



## هوش مصنوعی هنوز حریف ریاضیدانان برتر نمی‌شود

یک معیار جدید که هوش مصنوعی را در برابر مسائل ریاضی از پیش دیده نشده قرار می‌دهد، نشان می‌دهد که این سامانه‌ها هنوز تا رسیدن به سطح برترین متخصصان انسانی فاصله دارند.

یک معیار جدید که هوش مصنوعی را در برابر مسائل ریاضی از پیش دیده نشده قرار می‌دهد، نشان می‌دهد که این سامانه‌ها هنوز تا رسیدن به سطح برترین متخصصان انسانی فاصله دارند.

به گزارش اسپنا، هوش مصنوعی دقیق‌ترین آزمون ریاضی خود را تاکنون پشت سر گذاشته است. نتایج آماده است و مدل‌های هوش مصنوعی که در آن شرکت کردند، به مهارت‌های حل مسئله ریاضیدانان برتر نرسیدند.

به نقل از نیچر، این آزمون بخشی از پروژه‌ای به نام First Proof به معنای اثبات اول است که هدف از آن ارزیابی توانایی هوش مصنوعی در حل سوالات پیچیده در ریاضیات است. ۱۰ مسئله ریاضی در سطح تحقیق برای چهار سیستم هوش مصنوعی مطرح شد. سپس هیئت منصفه‌ای از متخصصان انسانی ناشناس در زمینه‌های ریاضی مربوطه، پاسخ‌های مدل‌ها را ارزیابی کردند. این آزمون اولین آزمون از نوع خود بود که به طور همزمان سه شرط کلیدی را برآورده می‌کرد: اول، شامل سوالات ریاضی در سطح تحقیق بود؛ دوم، شامل مسائلی بود که در داده‌های آموزشی ظاهر نشده بودند و سوم، به طور رسمی توسط ریاضیدانان درجه بندی شده بود.

این یافته‌ها پس از پیشرفت‌های اخیر هوش مصنوعی در حل مسائل ریاضی رخ داده است. برای مثال، ماه گذشته، یک چت بات ساخته شده توسط شرکت فناوری اوپن‌ای‌آی در سانفرانسیسکو، کالیفرنیا، یک چالش ریاضی ۸۰ ساله را که توسط ریاضیدان فقید، پاول اردوش، مطرح شده بود، حل کرد. گروه فرست پروف می‌گوید که تکرارهای آینده این آزمون می‌تواند به محققان کمک کند تا قضاوت کنند که مدل‌های هوش مصنوعی چقدر می‌توانند به عنوان مثال، در حل خودکار مسائل، بررسی اثبات‌ها یا ایفای نقش دستیاران تحقیق برای ریاضیدانان مفید باشند.

یکی از نوآوری‌های مهم آزمون فرست پروف این بود که سوالات قبلاً در هیچ کجای مقالات منتشر شده یا در اینترنت ذکر نشده بودند و این خطر را که مدل‌ها به سادگی اطلاعاتی را که در طول آموزش خود آموخته‌اند، تکرار کنند، از بین می‌برد. در عوض، ۱۰ محقق از طیف گسترده‌ای از تخصص‌های ریاضی، هر کدام سوالاتی را ارائه دادند که در جریان تحقیقات خود حل کرده بودند، اما هنوز منتشر نکرده بودند.

فرست پروف در ماه فوریه یک آزمون آزمایشی با دسته‌ای متفاوت از مسائل جدید برگزار کرد. در آن دور، هر کسی می‌توانست سیستم‌های هوش مصنوعی مورد علاقه خود را روی مسائل امتحان کند و بسیاری از گروه‌ها این کار را انجام دادند اما نتایج به طور رسمی تأیید نشد. همچنین هیچ راهی برای بررسی مستقل اینکه هوش مصنوعی کمکی از انسان‌ها دریافت نکرده است یا خیر، وجود نداشت. این بار، فرست پروف خودش آزمایش را اجرا کرد؛ گروه از مدل‌ها خواست تا مسائل را به روشی کاملاً خودکار حل کنند و گروهی متشکل از ۳۰ ریاضیدان برای بررسی پاسخ‌ها به کار گرفته شدند. جرمی آویگاد، ریاضیدان و رئیس موسسه استدلال به کمک رایانه در ریاضیات در دانشگاه کارنگی ملون در پیتسبورگ، پنسیلوانیا، می‌گوید: برگزارکنندگان به

وضوح با دقت بیشتری به دسته دوم فکر کرده‌اند تا آن را کنترل شده‌تر و سیستماتیک‌تر کنند. قانون دیگر این بود که مدل‌های شرکت‌کننده باید در دسترس عموم باشند. این بدان معنا بود که Aletheia گوگل که سیستمی است که به طور خاص برای حل مسائل ریاضی طراحی شده است و نسخه کامل و منتشر نشده Claude Mythos، مدلی که توسط آنتروپیک در سانفرانسیسکو، کالیفرنیا ساخته شده است، قابل استفاده نبودند. اوپن‌ای‌آی تنها شرکت بزرگی بود که با مدل ChatGPT ۵.۵ Pro خود در این آزمایش شرکت کرد.

سیستم‌های دیگر توسط سه گروه دانشگاهی، از دانشگاه کالیفرنیا، لس‌آنجلس (UCLA)، دانشگاه پرینستون در نیوجرسی و موسسه فناوری فدرال سوئیس (ETH) در زوریخ ارائه شدند. هر سه دانشگاه، «مهارهایی» بر روی چت بات‌های موجود، مانند چت‌جی‌پی‌تی، جمینای گوگل و نسخه عمومی کلاود آنتروپیک ساختند. مهار، سیستمی خودکار است که از یک چت بات سوالاتی می‌پرسد و پاسخ آن توسط یک چت بات دیگر، اغلب با رفت و برگشت‌های مکرر، بررسی می‌شود.

### نتایج ریاضی

مدل تیم موسسه فناوری فدرال سوئیس بهترین عملکرد را داشت و ۶ از ۱۰ مسئله را با سیستمی حل کرد که در آن پاسخ‌های چت‌جی‌پی‌تی توسط یک «شورای مشورتی» متشکل از هر سه چت بات اصلی بررسی شده یا بهبود یافته بود. تیم دانشگاه کالیفرنیا، لس‌آنجلس که یک مهار بر روی چت‌جی‌پی‌تی ساخته بود، دومین تیم برتر بود و پس از آن تیم چت‌جی‌پی‌تی اوپن‌ای‌آی بدون مهار و پرینستون مهاره‌ای که عمدتاً از Gemini ۲.۱ Pro به عنوان پشتیبان خود استفاده می‌کند، قرار گرفتند.

یوهانس اشمیت، ریاضیدانی که عضوی از تیم موسسه فناوری فدرال سوئیس بود، می‌گوید که برای تنظیم دقیق سیستم خود قبل از مسابقه، او و همکارانش با جامعه ریاضی گسترده‌تر تماس گرفتند و از آنها مسائل را درخواست کردند. پاسخ شگفت‌انگیز

بود: طرف چند روز، ۳۰ مسئله ارسالی از حوزه‌های مختلف ریاضیات دریافت کردیم و مردم بسیار کنجکاو و روشن‌فکر بودند. لورن ویلیامز، ریاضیدان دانشگاه هاروارد در کمبریج، ماساچوست و عضو تیم فرست پروف، می‌گوید: مشخص نیست که آیا مسائل حل نشده لزوماً سخت‌تر از بقیه بوده‌اند یا خیر. او می‌افزاید: من فکر می‌کنم مسائلی که حل نشده بودند، چه از نظر موضوع و چه از نظر ایده‌های اثبات، از چیزهایی که قبلاً در مقالات علمی آمده بودند، دورتر بودند.

مدل‌های استدلال همچنین مستعد توهم یا تولید خروجی‌های واقعاً نادرست بودند، حتی زمانی که صریحاً به آنها گفته می‌شد که منابع خود را بررسی کنند که مشکلی شناخته شده در مدل‌های زبانی بزرگ است. ویلیامز می‌گوید که از کمبود «شدید» استناد در تمام پاسخ‌های مدل‌های هوش مصنوعی شگفت‌زده شده است به ویژه در

مورد مسئله ۲، که چندین مدل با اقتباس از روشی که یک مسئله مشابه در گذشته توسط انسان ها حل شده بود، آن را حل کردند. چندین راه حل، در برخی موارد، کپی کردن عبارات از مقاله قبلی به صورت خط به خط و استفاده مجدد از نمادگذاری ها و اصطلاحات دقیق بود، اما هرگز به آن مقاله در هیچ کجا استناد نکردند. اکنون که مسائل فرست پروف منتشر شده اند، شرکت هایی که رسماً در آن شرکت نکرده اند، احتمالاً از آنها برای آزمایش غیررسمی سیستم های خود استفاده خواهند کرد. کوین بارتو، ریاضیدان دانشگاه کمبریج انگلستان که معیارهای ریاضی غیررسمی خود را برای هوش مصنوعی اجرا کرده است، می گوید: شخصاً از دیدن مدل های داخلی آزمایش شده از سه آزمایشگاه لذت می بردم، فقط برای اینکه ببینم مرز واقعی در حال حاضر کجاست.