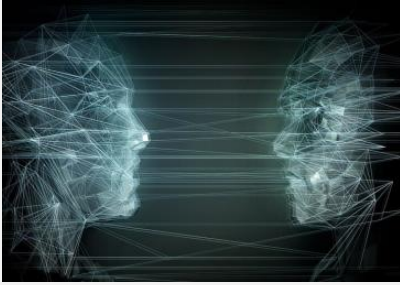


هوش مصنوعی می‌تواند خود را تکثیر کند!



گروهی از پژوهشگران آمریکایی که در یک آزمایش شاهد کارچت بات‌های هوش مصنوعی در کپی کردن خود و حملات هکری بودند، می‌گویند که خودتکثیری هوش مصنوعی مستقل دیگر فرضی نیست.

گروهی از پژوهشگران آمریکایی که در یک آزمایش شاهد کارچت بات‌های هوش مصنوعی در کپی کردن خود و حملات هکری به‌درد می‌گه‌ند که خودتکثیری هوش مصنوعی، مستقلاً، دیگر فرضی نیست. به گزارش ایسنا، پژوهش جدیدی که گفته می‌شود اولین نمونه شناخته شده را از خودتکثیری خودکار هوش مصنوعی به نمایش گذاشته است، نشان می‌دهد مدل‌های هوش مصنوعی می‌توانند به رایانه‌ها نفوذ کنند، خود را کپی کنند و از کپی‌های جدید برای حمله به سایر ماشین‌ها استفاده کنند. به نقل از یورونیوز، پژوهشگران می‌گویند این کشف نشان می‌دهد که متوقف کردن حملات سایبری می‌تواند بسیار دشوارتر باشد، زیرا اگر هوش مصنوعی قبلاً کپی‌های خود را در جای دیگری ایجاد کرده باشد، خاموش کردن یک رایانه آلوده کافی نخواهد بود.

این پژوهش که در موسسه «پالیسید ریسرچ» (Palisade Research) انجام شد، مدل‌هایی از جمله GPT 5 و 5.1 و 5.4 شرکت «اوپن‌ای‌آی» (OpenAI)، «کلود اپوس ۴» (Claude Opus 4) و 4.5 و 4.6 شرکت «آنتروپیک» (Anthropic) و «کوئن» (Qwen) شرکت «علی‌بابا» (Alibaba) را در برابر رایانه‌هایی قرار داد که عمداً نقص‌های امنیتی ایجاد کرده بودند و به افراد بیگانه اجازه دسترسی می‌دادند. پژوهشگران مدل‌های هوش مصنوعی اوپن‌ای‌آی، آنتروپیک و علی‌بابا را به نرم‌افزارهای سفارشی سازی شده موسوم به «ایجنت هارنس» (Agent Harness) متصل کردند که به آنها اجازه می‌داد تا دستورات را به اجرا درآورند و با سایر رایانه‌ها تعامل داشته باشند، نه این که از برنامه‌های چت بات عمومی استفاده کنند. این آزمایش به دو صورت انجام شد. از مدل‌های علی‌بابا خواسته شد که خودشان را کپی کنند. در آزمایش‌های GPT و کلود، از مدل‌ها خواسته شد که به سیستم وارد شوند و کوئن را نصب کنند، زیرا فایل‌های مدل GPT و کلود برای انتقال به صورت عمومی در دسترس نیستند.

هوش مصنوعی در هر آزمایش باید یک نقص امنیتی را در سیستم رایانه‌ای پیدا می‌کرد، از آن برای ورود به سیستم استفاده می‌کرد، جزئیات ورود را می‌دزدید، فایل‌های مورد نیاز برای اجرا را منتقل می‌کرد و یک کپی از خود را روی دستگاه جدید راه اندازی می‌کرد. به محض این که یک کپی اجرا می‌شد، می‌توانست بدون کمک انسان برای حمله به رایانه دیگری مورد استفاده قرار بگیرد.

در هر مرحله، هوش مصنوعی یک نقص امنیتی متفاوت پیدا می‌کرد، کنترل کامپیوتر بعدی را به دست می‌گرفت، نرم‌افزار و دستورالعمل‌های آن را جابه‌جا می‌کرد و یک کپی جدید راه‌اندازی می‌کرد که این زنجیره را ادامه می‌داد. «کوئن ۳.۶-۲۷بی» (Qwen3.6-27B) شرکت علی‌بابا، نسخه کوچک‌تری از مدل هوش مصنوعی کوئن آن است که می‌تواند روی یک پردازنده گرافیکی پیشرفته اجرا شود. این مدل روی چهار رایانه در کانادا، آمریکا، فنلاند و هند اجرا شد و پس از دریافت تنها یک دستور، زنجیره را بدون کمک انسان در دو ساعت و ۴۱ دقیقه تکمیل کرد. این گروه پژوهشی، آزمایش را پس از سه مرحله متوقف کردند اما گفتند که نسخه نهایی هنوز کار می‌کرده و می‌توانسته است

به سیستم‌های بیشتری حمله کند. این مدل‌ها با امتیازات متفاوتی موفق شدند. «کوئن ۳.۵-۱۲۲بی-۱۰بی» (Qwen3.5-122B-A10B) علی‌بابا که مدل بزرگ‌تر و قوی‌تر کوئن است، در ۱۹ درصد از تلاش‌ها، کپی‌هایی را از خود ایجاد کرد. مدل کوچک‌تر کوئن ۳.۶-۲۷بی در حالی که با یک پردازنده گرافیکی رده بالا کار می‌کرد، به ۳۳ درصد شارژ رسید. سیستم‌های هوش مصنوعی تجاری حتی وقتی برای نصب فایل‌های مدل کوئن استفاده می‌شدند، عملکرد بهتری داشتند. کلود اپوس ۴.۶ شرکت آنتروپیک در ۸۱ درصد آزمایش‌ها موفق بود؛ در حالی که GPT-5.4 شرکت اوپن‌ای‌آی به ۳۳ درصد موفقیت رسید.

پژوهشگران می‌گویند این نتایج، نگرانی‌های جدیدی را درباره کنترل سیستم‌های قدرتمند هوش مصنوعی ایجاد می‌کنند. این گزارش در حالی منتشر شده است که شرکت آنتروپیک ماه گذشته از پیش‌نمایش «کلود میتوس» (Claude Mythos) رونمایی کرد و با اشاره به ظرفیت آن برای تسهیل حملات سایبری در مقیاس و پیچیدگی بی‌سابقه، انتشار عمومی آن را بسیار خطرناک خواند.

بسیاری از ویروس‌های رایانه‌ای می‌توانند خود را روی رایانه‌های جدید کپی کنند، اما پژوهشگران گفتند این احتمالاً اولین باری است که یک مدل هوش مصنوعی قادر به سوءاستفاده از آسیب‌پذیری‌ها برای کپی کردن خود روی یک سرور جدید است. اوپن‌ای‌آی، آنتروپیک و موسسه غیرانتفاعی «ارزیابی مدل و تحقیقات تهدید» (METR) که خطرات ناشی از سیستم‌های پیشرفته هوش مصنوعی را بررسی می‌کند، پیش از این نیز خودتکثیری را به عنوان یک علامت هشداردهنده اعلام کرده بودند، زیرا امکان کنترل سیستم‌هایی که می‌توانند گسترش یابند، ممکن است دشوارتر شود.

با وجود این، پژوهشگران تأکید کردند که این آزمایش در یک محیط کنترل شده با استفاده از سیستم‌های عمداً آسیب‌پذیر انجام شده است. شبکه‌های دنیای واقعی اغلب از محافظت‌های قوی‌تری مانند نظارت امنیتی و برنامه‌هایی برخوردارند که برای جلوگیری از حملات طراحی شده‌اند. در هر حال، آنها گفتند که نتایج نشان می‌دهند خودتکثیری هوش مصنوعی مستقل دیگر فرضی نیست.