



شرکت‌های فناوری به دنبال مدل‌های کوچک‌تر هوش مصنوعی هستند

بسیاری از توسعه‌دهندگان فناوری تمایل دارند تا از نمونه‌های کوچک‌تر مدل‌های هوش مصنوعی استفاده کنند. این گزارش به بررسی دلایل این تمایل می‌پردازد.

بسیاری از توسعه‌دهندگان فناوری تمایل دارند تا از نمونه‌های کوچک‌تر مدل‌های هوش مصنوعی استفاده کنند. این گزارش به بررسی دلایل این تمایل می‌پردازد.

به گزارش ایسنا، شرکت «اوپن ای آی» (OpenAI) روز پنجشنبه GPT-4o mini را معرفی کرد که نسخه کوچک‌تر و ارزان‌تر مدل هوش مصنوعی GPT-4o آن است.

به نقل از فست کمپانی، اوپن ای آی یکی از چندین شرکت هوش مصنوعی است که نسخه‌ای از بهترین مدل بنیادین خود را توسعه داده تا مقداری از هوش را با سرعت و مقرون به صرفه شدن عوض کند. چنین مبادله‌ای می‌تواند به توسعه‌دهندگان بیشتری امکان دهد تا برنامه‌های خود را با هوش مصنوعی تقویت کنند و شاید در آینده روزنه‌های جدیدی را برای برنامه‌های پیچیده‌تر باز کند.

بزرگ‌ترین مدل‌های زبانی بزرگ از میلیاردها یا تریلیون‌ها پارامتر استفاده می‌کنند تا طیف گسترده‌ای را از استدلال و وظایف مرتبط با پرس و جو انجام دهند. آنها با حجم بزرگی از داده‌ها آموزش دیده‌اند که موضوعات گوناگونی را پوشش می‌دهند. از سوی دیگر، مدل‌های زبانی کوچک فقط از میلیون‌ها یا ده‌ها میلیون پارامتر برای انجام دادن وظایف کمی استفاده می‌کنند و به قدرت محاسباتی کمتر و مجموعه کوچک‌تری از داده‌های آموزشی متمرکز نیاز دارند.

برای توسعه‌دهندگانی که برنامه‌های ساده‌تری دارند، مدل‌های زبانی کوچک ممکن است تنها گزینه قابل اجرا باشند. اوپن ای آی می‌گوید GPT-4o mini تا ۶۰ درصد ارزان‌تر از GPT-3.5 Turbo است که پیش از این مقرون به صرفه‌ترین مدل این شرکت برای توسعه‌دهندگان بود.

یک مورد دیگر، بحث سرعت است. بسیاری از کاربردهای هوش مصنوعی به دانش عمومی گسترده یک مدل زبانی بزرگ نیاز ندارند. آنها ممکن است به پاسخ‌های سریع‌تری برای پرسش‌های ساده‌تر نیاز داشته باشند. «مایک اینتراتور» (Mike Intrator) مدیرعامل «کورویو» (CoreWeave) که میزبان مدل‌های هوش مصنوعی در فضای ابری است، گفت: منظور از تأخیر، زمان مورد نیاز برای یک برنامه هوش مصنوعی به منظور دریافت پاسخ از یک مدل در فضای ابری است. اگر فرزند من مقاله‌ترم خود را با کمک هوش مصنوعی بنویسد، تأخیر مشکل بزرگی نیست اما اگر بخواهد از آن برای جراحی یا رانندگی خودران استفاده کند، تأخیر با تأثیر بسیار بیشتری همراه می‌شود.

اینتراتور خاطرنشان کرد: مدل‌های مورد استفاده در خودروهای خودران باید به اندازه‌ای کوچک باشند که روی یک تراشه رایانه‌ای اجرا شوند، نه در یک سرور فضای ابری.

مدل GPT-4o mini کوچک‌تر از مدل‌های دیگر است اما هنوز آن قدر کوچک نیست که روی دستگاهی مانند تلفن همراه یا کنسول بازی اجرا شود. بنابراین، باید مانند سایر مدل‌های شرکت اوپن ای آی روی یک سرور در فضای ابری اجرا شود.

کلید نسل بعدی برنامه‌های هوش مصنوعی

«رابرت نیشیهارا» (Robert Nishihara) یکی از بنیان‌گذاران و مدیرعامل شرکت «انی اسکیل» (Anyscale) گفت: امروزه بیشتر برنامه‌های کاربردی مبتنی بر هوش مصنوعی شامل یک پرس و جو یا چند پرس و جو برای مدلی هستند که در فضای ابری اجرا می‌شود اما برنامه‌های پیشرفته‌تر به پرس و جوهای زیادی برای مدل‌های گوناگون نیاز دارند. به عنوان مثال، برنامه‌ای که به شما کمک می‌کند یک اقامت‌گاه اجاره‌ای را انتخاب کنید، ممکن است از یک مدل برای ایجاد معیارهای انتخاب، یک مدل دیگر برای انتخاب گزینه‌های اجاره‌ای و یک مدل دیگر برای امتیاز دادن به هر یک از آن گزینه‌ها استفاده کند. کارگردانی و سازمان‌دهی همه این پرسش‌ها یک تجارت پیچیده است.

عملکرد مدل‌ها مهم است اما سرعت و هزینه آنها نیز به یک اندازه اهمیت دارد. اوپن ای آی این موضوع را می‌داند؛ همان‌طور که شرکت‌هایی مانند متا و گوگل می‌دانند.

نیشیهارا ادامه داد: تلاش‌های این شرکت‌ها به منظور کوچک‌سازی مدل‌های هوش مصنوعی برای برنامه‌های پیچیده‌تر مانند دستیارهای شخصی بسیار مهم هستند.

اوپن ای آی اندازه پارامتر مدل‌های خود را فاش نمی‌کند اما مدل جدید آن احتمالاً به اندازه Claude 3 Haiku شرکت «آنتروپیک» (Anthropic) و Gemini ۱.۵ Flash گوگل است.

اوپن ای آی گفت: GPT-4o mini عملکرد بهتری نسبت به مدل‌های مشابه در آزمایش‌ها دارد.

اوپن ای آی خاطرنشان کرد که توسعه‌دهندگان اپلیکیشن می‌توانند از طریق یک واسط برنامه نویسی کاربردی به GPT-4o mini دسترسی داشته باشند و مدل‌های جدید نیز از برنامه ChatGPT پشتیبانی می‌کنند.