

به نام خدا



دانشگاه تهران
دانشکدگان فنی
دانشکده مهندسی برق و کامپیوتر



درس سیستم های چند عاملی پیشرفته پروژه ای درس مرحله ای ۱

نام و نام خانوادگی: ارشیا مخلص

شماره دانشجویی: 810104247

آذر ماه ۱۴۰۴

فهرست

سوال 1.	3	سوال 1.....
مقدمه	3
LLM-based Agents	3
همکاری چندعاملی(Multi-Agent Collaboration)	3
چارچوب ها و پلتفرم ها(Frameworks & Benchmarks)	3
سامانه های کاربردی دامنه محور	4
شبیه سازی اجتماعی و جوامع عامل ها	4
ارزیابی و بهبود مدل	4
جمع بندی و مسیر های آینده	4
سوال 2.	6	سوال 2.....
هدف کلی مقاله چیست؟	6
عامل ها چه هستند و محیط بازی چگونه است؟	6
عامل ها چگونه با هم ارتباط برقرار می کنند؟	6
آیا عامل ها با هم همکاری می کنند؟ چگونه و چرا؟	6
منظور از Theory of Mind در این مقاله چیست؟	6
حافظه در این مقاله چگونه مدیریت و نگهداری می شود؟	7
نتایج اصلی مقاله چیست؟	7
یک محیط/بازی مشابه که بتوان آزمایش را تکرار کرد چیست؟	7
دو شکل مقاله چه چیزی را نشان می دهد؟	8
سوال 3.	9	سوال 3.....
هدف کلی این مقاله چیست؟	9
عامل ها چه هستند و محیط بازی چگونه است؟	9
عامل ها چگونه با هم ارتباط برقرار می کنند؟	9
آیا عامل ها با هم همکاری می کنند؟ چگونه و چرا؟	9
منظور از Theory of Mind در این مقاله چیست؟	10
حافظه در این مقاله چگونه مدیریت و نگهداری می شود؟	10
نتایج اصلی مقاله چیست؟	10
یک محیط/بازی مشابه که بتوان آزمایش را در آن تکرار کرد چیست؟	11
دو شکل مقاله چه چیزی را نشان می دهد؟	11

سوال 1.

در سال های اخیر، عامل های مبتنی بر LLM از یک «چت بات قوی» به اجزای هوشمندی تبدیل شده اند که می توانند مثل یک عضو تیم فکر کنند، برنامه بروزرساند، ابزار استفاده کنند و در کنار عامل های دیگر کار جمعی انجام دهند. مخزن هایی مثل Awesome Multi-Agent Papers و Multi-Agent-Papers تصویر روشنی می دهند از این که جامعه تحقیقاتی چطور از عامل تکی شروع کرده و به طراحی جوامع بزرگ عامل ها، چارچوب های چندعاملی و بنچمارک های پیچیده برای سنجش همکاری، رقابت و هماهنگی رسیده است.

مقدمه

در سطح پایه، LLM-based agent یعنی مدلی که فقط متن تولید نمی کند، بلکه در یک «پوسته عامل» قرار می گیرد و می تواند اهداف را بفهمد، کار را به گام های کوچک بشکند، با محیط و ابزارها تعامل کند و نتایج را ارزیابی کند. هرچه وظیفه پیچیده تر و بلندمدت تر می شود، یک عامل تنها کمتر کافی است؛ به همین دلیل کارها به سمت چندعامل رفته اند تا مثل یک تیم انسانی، نقش ها و مسئولیت ها تقسیم شوند و از تعامل، نقد متقابل و تخصصی سازی سود ببرند.

LLM-BASED AGENTS پایه

در لایه‌ی پایه، کارهایی مثل ReAct نشان می دهند چگونه می توان استدلال زنجیره‌ای و عمل (مثلًاً وب جستجو، اجرای کد، فراخوانی API) را در یک حلقه ادغام کرد تا عامل بتواند در محیطی مثل ALFWorld یا WebShop یا کار انجام دهد. فهرست هایی مثل Awesome-Agent-Papers و مخازن LLM_MultiAgents_Survey_Papers«» مجموعه‌ای بزرگ از این عامل‌های تک نقشی را مستند می کند که واحد اصلی بسیاری از سیستم‌های چندعاملی بعدی هستند.

همکاری چندعاملی (MULTI-AGENT COLLABORATION)

در سطح بعد، کارها به همکاری چندعامل می‌رسند؛ جایی که چند LLM یا چند شخصیت از یک LLM در قالب نقش های متفاوت کنار هم قرار می‌گیرند. نمونه‌های شاخص در README اولیه شامل MetaAgents برای هماهنگی اجتماعی و تشکیل تیم در سناریوی نمایشگاه کار، GameGPT برای توسعه بازی با نقش‌هایی مثل مدیر، مهندس و بازیگران، Hanabi برای ارزیابی توان هماهنگی و Theory of Mind در بازی‌هایی مثل LLM-Coordination Overcooked و Co-NavGPT، AutoGen Multi-Agent Collaboration، AutoAgents CAMEL، AutoGen و الگوهای عمومی تری از گفتگو، نقد متقابل، self-play و تقسیم کار میان عامل‌ها ارائه می‌دهند.

چارچوب‌ها و پلتفرم‌ها (FRAMEWORKS & BENCHMARKS)

برای اینکه این ایده‌ها فقط در مثال‌های پراکنده نمانند، مجموعه‌ای از چارچوب‌ها و بنچمارک‌ها طراحی شده اند: AutoGen و AgentScope و MALLM بستری‌های عمومی برای تعریف عامل‌ها، نقش‌ها و جریان مکالمه و ابزارها فراهم می‌کنند، در حالی که MultiAgentBench، AgentBench، Collab-Overcooked و SmartPlay روی سنجش توان عامل‌ها در همکاری، رقابت، برنامه‌ریزی و تعامل با محیط تمرکز دارند. این

بنچمارک ها معمولاً سناریوهای بازی، حل مسئله، QA، استفاده از ابزار و ارزیابی امنیت/ریسک را در بر می‌گیرند و معیارهایی مثل موفقیت در کار، پایداری همکاری و کیفیت استدلال را می‌سنجد.

سامانه های کاربردی دامنه محور

بخش بزرگی از لیست به کاربردهای دامنه ای اختصاص دارد؛ در نرم افزار و کدنویسی، کارهایی مثل ChatDev، AutoSafeCoder و CodeR، MetaGPT نیازمندی، طراحی، پیاده سازی و دیباگ کنار هم قرار می‌دهند. در پژوهشی، Agent Hospital، MDAgents و Zodiact، کارهای مشابه از عامل‌ها برای تصمیم‌سازی تشخیصی، آموزش پزشکی و شبیه سازی بیمارستان استفاده می‌کنند. در رباتیک و ناویری، RoCo و Co-NavGPT و MALMM با تکیه بر LLM/VLM به عنوان برنامه ریز سطح بالا، چند ربات را برای ناویری و دستکاری هماهنگ می‌کنند. حوزه‌های دیگر مثل سیستم‌های توسعه‌گر (AgentCF)، مالی (TwinMarket)، حقوق و دادگاه‌های شبیه سازی شده (AgentCourt) و برنامه ریزی شهری (urban planning agents) هم نشان می‌دهند که الگوی چندعاملی تقریباً به هر حوزه‌ای که تصمیم‌گیری و تعامل پیچیده دارد قابل تعمیم است.

شبیه سازی اجتماعی و جوامع عامل‌ها

گروه دیگری از کارها هدفشان حل یک وظیفه مهندسی نیست، بلکه مطالعه رفتار جمعی و اجتماعی عامل‌ها است؛ TwinMarket، OASIS، AgentSociety، Generative Agents و ده‌ها کار مشابه، شهرها، بازارها یا جوامع بزرگ مصنوعی می‌سازند که در آن صدھا تا میلیون‌ها عامل LLM زندگی روزمره، تصمیم‌اقتصادی، کنش سیاسی یا تعامل اجتماعی را شبیه سازی می‌کنند. این خط پژوهشی برای مطالعه شکل‌گیری هنجارها، ظهور همکاری پایدار، تکامل فرهنگی و حتی هم ترازی اجتماعی مدل‌ها اهمیت دارد و نشان می‌دهد LLM‌ها فراتر از حل مسائل تکی، می‌توانند واحدهای سازنده «جامعه‌های مصنوعی» باشند.

ارزیابی و بهبود مدل

برای نگه داشتن این سیستم‌ها در مسیر درست، بخش بزرگی از کارها روی ارزیابی و بهبود مرکز دارد؛ رویکردهایی مثل AutoDefense و ChatEval، Agent-as-a-Judge، multi-agent debate و کنند از خود عامل‌ها به عنوان منتقد، داور یا مهاجم/دفاع کننده استفاده کنند تا هم کیفیت استدلال و factuality بالا برود و هم ریسک‌هایی مثل jailbreak یا تصمیم غلط شناسایی شود. بنچمارک‌هایی مثل AgentBench و Decrypto، Collab-Overcooked این سامانه‌ها تا چه حد می‌توانند نیات دیگران را درک کنند، برنامه مشترک بسازند و در محیط‌های پر ریسک قابل اعتماد عمل کنند.

جمع بندی و مسیرهای آینده

در مجموع، این README‌ها و فهرست‌ها نشان می‌دهند که LLM-based multi-agent systems از «چند مدل حرف زن» فراتر رفته‌اند و در حال تبدیل شدن به معماری‌های جدی برای حل مسائل پیچیده، شبیه سازی جوامع و خودکارسازی جریان‌های کاری هستند، اما هنوز چالش‌های مهمی مثل مقیاس پذیری، پایداری همکاری، Theory of Mind قابل اعتماد، ایمنی و ارزیابی دقیق باز هستند. همین چالش‌ها نقطه شروع خوبی برای پروژه‌های جدید است: از طراحی پروتکل‌های بهتر برای ارتباط عامل‌ها، تا ساخت بنچمارک‌های دقیق‌تر برای هماهنگی و طراحی چارچوب‌هایی که بتوانند این سیستم‌ها را در کاربردهای واقعی امن و پایدار نگه دارند.

سوال 2.

هدف کلی مقاله چیست؟

هدف کلی این مقاله ارزیابی توانایی عامل های مبتنی بر LLM برای همکاری چندعاملی در یک مأموریت متغیر جستجو و خنثی سازی بمب است، به خصوص از نظر برنامه ریزی مشترک و استنتاج Theory of Mind (استدلال درباره وضعیت ذهنی دیگر عامل ها). مقاله می خواهد نشان دهد که این عامل ها بدون آموزش ویژه و به صورت صفر-شات می توانند رفتاری نزدیک به روش های پیشرفته MARL و برنامه ریزی داشته باشند، اما هنوز به خاطر خطای مدیریت کانتکست بلندمدت و هالوسینیشن از وضعیت محیط، برنامه ریز یا هم تیمی ایده آلی نیستند.

عامل ها چه هستند و محیط بازی چگونه است؟

در این تحقیق سه عامل LLM محور با نام های Charlie و Bravo و Alpha تعریف شده اند که هر کدام نقش یک متخصص در تیم خنثی سازی بمب را بازی می کنند و با استفاده از یک مدل زبانی) مثل GPT-3.5 یا GPT-4) یا تصمیم می گیرند. محیط بازی یک دنبیای متغیر است که روی یک گراف از اتاق ها (نودها) ساخته شده؛ در هر اتاق ممکن است بمب های چندمرحله ای رنگی وجود داشته باشد و هر عامل فقط اتاق فعلی، ابزارهای خودش، پیام های تیم و وضعیت کلی امتیاز را می بیند و می تواند بین حرکت به اتاق های مجاور، بازرسی بمب و استفاده از سیم برهای رنگی انتخاب کند تا بمب ها را به طور هماهنگ خنثی کنند.

عامل ها چگونه با هم ارتباط برقرار می کنند؟

ارتباط بین عامل ها کاملاً متغیر است: در هر راند، هر عامل علاوه بر انتخاب اکشن، یک پیام متغیر برای تیم می فرستد که بلافاصله از طریق کانال چت داخلی به صورت broadcast در مشاهدات دور بعد همه اعضاء نمایش داده می شود. خود محیط بازی و یک لایه ای «text game interface» این پیام ها را از متن طبیعی به اعمال انتزاعی) مثلاً «X» «Y» «inspect bomb» یا «move to room» رمزگشایی می کند و در عوض، وضعیت محیط را دوباره به صورت توضیح متغیر (محتوای اتاق، نتیجه ای عمل، موقعیت هم تیمی ها، پیام های قبلی) به عامل ها بر می گرداند.

آیا عامل ها با هم همکاری می کنند؟ چگونه و چرا؟

بله، همکاری صریح و ضمنی هر دو شکل می گیرد: عامل ها وظایف را بین خود تقسیم می کنند، یک عامل) معمولاً نقش رهبر را می گیرد و به بقیه می گوید به کدام اتاق بروند یا کدام بمب را بررسی/خنثی کنند و دیگران هم غالباً طبق این پیشنهادها عمل می کنند. این همکاری به این دلیل ضروری است که اطلاعات و ابزارها بین عامل ها پخش شده (هر کس تعدادی سیم بر رنگی دارد و فقط اتاق فعلی را می بیند)، بنابراین برای رسیدن به حداکثر امتیاز باید هم در فضا (پخش شدن در اتاق های مختلف) و هم در توالی برش های رنگی بمب ها با هم هماهنگ شوند؛ نتایج نشان می دهد تیم های مبتنی بر GPT-4 با این همکاری می توانند همه بمب ها را در زمان محدود خنثی کنند.

منظور از THEORY OF MIND در این مقاله چیست؟

Theory of Mind در این مقاله یعنی توانایی عامل برای استدلال درباره «آن چه خود و دیگران می دانند یا نمی دانند»، در سه سطح:

ا توصیف وضعیت دانشی و باوری خود عامل.
Introspection:

استنتاج این که یک هم تیمی در مورد وضعیت جهان چه باوری دارد) مثلاً آیا Bravo می داند First-order ToM: بمب در اتاق ۳ است یا نه.)

استنتاج این که هم تیمی درباره وضعیت ذهنی خود عامل چه فکری می کند) مثلاً آیا Charlie می داند که Alpha از محتوای اتاق ۵ خبر دارد.)

در طول مأموریت، بعد از بعضی اکشن ها از عامل ها سؤالات ToM پرسیده می شود) الهام گرفته از تست Sally و پاسخ های متى آن ها توسط انسان ها با توجه به تاریخچه کامل تعاملات برچسب گذاری می شود.) Anne)

حافظه در این مقاله چگونه مدیریت و نگهداری می شود؟

حافظه ای عامل ها از دو بخش اصلی تشکیل شده است:

History در ورودی مدل: بخشی از تاریخچه ای راندهای قبلی (قوانین بازی، چند راند اخیر شامل مشاهدات، اعمال و پیام ها) تا حد ظرفیت ورودی مدل در prompt نگه داشته می شود.

Belief State صریح) در حالت GPT-4 + Belief: بعد از هر مشاهده، عامل با پرامپت خاصی وادر می شود یک متن ساخت یافته به روز کند که باور های مهمش درباره ای محیط) مثلاً «Bomb 1 is in Room 0 with sequence Red», «Room 3 connected to Room 0») بعد تکرار می شود تا محدودیت طول context جبران شود و خطاهای ناشی از فراموشی و هالوسینیشن کاهش یابد.

نتایج اصلی مقاله چیست؟

از نظر عملکرد وظیفه: تیم مبتنی بر ChatGPT (gpt-3.5) همه بمب ها را خنثی نمی کند و میانگین امتیاز حدود ۴۳ می گیرد، در حالی که تیم ۴ GPT-4+Belief GPT-4+Belief امتیاز کامل ۹۰ می گیرند؛ GPT-4+Belief تعداد راندها را نسبت به GPT-4 عادی تقریباً نصف می کند ($28/3 \rightarrow 12/3$) و درصد اکشن معنبر را از حدود ۷۲% به ۸۶% می رساند، که آن را از نظر کارایی به MAPPO نزدیک می کند.

نسبت به خط مبنای CBS Planner: (برنامه ریز مرکز با اطلاعات کامل) در ۶ راند به جواب بهینه می رسد و MAPPO پس از میلیون ها گام آموزش در حدود ۱۱ راند، در حالی که عامل های LLM بدون آموزش خاص و با تصمیم گیری کاملاً غیرمرکز به عملکردی نسبتاً نزدیک می رسد.

از نظر ToM: دقت ChatGPT در introspection حدود ۷۹٪، برای first-order حدود ۴۲٪ و برای second-order حدود ۱۲٪ است؛ در GPT-4 این اعداد به حدود ۸۰٪، ۶۰٪ و ۶۴٪ می رسد و در GPT-4+Belief به حدود ۹۷٪، ۸۰٪ و ۶۹٪ افزایش می یابند، که نشان می دهد نمایش صریح باورها هم عملکرد تیمی و هم استنتاج ToM را بهبود می دهد.

از نظر کیفی: الگوهایی مثل رهبری emergent Alpha (دستور می دهد، بقیه تبعیت می کنند)، کمک به هم تیمی، حل تعارض و اشتراک اطلاعات در رفتار تیم های GPT-4 دیده می شود، ولی دو نوع شکست سیستماتیک هم مشاهده می شود: نادیده گرفتن محدودیت های بلندمدت (مثلاً حرکت به اتاق غیرمجاور) و هالوسینیشن درباره وضعیت بازی (مثلاً فرض وجود بمبی که نیست یا مرحله ای که دیده نشده).

یک محیط/بازی مشابه که بتوان آزمایش را تکرار کرد چیست؟

یک محیط مشابه می تواند هر سناریوی جست وجو و نجات متى یا شبیه سازی شده ای باشد که در آن: جهان به اتاق ها/نواحی گراف مانند تقسیم شده،

اطلاعات و ابزار بین عامل‌ها توزیع شده،

موفقیت نیاز به توالی درست اعمال و هماهنگی بین اعضای تیم دارد.

برای مثال می‌توان یک «ماموریت نجات گروگان» متنی طراحی کرد که در آن چند عامل باید در یک ساختمان، اتفاق‌ها را جست و جو کنند، کلیدها و ابزارها را تقسیم کنند، درها را باز کنند و گروگان‌ها را نجات دهند؛ هر عامل فقط اتفاق فعلی خود را می‌بیند و از طریق چت متنی با بقیه هماهنگ می‌شود، دقیقاً مشابه ساختار گرافی و ارتباطی این مقاله. یک گزینه دیگر بسط همین سناریو بمب خنثی کنی به تعداد بیشتر اتفاق‌ها و بمب‌ها است، یا پیاده‌سازی آن در یک محیط متنی مانند ALFWorld یا TextWorld با قوانین مشابه.

دو شکل مقاله چه چیزی را نشان می‌دهند؟

شکل اصلی فریم ورک (Figure 1) این شکل معماری کلی سیستم را نشان می‌دهد؛ در سمت راست «Task Environment» (محیط بمب‌ها و اتفاق‌ها) قرار دارد، که از طریق «Text Game Interface» به توصیف‌های متنی و بازخورد درباره اعمال تبدیل می‌شود؛ در سمت چپ سه عامل (Alpha، Bravo، LLM Charlie) دیده می‌شوند که هر کدام با مدل زبانی belief state متنی و تاریخچه‌ی تعامل خود کار می‌کنند و در هر راند بر اساس observation، belief state را به روزرسانی کرده، اکشن و پیام به تیم تولید می‌کنند. این شکل جریان حلقه‌ای: محیط → مشاهده متنی → به روزرسانی باور → انتخاب عمل و پیام → رمزگشایی به عمل → به روزرسانی محیط را به صورت شماتیک نمایش می‌دهد.

شکل رفتار و خطاهای (Figure 2) این شکل چند پنل مثالی از لاگ تعامل را نشان می‌دهد؛ در پنل بالا-چپ نمونه ای از شکست در مدیریت کانتکست بلندمدت (درخواست حرکت به اتفاق غیرمجاور و خطای محیط) آمده است، در بالا-راست نمونه ای از رفتارهای تعاملی مثبت مثل رهبری emergent و تقسیم وظایف در پیام‌های Alpha و Bravo و Charlie را می‌بینیم، و در پنل‌های پایین، پاسخ‌های GPT-4+Belief و ChatGPT به پرسش‌های GPT-4+Belief (second-order، first-order، introspection) کنار هم گذاشته شده تا نشان دهد ToM چطور با توجه به تاریخچه‌ی موقعیت و پیام‌ها، استنتاج دقیق‌تری از «چه کسی چه چیزی می‌داند» انجام می‌دهد.

سوال 3.

هدف کلی این مقاله چیست؟

هدف کلی مقاله این است که نشان دهد می توان با استفاده از LLM های از پیش آموزش دیده و چند مثال زنجیره استدلال (chain-of-thought)، آن ها را به عامل هایی تبدیل کرد که در بازی های استراتژیک بتوانند به صورت انعطاف پذیر و قابل تعمیم، درباره ای وضعیت ها، پاداش ها و باورهای دیگران استدلال کنند، بدون این که نیاز به آموزش مجدد روی هر بازی جدید باشد. نویسنگان می خواهند نشان دهند که با یک «prompt compiler» قادر است در انواع بازی های ماتریسی و مذاکره، استراتژی های نزدیک به نظریه ای بازی و انسان مانند تولید کند.

عامل ها چه هستند و محیط بازی چگونه است؟

در این کار عامل ها همان LLM هایی هستند) مثل code-davinci-002 و نسخه های متغیر davinci/curie که نقش بازیکن را در بازی های استراتژیک بر عهده می گیرند و با تکیه بر پرامپت های حاوی دموهای استدلال، حرکت بعدی خود را انتخاب می کنند. دو نوع محیط اصلی تعریف می شود:

بازی های ماتریسی (Matrix Games): بازی های 2×2 و بزرگ تر (3×3 ، 3×4 ، چندبازیکنه) مثل Prisoner's Dilemma، Chicken، Stag Hunt، Matching Pennies آن ها هر بازیکن یک اکشن از مجموعه اکشن های محدود انتخاب می کند و ترکیب اکشن ها یک پاداش عددی برای هر بازیکن تولید می کند.

بازی های مذاکره (Deal or No Deal): سناریوهایی که دو بازیکن باید یک «گلدان» شامل چند نوع آیتم (کلاه، کتاب، توب) را بین خود تقسیم کنند، در حالی که هر کدام ارزش متفاوتی برای هر آیتم دارند و از ارزش طرف مقابل خبر ندارند؛ در برخی آزمایش ها خود مدل نقش «broker» (داور) را دارد و باید یک تقسیم «منصفانه» پیشنهاد دهد.

عامل ها چگونه با هم ارتباط برقرار می کنند؟

در بیشتر آزمایش ها، ارتباط بین عامل ها ضمنی است و از طریق ساختار بازی و مدل سازی ذهنی طرف مقابل انجام می شود؛ یعنی LLM در استدلال متغیر خود در نظر می گیرد که «حریف چه هدفی دارد و چه اکشنی خواهد زد» و بر همان اساس بهترین پاسخ را انتخاب می کند (level-k / Theory-of-Mind سبک). در محیط مذاکره، ارتباط به صورت دیالوگ رفت و برگشتی است: انسان یا عامل مقابل یک پیشنهاد متغیر برای تقسیم آیتم ها می دهد، LLM ای استدلال متغیر درباره ارزش ها و باورها پاسخ می دهد (accept/reject/propose)، و این رفت و برگشت در چند دور ادامه می یابد؛ در سناریوی «broker» هم LLM براساس اطلاعات دو طرف، پیشنهاد متن محور جدید تولید می کند.

آیا عامل ها با هم همکاری می کنند؟ چگونه و چرا؟

این مقاله هم سناریوهای رقابتی و هم تعاقنی را پوشش می دهد، اما در هر دو حالت نوعی «استدلال متقابل» بین عامل ها وجود دارد که به همکاری یا سازش منجر می شود. در برخی بازی های ماتریسی مثل «welfare» یا وقتی هدف «کمینه کردن تفاوت پاداش ها (equality) است، LLM عملاً به دنبال حداکثرسازی رفاه

مشترک یا تقسیم منصفانه پاداش است، که نوعی همکاری است. در بازی مذاکره‌ی واقعی هم نتایج نشان می‌دهد که عامل نه تنها پاداش خود را بالا می‌برد، بلکه پاداش انسان را هم کمی بهتر از baseline می‌کند و در ارزیابی کاربران «منطقی، سازش‌گر و شبیه انسان» دیده می‌شود، یعنی به طور ضمنی همکاری/سازش را یاد می‌گیرد تا به توافق برسد.

منظور از THEORY OF MIND در این مقاله چیست؟

هر چند مقاله مستقیماً روی واژه‌ی «Theory of Mind» تمرکز ندارد، ولی مفهوم آن را در قالب سه مؤلفه‌ی استراتژیک پیاده می‌کند:

Search: جست‌وجو در درخت بازی با فرض‌هایی درباره‌ی این که حریف چه هدفی دارد و چه اکشنی انتخاب می‌کند.

Value Assignment: این که مدل بفهمد خود بازیکن و حریف چه پاداشی از هر حالت می‌گیرند (مثلاً حریف می‌خواهد «max» خودش را زیاد کند یا «کمک» کند).

Belief Tracking: نگه داشتن باور روی اطلاعات پنهان مثل حالت جهان یا ترجیحات حریف، مثلاً از روی این که «حریف ۹ بار از ۱۰ بار کتاب را به توب ترجیح داده، پس کتاب را بیشتر دوست دارد.» در عمل، این یعنی مدل تلاش می‌کند باورهای حریف و اهدافش را از روی رفتار و پاداش‌ها استنتاج کند و بر اساس آن استراتژی خود را تنظیم کند، که جوهره‌ی Theory of Mind در محیط‌های استراتژیک است.

حافظه در این مقاله چگونه مدیریت و نگهداری می‌شود؟

حافظه در این کار بیشتر در قالب «استدلال متنی درون پرامپت» است:

یک «prompt compiler» برای هر بازی یک دمو تولید می‌کند که در آن مراحل جست‌وجو، محاسبه‌ی امیدریاضی پاداش‌ها، مقایسه‌ی اکشن‌ها و به روزکردن باورها روی حالت‌های پنهان یا ارزش‌ها به صورت متن chain-of-thought نوشته شده است.

هنگام حل بازی جدید، LLM همین الگوی متنی را روی داده‌های جدید تکرار می‌کند: برای هر حالت/اکشن پاداش‌ها را حساب می‌کند، روی hidden states یا ترجیحات حریف belief به روز می‌کند و در چند دور propose → evaluate → revise» (به خصوص در مذاکره) پیشنهاد را اصلاح می‌کند؛ این متن استدلالی کوتاه مدت حافظه‌ی کاری عامل است.

در تنظیمات پیچیده‌تر، از «factored reasoning» و ابزارهایی مثل (...)(...) mean و (...)(...) search به صورت نمادین در متن استفاده می‌شود تا استدلال روی زیرمسئله‌ها و زیر درخت‌ها شکسته شود و مشکل محدودیت context کمتر شود.

نتایج اصلی مقاله چیست؟

در بازی‌های ماتریسی با پاداش‌های جدید، LLM‌ها پرامپت استراتژیک توانست در تقریباً همه‌ی موارد بهترین پاسخ نظریه‌ی بازی را پیدا کند (دقیقاً ۱/۰۰ دقیقه ۰-shot، ۲-shot و ۴-shot معمولی) و ساده عملکرد خیلی ضعیف تری داشتند.

در بازی‌هایی با قوانین پیچیده‌تر (درخت‌های دو مرحله‌ای، ماتریس‌های بزرگ‌تر، چند بازیکن)، استفاده از «factored search» باعث شد مدل بتواند به این ساختارهای جدید بدون مثال اضافی تعمیم بدهد و با دقت بالا

بهترین پاسخ را انتخاب کند، در حالی که بدون فاکتور کردن، دقت به خاطر محدودیت context و پیچیدگی افت می‌کرد.

در سناریوهای با اهداف جدید (مثلًا daxity، welfare، max/help)، مدل فقط با دیدن دموهایی که هدف را توضیح می‌دهند و بدون آموزش مجدد توانست در اکثر موارد با دقت نزدیک ۱۰۰٪ بهترین پاسخ را برای هدف جدید انتخاب کند، که نشان دهنده‌ی تعمیم به «اهداف انزواعی» است.

در نقش «Deal or No Deal broker»، عامل استدلالی توانست تفاوت با تقسیم ایده‌آل عادلانه را بیش از دو برابر نسبت به baseline کاهش دهد (یعنی پیشنهادات بسیار نزدیک به optimum equality یا Rawlsian fairness تولید کند).

در نقش مذاکره کننده‌ی مستقیم با انسان‌ها، روش پیشنهادی هم پاداش خود و هم پاداش انسان را کمی بالا برد و در ارزیابی کاربران، از نظر «شبیه انسان بودن، معقول بودن، سازش‌گر و کمتر تهاجمی بودن» از baselines نمره‌گرفت، یعنی هم از نظر کارایی و هم از نظر کیفیت اجتماعی رفتار، برتر بود.

یک محیط/بازی مشابه که بتوان آزمایش را در آن تکرار کرد چیست؟

هر محیط نظریه‌ی بازی که با چند بازیکن، جدول پاداش و شاید اطلاعات پنهان تعریف می‌شود، مناسب تکرار این سبک آزمایش است؛ مثلًا:

نسخه‌ای از بازی «Diplomacy» یا «Hanabi» که به جای اجرای کامل، موقعیت‌ها به صورت توصیف متنی payoff‌های ساده شده به مدل داده شود تا بهترین حرکت یا مذاکره را انتخاب کند.

یک بازار ساده دو یا سه کالایی که چند عامل باید قیمت گذاری یا پیشنهاد خرید/فروش بدنه؛ LLM می‌تواند با همان ساختار search / value / belief استراتژی قیمت یا پیشنهاد را پیدا کند.

سناریوهای ساده‌ی مزایده (auction) که در آن چند bidder با ارزش‌های پنهان در یک حراج تک شیء یا چندشیء شرکت می‌کنند و مدل باید استراتژی پیشنهاد را انتخاب کند.

نکته‌ی اصلی این است که بتوان محیط را به صورت «درخت یا ماتریس پاداش» و چند متن دمو برای جست‌جو و محاسبه‌ی ارزش‌ها به LLM داد تا همان روش روی آن اعمال شود.

دو شکل مقاله چه چیزی را نشان می‌دهند؟

شکل ۱: در سمت چپ یک درخت بازی جزئی مشاهده‌ای با دو بازیکن Bob و Gopher و حالت‌های پنهان s_1, s_2 و اکشن‌های a_1, a_2, b_1, b_2 نشان داده شده که در برگ‌های آن پاداش‌های دو بازیکن نوشته شده است؛ در سمت راست، یک نمونه از استدلال متنی مرتبط آورده شده که در آن با رنگ‌های مختلف سه جزء (search «جست‌جوی شاخه‌ها»، value «محاسبه‌ی امیدریاضی پاداش‌ها» و belief «belief دادن باور روی حالت جهان از روی اقدام حریف») هایلایت شده تا نشان دهد prompt compiler چگونه chain-of thought را ساختار می‌دهد.

شکل ۳ (برای Deal or No Deal): پنل (a) یک اپیزود نمونه را نشان می‌دهد که در آن تعداد آیتم‌ها و ارزش‌های Alice و Bob نمایش داده شده و پیشنهاد اولیه حریف مشخص است؛ پنل (b) نمونه‌ای از استدلال متنی مربوط به همین اپیزود را نشان می‌دهد که در آن با رنگ صورتی بخشی که روی «ارزش‌ها» حساب می‌کند، با نارنجی «جست‌جوی پیشنهادهای جایگزین» و با خاکستری «به روزرسانی باورها درباره ترجیحات حریف» مشخص شده است. این شکل‌ها کمک می‌کنند بصری بینیم که سه مؤلفه‌ی search / value / belief چگونه داخل متن زنجیره استدلالی جاگذاری شده‌اند.