

خیال خام کنترل ایمن هوش مصنوعی

دکتر رومن وی یامپولسکی از کارشناسان برجسته هوش مصنوعی می‌گوید ما با یک رویداد تقریباً تضمین شده با پتانسیل ایجاد یک فاجعه وجودی روبرو هستیم.



دکتر رومن وی یامپولسکی از کارشناسان برجسته هوش مصنوعی می‌گوید ما با یک رویداد تقریباً تضمین شده با پتانسیل ایجاد یک فاجعه وجودی روبرو هستیم.

به گزارش ایسنا، یک متخصص در زمینه هوش مصنوعی، بررسی کاملی روی ادبیات علمی در مورد هوش مصنوعی انجام داده و به این نتیجه رسیده است که هیچ مدرکی برای حمایت از این تصور وجود ندارد که هوش مصنوعی را می‌توان به طور ایمن کنترل کرد.

علاوه بر این، حتی اگر مکانیسم‌ها یا روش‌هایی برای کنترل هوش مصنوعی پیشنهاد شده باشد، دکتر رومن وی یامپولسکی (Roman V. Yampolskiy) معتقد است که این اقدامات برای تضمین ایمنی کافی نخواهند بود.

به نقل از آی‌ای، دکتر یامپولسکی دانشیار دانشگاه لوئیزیولا (Louisville) است که نتایج تحقیقات خود را در کتابی با عنوان «هوش مصنوعی غیر قابل توضیح، غیر قابل پیش‌بینی، غیر قابل کنترل» منتشر کرده است.

کار او در زمینه «ایمنی هوش مصنوعی» تا حدی توسط ایلان ماسک و موسسه آینده‌زنگی (FLI) که توسط مکس تیگمارک استاد دانشگاه MIT و زان تالین یکی از بنیانگذاران اسکایپ تاسیس شده است، پشتیبانی شده است.

موسسه آینده‌زنگی همان موسسه‌ای است که برای توسعه ابزارهای هوش مصنوعی قدرتمندتر از GPT-4 درخواست کرد این پروژه با یک توقف شش ماهه روبرو شود.

نامه سرگشاده این موسسه توسط بیش از 33 هزار نفر از جمله ایلان ماسک، استیو ورنیاک هم‌بنیانگذار اپل، یوشوا بنجیو (Yoshua Bengio) پدرخوانده هوش مصنوعی و سایر کارشناسان امضا شده است.

نگرانی در مورد مدیریت و تنظیم مقررات در مقابل خطرات هوش مصنوعی

با نفوذ ابزارهای جدید هوش مصنوعی در بازار طی یک سال گذشته، دکتر یامپولسکی پیشنهاد می‌کند که سیستم‌های هوش مصنوعی پیشرفته به دلیل غیر قابل پیش‌بینی بودن و استقلال ذاتی، با وجود مزایای بالقوه‌شان، همیشه خطرانی را به همراه خواهند داشت.

دکتر یامپولسکی در یک بیانیه مطبوعاتی گفت: چرا بسیاری از پژوهشگران تصور می‌کنند که مشکل کنترل هوش مصنوعی قابل حل است؟ تا آنجا که ما می‌دانیم، هیچ گواهی برای آن وجود ندارد، هیچ مدرکی وجود ندارد. قبل از شروع تلاش برای ساخت یک هوش مصنوعی کنترل شده، مهم است که نشان دهیم مشکل قابل حل است.

وی افزود: این همراه با آماري که نشان می‌دهد توسعه آبرهوش مصنوعی یک رویداد تقریباً تضمین شده است، نشان می‌دهد که ما باید از تلاش‌های ایمنی قابل توجهی در زمینه هوش مصنوعی حمایت کنیم.

توسعه ابرهوش مصنوعی اجتناب ناپذیر است

از آنجایی که هوش مصنوعی در حالت ابرهوشمندی می‌تواند به تنهایی یاد بگیرد، تطبیق یابد و به صورت نیمه مستقل عمل کند، اطمینان از ایمن بودن آن، به ویژه با وجود افزایش قابلیت‌های آن به طور فزاینده‌ای چالش برانگیز می‌شود.

می‌توان گفت که هوش مصنوعی فوق‌هوشمند «ذهن» خود را خواهد داشت. پس چگونه آن را کنترل کنیم؟ آیا قوانین رباتیک آیزاک آسیموف که رفتار اخلاقی بین انسان و ربات را برقرار می‌کند، هنوز در دنیای امروز اعمال می‌شود؟

دکتر یامپولسکی که حوزه اصلی مورد علاقه او ایمنی هوش مصنوعی است، توضیح داد: ما با یک رویداد تقریباً تضمین شده با پتانسیل ایجاد یک فاجعه وجودی روبرو هستیم.

وی افزود: جای تعجب نیست که بسیاری این را مهم‌ترین مشکلی که بشریت تا به حال با آن روبرو بوده است، توصیف می‌کنند. نتیجه می‌تواند رفاه یا انقراض بشر باشد.

یکی از چالش‌های اصلی، تصمیم‌گیری‌های بالقوه و شکست موجودات فوق‌هوشمند است که پیش‌بینی و تهدید آنها را دشوار می‌کند.

علاوه بر این، فقدان قابلیت توضیح در تصمیمات هوش مصنوعی نگرانی‌هایی را ایجاد می‌کند، زیرا درک این تصمیمات برای کاهش تصادفات و تضمین نتایج بدون سوگیری، به ویژه در زمینه‌های حیاتی مانند مراقبت‌های بهداشتی و مالی که هوش مصنوعی در آن نفوذ کرده است، بسیار مهم است.

با افزایش استقلال هوش مصنوعی، کنترل انسان روی آن کاهش می‌یابد که منجر به نگرانی‌های ایمنی می‌شود. دکتر یامپولسکی ایجاد تعادل بین توانایی و کنترل هوش مصنوعی را پیشنهاد و اذعان می‌کند که ابرهوش ذاتاً فاقد قابلیت کنترل است.

وی همچنین درباره اهمیت همسویی هوش مصنوعی با ارزش‌های انسانی بحث می‌کند و روش‌هایی را برای به حداقل رساندن خطرهای مانند اصلاح‌پذیری، شفافیت هوش مصنوعی و طبقه‌بندی به عنوان قابل کنترل یا غیرقابل کنترل پیشنهاد می‌کند.

کند.

با وجود این چالش‌ها، دکتر یامپولسکی ادامه تحقیقات و سرمایه‌گذاری در ایمنی و امنیت هوش مصنوعی را تشویق و تأکید می‌کند که اگرچه دستیابی به هوش مصنوعی صد درصد ایمن ممکن است دست نیافتنی باشد، اما تلاش برای بهبود ایمنی آن همواره ارزشمند است.