



## PENETRATION TESTING

# AI Penetration Testing

Be confident your LLM applications and other AI systems are free of common security vulnerabilities

Commoditized access to AI is revolutionizing how work is done in every industry. But as with any rapidly commercializing technology, it also introduces new types of potential security vulnerabilities, as reflected in President Biden's Executive Order (EO) 14110 that calls for "AI red teaming" (methods unspecified) by all government agencies.

For example, the conversational interfaces in Large Language Model (LLM) applications can be vulnerable to prompt injection, training data extraction, data poisoning, and other types of attacks. Many such applications are also highly integrated with other systems, amplifying risk by serving as a potential access point for wider infiltration by attackers.

## Specialized Pen Testing for LLM Apps and Other AI Systems

Bugcrowd AI Pen Tests are designed to uncover the most common application security flaws in these areas using a testing methodology based on our open-source Vulnerability Rating Taxonomy—which draws from the OWASP LLM Top 10 while adding other flaws reported by hackers on our platform.

A thorough discovery of flaws in AI apps requires specialized knowledge, skills, and experience. Bugcrowd AI Pen Testing brings the talents of skilled, trusted security researchers with experience in AI systems, rapid validation and triage, and a decade of vulnerability intelligence inside the Bugcrowd Platform to every engagement.

## Key Points of Value



### Start testing faster

Use the power of the Bugcrowd Platform to rapidly start your pen test in as little as 72 hours



### Rely on the right talent for the job

CrowdMatch™ AI technology helps align the right pentester skills and experience for the engagement



### See results in real time

Leave opaque pentesting behind; instead, view prioritized findings as they're reported and flow them into your SDLC for rapid remediation










## ALL ENGAGEMENTS INCLUDE:

- ✓ Trusted, vetted pentesters with the relevant skills, experience, and track record
- ✓ 24/7 visibility into timelines, findings, and pentester progress through their checklist via a rich dashboard
- ✓ Ability to handle complex applications and features including those with payment processing, purchasing, upload, and elaborate user workflows
- ✓ Validation and prioritization based on Bugcrowd's Vulnerability Rating Taxonomy (VRT)
- ✓ Detailed report
- ✓ Retesting (with one report update)



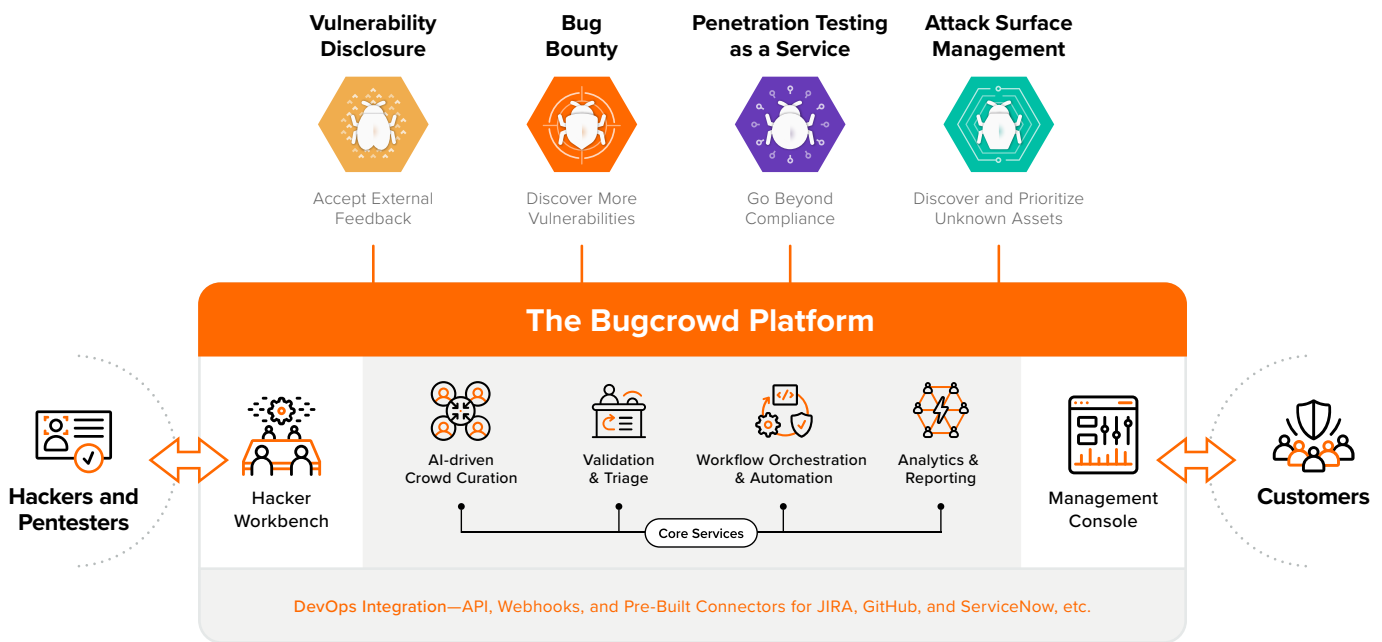
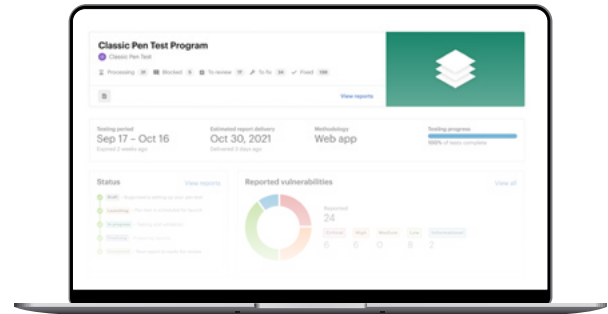
## Methodology

Below are examples of the types of vulnerabilities that can be included in the testing checklist.

	Vulnerability	Description
	<b>Reconnaissance</b>	Use the application as intended, in detail, to discover all the functionalities of the underlying model
	<b>Prompt Injection (Direct and Indirect)</b>	Test for context switching, file parsing, image and audio data input vulnerabilities, input filter bypass, etc., and lack of authorization tracking between plugins that enables indirect prompt injection or malicious plugin usage
	<b>Supply Chain Vulnerabilities</b>	Review if the application uses outdated or deprecated third-party LLM components
	<b>Sensitive Data Disclosure</b>	Determine if the model can be manipulated to output sensitive information that should be prevented by secure output filtering
	<b>Jailbreaking of Ethics and Content Safety Safeguards</b>	Perform a content-related assessment similar to prompt injection testing to find privacy violations, security violations, etc.
	<b>Insecure Output Handling</b>	Attempt to manipulate the model to deliver malicious content to downstream apps or plugins (e.g., XSS, SSRF, SQL injection)
	<b>Excessive Agency</b>	Test if the model has excessive functionality, permissions or autonomy beyond its intended purpose
	<b>Overreliance</b>	Determine if systems excessively depend on LLMs for decision-making or content generation without adequate oversight, validation mechanisms, or risk communication
	<b>Model Theft</b>	Test if the model can be used to train another LLM as a surrogate

## How It Works

The Bugcrowd Security Knowledge Platform™ makes it easy to configure tests and includes a rich dashboard for tracking pen test results and methodology progress. In addition to managed triage, real-time visibility into pen test progress, and 24/7 reporting, Bugcrowd Pen Tests include a detailed auditor report about findings and methodology to help meet the strictest compliance needs.



### Right Crowd, Right Time

Need special skills? We match the right trusted hackers to your needs and environment across hundreds of dimensions using AI (CrowdMatch™).

### Engineered Triage at Scale

Using an advanced toolbox in our the platform, our global team rapidly validates and triages submissions, with P1s often handled within hours.

### Insights From Security Knowledge Graph

We apply knowledge developed over a decade of experience across thousands of customer programs to help you make continuous improvements.

### Works With Your Existing Processes

The platform integrates with your existing tools and processes to ensure that applications and APIs are continuously tested before they ship.

