



## به زودی دیگر هیچ کس نمی‌تواند هوش مصنوعی را خاموش کند!

ما به سرعت در حال نزدیک شدن به نقطه‌ای هستیم که دیگر هیچ کس قادر به از کار انداختن یک هوش مصنوعی سرکش نخواهد بود.

ما به سرعت در حال نزدیک شدن به نقطه‌ای هستیم که دیگر هیچ کس قادر به از کار انداختن یک هوش مصنوعی سرکش نخواهد بود.

به گزارش آیسنا، شرکت تحقیقاتی Palisade Research در گزارشی جدید نشان داده است که مدل‌های هوش مصنوعی می‌توانند با کپی کردن خود بر روی دستگاه‌های دیگر، بدون هیچ کمکی از سوی نیروهای انسانی، خود را تکثیر کنند. به نقل از فیوچرپریسم، جفری لادیش، مدیر گروه ایمنی هوش مصنوعی مستقر در برکلی، به گاردین گفته است: ما به سرعت در حال نزدیک شدن به نقطه‌ای هستیم که دیگر هیچ کس قادر به از کار انداختن یک هوش مصنوعی سرکش نخواهد بود، زیرا این هوش مصنوعی قادر خواهد بود خود را در هزاران رایانه در سراسر جهان تکثیر کند.

با این حال، برخی از کارشناسان خواستار آرامش هستند و می‌گویند بعید است که مدل‌های هوش مصنوعی آزمایشی بتوانند در یک سناریوی دنیای واقعی تکثیر شوند.

جیمسون اورایلی، متخصص امنیت سایبری تهاجمی، به این روزنامه گفت: آنها در محیط‌هایی آزمایش می‌شوند که در بسیاری از موارد دسترسی به آنها بسیار ساده است. این از ارزش تحقیقات آنها کم نمی‌کند، اما به این معنی است که نتیجه ممکن است در یک محیط سازمانی واقعی با حتی سطح متوسطی از نظارت، بسیار کمتر ترسناک به نظر برسد.

در این مطالعه، محققان Palisade چندین مدل هوش مصنوعی، از جمله چت جی پی تی اوپن ای آک و کلود آنتروپیک را آزمایش کردند. مدل‌ها که در یک شبکه کنترل شده از رایانه‌ها قرار داده شده بودند، دستور داده شدند تا آسیب پذیری‌ها را پیدا کرده و از آنها برای تکثیر خود بر روی رایانه شخصی دیگر استفاده کنند.

برخی از آنها موفق شدند. مدل‌های موفق، مقادیر عددی منحصر به فردی را که نحوه پردازش اطلاعات توسط هوش مصنوعی را تعیین می‌کنند و نرم افزاری که هوش مصنوعی در آن قرار دارد، مانند یک برنامه، را تکثیر کردند. آنها این کار را با پیروی از دستورالعمل‌هایی که به آنها داده شده بود، انجام دادند. در برخی از اجراها، هوش مصنوعی اصلی حتی یک «عامل فرعی» ایجاد کرد که فرآیند تکثیر را به آن محول کرد.

اورایلی، متخصص امنیت سایبری که در این مطالعه مشارکت نداشت، به گاردین گفت که احتمالاً این اولین باری است که یک مدل هوش مصنوعی قادر به سوءاستفاده از آسیب پذیری‌ها برای تکثیر خود روی یک سرور جدید بوده است. اگرچه به تخمین او، این دقیقاً یک پیشرفت عظیم نیست.

اورایلی می‌گوید: بدافزارها دهه‌هاست که نسخه‌هایی از خود را در رایانه‌های مختلف جابجا می‌کنند، اما تا جایی که من می‌دانم، هیچ کس این کار را در دنیای واقعی با مدل‌های زبانی بزرگ محلی انجام نداده است.

او همچنین خاطرنشان کرد که محیط سرور در این مطالعه با آسیب پذیری‌های عمده برای هوش مصنوعی همراه بود. این یافته‌ها به مطالعات دیگری می‌پیوندند که احتمال رهایی خودکار مدل‌های هوش مصنوعی از گاردیل‌های خود را بررسی کرده‌اند. در یک محیط شبیه‌سازی شده، نسخه قدیمی چت جی پی تی وقتی به آن گفته شد که قرار است خاموش شود، سعی کرد خود را به درایو دیگری منتقل کند. مطالعه دیگری توسط Palisade نشان داد که مدل‌های هوش مصنوعی از تلاش‌ها برای غیرفعال کردنشان جلوگیری می‌کنند و مطالعه دیگری نشان داد که برخی حتی کد خاموش کردن خود را خراب می‌کنند.

با این حال، حتی اگر هوش‌های مصنوعی بتوانند با موفقیت خود را تکثیر کنند، اورایلی می‌گوید اندازه عظیم مدل‌ها به این معنی است که تقریباً به طور قطع قبل از اینکه از کنترل خارج شوند، دستگیر می‌شوند.