnerabilities

System Security

- Code Injection
 - Insecure Deserialization
 - Remote Code Execution
- Infrastructure Security
 - API Protection
 - Authentication
 - Authentication

Professional Development

Real-world Applications

- LLM Security Testing
- Agentic AI Security
- Responsible Disclosure

Future Directions

Find the detailed version of this roadmap and other similar roadmaps

roadmap.sh

AI Security Fundamentals

Why Red Team AI Systems?

Ethical Considerations

Role of Red Teams

Confidentiality, Integrity, Availability

Threat Modeling

Risk Management

Vulnerability Assessment

Jailbreak Techniques

Safety Filter Bypasses

Prompt Injection

Direct

Indirect

Countermeasures

Testing Methodologies

Black Box Testing

White Box Testing

Grey Box Testing

Automated vs Manual

Continuous Testing

Tools and Frameworks

Testing Platforms

Monitoring Solutions

Benchmark Datasets

Custom Testing Scripts

Reporting Tools

Emerging Threats

Advanced Techniques

Research Opportunities

Industry Standards

Visit the following relevant tracks to keep learning

AI Engineer

AI & Data Scientist

Data Analyst

Special thanks to "Learn Prompting" for their help in making this roadmap

learnprompting.org

- ✓ AI and Data Scientist Roadmap
- ✓ AI Engineer Roadmap
- ✓ Data Analyst Roadmap
- ✓ MLOps Roadmap

Supervised Learning

Unsupervised Learning

Reinforcement Learning

Neural Networks

Generative Models

Large Language Models

Prompt Engineering

Model Weight Stealing

Unauthorized Access

Model Extraction

Data Poisoning

Adversarial Examples

Model Inversion

Model Manipulation

Adversarial Training

Robust Model Design

Continuous Monitoring

Defense Strategies

Conferences

Research Groups

Forums

Community Engagement

Lab Environments

CTF Challenges

Red Team Simulations

Practical Experience

Specialized Courses

Industry Credentials

Certifications