Firewall for AI (AI 防火墙)

Akamai Firewall for AI (AI 防火墙) 是一款专门设计的安全解决方案,旨在保护 AI 赋能的应用程序、大型语言模型 (LLM) 和 AI 驱动的 API,使其免受新兴网络威胁的侵害。通过保护传入的 AI 查询和传出的 AI 响应,防火墙消除了生成式 AI 引入的安全漏洞。

这款防火墙提供实时检测、基于策略的实施和自适应安全措施等功能,可以防御提示注入、敏感数据泄露、对抗性漏洞利用以及专门针对 AI 的拒绝服务 (DoS) 攻击。

Firewall for AI (AI 防火墙) 可无缝融入边缘位置、云端、混合环境和本地环境,从而确保在所有位置实现一致的安全性、保护能力、治理措施和合规性,同时保持高性能。

防御针对 AI 的威胁

Firewall for AI (AI 防火墙) 针对 AI 驱动的应用程序提供全面的安全防护功能,可以识别并抵御与 AI 相关的漏洞,而这是传统安全工具无法解决的问题。

- 防御提示注入——防止攻击者通过欺骗性输入来操纵 AI 模型。
- **防止数据丢失 (DLP)**——检测并拦截 AI 生成响应中的敏感数据泄露,同时防止在请求中接收敏感数据。
- 有毒和有害内容过滤——标记仇恨言论、错误信息和冒犯性内容,然后再发送出去。
- 对抗性 AI 安全性——防御远程代码执行、模型后门以及数据投毒攻击。
- **缓解拒绝服务攻击**——通过控制过多的查询使用和模型过载,缓解 AI 驱动的 DoS 攻击。

对企业的好处

统一 AI 安全态势

打造跨边缘位置、云端、混合环境和 本地环境的标准化 AI 安全性

- ♠ 自动化 AI 威胁检测 针对 AI 的安全保护, 无需手动调整 规则
- **无缝 WAAP 集成** 提供 AI 防御功能,扩展 Web 应用 程序和 API 保护 (WAAP)
- 防止 AI 滥用和法律风险 阻止数据泄露、IP 盗窃和监管违规行为
- **简化 AI 安全防护** 无需内部工程师手动实施安全策略
- **◇ 多云灵活性** 跨多个环境保护 Al 工作负载
- 企业级 AI 保护功能 以 Akamai 全球威胁情报为后盾



1

灵活的部署选项

Firewall for AI (AI 防火墙) 提供多种部署模型,可针对不同 AI 架构和云环境量身定制。

部署模型	说明
Akamai 边缘集成	在 Akamai 边缘位置实施低延迟的安全措施, 以内联方式保护 AI 应用程序。
REST API	通过基于 API 的风险检测和评分,扫描 AI 输入与输出。
反向代理部署(未来功能)	通过 Akamai 的安全代理来路由 Al 流量, 提供深层次的检测和过滤。

利用这种灵活性,企业可以为部署在任何位置的 LLM 提供安全保护,这包括多云环境、混合环境和本地环境等。

工作原理

AI 流量分析

该防火墙监控并分析 AI 互动,检查传入的用户提示及 AI 生成的输出,检测其中可能存在的威胁,然后再将流量发送给模型或最终用户。通过分析 AI 查询响应循环,该防火墙可以高效地阻止安全风险,同时保持应用程序性能。

风险评分和自适应威胁响应

在评估 AI 互动时,该防火墙会考虑多种安全指标,包括提示注入、敏感数据泄露和对抗性漏洞利用。

安全措施实施操作

Firewall for AI (AI 防火墙) 可根据不同风险评分和客户风险偏好采取三项重要安全措施:

- 监控:记录检测到的威胁供分析,而不干扰 AI 查询或响应。
- 修改: 内联调整 AI 生成的输出,删除或更改不安全的内容,同时保持自然的对话流。
- 拒绝: 阻止高风险输入内容进入 AI 模型,并防止将不安全的响应返回给用户。

安然无忧的合规性和治理

Firewall for AI (AI 防火墙) 可以帮助企业满足安全性和合规性标准。随着 AI 驱动的应用程序引入了新的监管难题,保持对数据隐私、模型完整性和安全风险的监管变得至关重要。

遵循监管要求

该防火墙可以帮助企业遵守隐私、安全与保障指导方针。通过实施针对 AI 的安全策略,企业可以降低与数据保护法规、合乎道德的 AI 使用以及企业治理规定相关的风险。



安全分析和日志记录

Firewall for AI (AI 防火墙) 提供了详细的审计日志和实时安全分析功能,使得安全团队可以深入了解 AI 安全事件。通过监控查询模式、威胁指标和响应行为,企业可以积极主动地检测异常情况、实施策略控制措施以及生成合规性报告。

企业级 AI 保护功能

该防火墙以 Akamai 的全球威胁情报为后盾,能够持续适应新出现的 AI 安全威胁。利用 AI 安全研究和威胁建模团队的实时数据洞察,企业可以保持有弹性的安全态势,同时确保 其 AI 应用程序以安全且负责任的方式运行。



如需了解更多信息,请咨询专家。

