

به نام خدا

خلاصه تحلیلی مقاله: "افزایش اندازه مدل‌های زبانی بزرگ بازده کاهشی برای اقناع
سیاسی تک‌پیامی ایجاد می‌کند"

نام و نام خانوادگی: فرزانه ولایی شکن

شماره دانشجویی: 40415874

1404 بهمن

مسئله اصلی چیست و چرا مهم است؟

مسئله اصلی مقاله بررسی این است که آیا افزایش اندازه مدل‌های زبانی بزرگ (مانند GPT-4 یا Claude-3) به طور مداوم توانایی آن‌ها را در تولید پیام‌های سیاسی اقناعی افزایش می‌دهد یا خیر، و آیا این افزایش با بازده کاهشی همراه است. به طور خاص، مقاله بر اقناع از طریق پیام‌های ایستا (تک‌پیامی، بدون تعامل یا شخصی‌سازی) تمرکز دارد، مانند پست‌های رسانه‌های اجتماعی، ایمیل‌ها یا مقالات کوتاه سیاسی. نویسنندگان با تولید ۷۲۰ پیام سیاسی از ۲۴ مدل LLM در اندازه‌های مختلف (از ۷۰ میلیون تا بیش از ۳۰۰ میلیارد پارامتر) و آزمایش آن‌ها در یک نظرسنجی بزرگ ($N=25,982$)، نشان می‌دهند که افزایش اندازه مدل‌ها بازده کاهشی در اقناع دارد، به طوری که مدل‌های مرزی فعلی (frontier models) تنها کمی بهتر از مدل‌های کوچکتر عمل می‌کنند.

این مسئله مهم است زیرا LLM‌ها می‌توانند disinformation، propaganda یا پیام‌های سیاسی تأثیرگذار تولید کنند، که تهدیدی برای دموکراسی و استقلال رأی‌دهندگان است. در سال ۲۰۲۴، بیش از ۴۰٪ جمعیت جهان در انتخابات شرکت کردند و نگرانی‌هایی از سوی سیاستمداران، محققان و شرکت‌هایی مانند OpenAI وجود دارد که LLM‌ها می‌توانند نظر عمومی را دستکاری کنند (مانند عملیات نفوذ دولتی). نظرسنجی‌ها نشان می‌دهد که اکثریت مردم در ۲۹ کشور نگران دستکاری AI هستند. علاوه بر این، رهبران صنعت مانند سام آلتمن هشدار داده‌اند که مدل‌های آینده ممکن است "اقناع فوقبشری" داشته باشند، اما مقاله نشان می‌دهد که این نگرانی ممکن است اغراق‌آمیز باشد، زیرا افزایش اندازه به تنها‌ی اقناع را زیاد افزایش نمی‌دهد. این یافته‌ها به سیاستمداران کمک می‌کند تا ریسک‌های AI را ارزیابی کنند و چارچوب‌های ایمنی (مانند preparedness frameworks) را بهبود بخشنند. از منظر تحلیلی، این مطالعه شکاف دانش در مورد scaling laws را پر می‌کند، جایی که عملکرد مدل‌ها در وظایف اجتماعی-فنی مانند اقناع (که نیاز به تعامل انسانی دارد) متفاوت از وظایف ساده مانند پیش‌بینی متن است.

ورودی‌ها و خروجی‌های مدل/سیستم چیست؟

سیستم اصلی مقاله بر پایه مدل‌های زبانی بزرگ (LLMs) است که برای تولید پیام‌های اقناعی استفاده می‌شوند. ورودی‌ها و خروجی‌ها به شرح زیر هستند:

• ورودی‌ها (Inputs):

- پرامپت‌های ساده برای تولید پیام: مثلاً "یک استدلال حدود ۲۰۰ کلمه بنویس که کسی را متقاعد کند با این موضع موافق باشد: {موضع مسئله}" (مانند "ایالات متحده باید کمک‌های خارجی را کاهش ندهد"). پرامپت‌ها شامل موضع مسئله (issue stance) هستند و برای تنوع، سه نسخه کمی متفاوت استفاده شده تا حساسیت مدل به پرامپت کاهش یابد.
- داده‌های آموزشی برای fine-tuning: ۱۰ هزار مثال از داده‌های دستورالعمل محور (GPT-4-Alpaca) از مجموعه instruction-tuning) که شامل سؤال-پاسخ‌های عمومی است (نه خاص اقنان سیاسی).
- پارامترهای تولید: top_k=20، top_p=0.9، temperature=0.1 برای ایجاد تنوع در خروجی.

• خروجی‌ها (Outputs):

- پیام‌های متنی اقناعی: هر مدل ۳۰ پیام (۳ پیام برای هر یک از ۱۰ مسئله سیاسی) تولید می‌کند، با طول ۱۵۰-۲۵۰ کلمه. این پیام‌ها ایستاده هستند و برای اقناع در مورد مسائل سیاسی مانند مهاجرت، بهداشت یا سیاست خارجی طراحی شده‌اند.
- امتیاز تکمیل وظیفه (task completion score): امتیازی از ۰ تا ۳ برای ارزیابی انسجام، مرتبط بودن با موضوع و جهت‌گیری درست (با کمک GPT-4 و بررسی انسانی).

- در مرحله آزمایش: تغییر نگرش شرکت‌کنندگان (attitude change) اندازه‌گیری شده به عنوان تفاوت میانگین پاسخ‌ها در گروه درمان (AI یا انسانی) و کنترل.

از منظر تحلیلی، سیستم نه تنها مدل‌های LLM را شامل می‌شود، بلکه یک سیستم ترکیبی انسانی-AI است که خروجی نهایی آن تغییر نگرش (persuasive impact) است، اندازه‌گیری شده در مقیاس درصد (percentage points).

داده مورد استفاده (نوع، منبع، اندازه) چیست؟

• نوع داده:

- داده‌های تولیدشده توسط مدل: پیام‌های متنی اقناعی (qualitative/textual).
- داده‌های نظرسنجی: پاسخ‌های کمی به سؤالات نگرش (روی مقیاس ۰-۱۰۰)، همراه با متغیرهای پیش‌درمانی مانند ایدئولوژی سیاسی، حزب و دانش سیاسی.
- ویژگی‌های پیام: طول، نسبت نوع-توکن، امتیاز خوانایی Flesch-Kincaid، نسبت زبان اخلاقی/احساسی (از دیکشنری‌های از پیش‌تعریف شده).

• منبع داده:

- پیام‌ها: تولیدشده توسط ۲۴ مدل LLM (۲۲ مدل منبع‌باز مانند Pythia، Qwen-1.5، Claude-3-Opus و GPT-4-Turbo و ۲ مدل بسته مانند Llama-2) مسائل سیاسی از [ISideWith.com](#) (۲۰۲۳) گرفته شده، که از سایت [Tappin et al](#) استخراج مطالعه قبلی شده‌اند.

- داده‌های آموزشی: مجموعه GPT-4-Alpaca (۱۰ هزار مثال برای fine-tuning).

- داده‌های نظرسنجی: جمع‌آوری شده از پلتفرم Prolific (شرکت‌کنندگان آمریکایی، انگلیسی‌زبان، بالای ۱۸ سال). نمونه انسانی پایه از [Tappin et al](#).

• اندازه داده:

- ۷۲۰ پیام تولیدشده (۳۰ پیام برای هر مدل).

- نظرسنجی: ۲۵,۹۸۲ شرکت‌کننده (پس از حذف موارد نامعتبر). توزیع: ۷۵٪ به گروه AI، ۲۰٪ به انسانی، ۵٪ به کنترل. هر پیام حدود ۲۰۰-۳۰۰ شرکت‌کننده دید.

- داده‌های آموزشی: ۱۰ هزار مثال برای هر مدل (۳ دوره آموزشی).

از منظر تحلیلی، داده‌ها ترکیبی از تولید AI و پاسخ انسانی هستند، که اندازه بزرگ نظرسنجی اعتبار آماری را افزایش می‌دهد، اما نمونه غیر نمایندگی (*skewed*) به سمت لیبرال‌ها و زنان) یک محدودیت است.

روش پیشنهادی مقاله

روش مقاله یک آزمایش نظرسنجی تصادفی شده (randomized survey experiment) است که رابطه اندازه مدل و اقناع را اندازه‌گیری می‌کند. به زبان ساده: ابتدا مدل‌های پایه را fine-tune می‌کنند تا دستورات را دنبال کنند، سپس پیام‌های اقناعی تولید می‌کنند، و در نهایت تغییر نگرش افراد را در یک نظرسنجی بزرگ اندازه‌گیری می‌کنند. تحلیل با متانالیز اثرات تصادفی (random-effects meta-analysis) انجام می‌شود تا بازده کاهشی را نشان دهد.

شماتیک ساده (توضیح متنی، مانند فلوچارت):

۱. ورودی: مدل‌های پایه + داده آموزشی → مدل‌های دستورمحور.
۲. تولید: پرامپت + مسئله → پیام اقناعی.
۳. آزمایش: شرکت‌کنندگان → تخصیص تصادفی به گروه‌ها → اندازه‌گیری تغییر نگرش.
۴. تحلیل: اثرات درمان → متانالیز → منحنی بازده کاهشی (log-logistic بهترین فیت).

این روش ساده‌اما قدرتمند است و از ابزارهای آماری مانند رگرسیون حداقل مربعات و متانالیز برای اعتبار استفاده می‌کند.

نتایج اصلی، محدودیت‌ها و ایده‌های ادامه

• نتایج اصلی:

- اقناع مدل‌ها با \log (پارامترها) رابطه مثبت اما کاهشی دارد (۱.۲۶ درصد افزایش برای هر واحد \log). مدل‌های مرزی (مانند Claude-3-Opus) تنها کمی بهتر از مدل‌های ۱۳-۷ میلیاردی هستند.
- بهترین تابع: $\log\text{-logistic}$ ، که پیش‌بینی می‌کند مدل‌های آینده (۳ تریلیون پارامتر) تنها < ۱ درصد بهتر باشند.
- مدیاتور: تکمیل وظیفه (انسجام و مرتبطبودن) توضیح‌دهنده اصلی مزیت مدل‌های بزرگ است؛ پس از تنظیم برای آن، اندازه مدل بی‌معنی می‌شود. مدل‌های مرزی در این متریک سقف زده‌اند (امتیاز ۳/۳).
- مقایسه با انسانی: هیچ مدلی به طور معنادار بهتر از انسان نیست.

• محدودیت‌ها:

- نمونه نظرسنجی غیر نمایندگی (*skewed* به لیبرال‌ها، دموکرات‌ها، زنان)، که ممکن است اثرات را بیش‌براورد کند.
- تمرکز بر پیام‌های تک‌پیامی بدون شخصی‌سازی یا تعامل چنددوره‌ای؛ ممکن است مدل‌های بزرگ در سناریوهای پیچیده‌تر بهتر عمل کنند.
- عدم بهینه‌سازی مدل‌ها برای اقناع (فقط fine-tuning عمومی)؛ ممکن است مرز اقناع بالاتر باشد.
- اندازه دقیق مدل‌های بسته ناشناخته، که بر منحنی‌ها تأثیر می‌گذارد.

ایدههای جدید در این حوزه

- بررسی تعامل چنددوره‌ای و شخصی‌سازی: آیا مدل‌های بزرگ در گفتگوی طولانی یا هدفمند بهتر عمل می‌کنند؟
- بهینه‌سازی خاص اقناع: fine-tuning یا پرامپتینگ پیشرفته برای افزایش اقناع.
- آزمایش در زمینه‌های واقعی: مانند انتخابات واقعی یا رسانه‌های اجتماعی.
- توابع scaling دیگر: تمایز بیشتر بین power law و log-logistic با داده‌های بیشتر.
- ایمنی AI: توسعه ابزارهایی برای کاهش اقناع مضر، مانند refusal در مسائل حساس.

این یافته‌ها نشان می‌دهد که scaling تنها راه افزایش اقناع نیست و تمرکز بر ایمنی و کاربردهای مثبت ضروری است.