

۱) مسئله اصلی چیست و چرا مهم است؟

این مقاله روی تشخیص شایعات/اطلاعات نادرست حوزه سلامت در فضای وب (به‌طور خاص در زبان چینی) تمرکز دارد. نویسنده‌گان می‌گویند با رشد سرعت انتشار محتواهای حوزه سلامت در اینترنت، «شایعات سلامت» در کنار اطلاعات درست پخش می‌شوند و می‌توانند تهدیدی برای سلامت عمومی باشند. مشکل بزرگ‌تر این است که (۱) اکثر دیتاست‌های موجود برای شایعات سلامت کوچک هستند و تنوع کافی برای آموزش مدل‌های عمیق را ندارند. و (۲) روش‌های سنتی تشخیص شایعه معمولاً فقط برچسب (شایعه/غیرشایعه) را خروجی می‌دهند، بدون اینکه دلیل یا استناد علمی ارائه کنند. برای کاربران نهایی، صرفاً دانستن اینکه یک خبر شایعه است کافی نیست؛ آن‌ها نیاز به دلایل و شواهد علمی دارند تا قانع شوند.

مقاله برای پرکردن این شکاف دو کار اصلی انجام می‌دهد:

- ساخت یک دیتاست بسیار بزرگ به نام HealthRCN با ۱.۱۲ میلیون نمونه داده شایعه/اطلاعات سلامت.
- ارائه یک سیستم به نام HRDE که یک LLM تقویت‌شده با بازیابی (RAG) است تا هم برچسب شایعه بودن/نبودن بدهد و هم تحلیل همراه با توضیح با استناد به منابع بازیابی شده ارائه کند.

۲) ورودی‌ها و خروجی‌های مدل/سیستم

ورودی سیستم:

- یک متن یا سوال وارد شده توسط کاربر که حاوی ادعایی در مورد سلامت است (مثلًا: "نوشیدن آب اکسیژنه می‌تواند تومورها را درمان کند"). این ورودی در پرامپت «User Input» قرار می‌گیرد.

روجی استاندارد شده:

باید حداقل شامل دو بخش باشد:

[Conclusion] : یکی از این سه برچسب می‌باشد

Rumor / Not rumor / Not related to health information

[Analysis] شامل استدلال و تحلیل درباره همان ورودی است.

اگر سیستم از اسناد بازیابی شده استفاده کند، در متن تحلیل ارجاع‌هایی مثل [1] می‌آورد و در انتهای یک بخش [References] شامل مشخصات منبع اضافه می‌شود (این بخش طبق مقاله در پس‌پردازش تکمیل می‌شود).

(3) داده‌های مورد استفاده (نوع، منبع، اندازه)

1-۳) دیتاست آموزشی: HealthRCN

این دیتاست برای رفع مشکل کمبود داده ساخته شده است.

- **اندازه:** ۱.۱۲ میلیون نمونه داده.
- **منبع اولیه:** (Crawling) پرسش‌های واقعی کاربران از وبسایت‌های مشاوره پزشکی معتبر (مانند [net.39](#)). حدود ۲ میلیون پرسش اولیه جمع‌آوری و پس از پاکسازی به ۱.۱۲ میلیون تقلیل یافت.
- **تولید داده‌های شایعه (Data Augmentation):** از آنجا که سوالات واقعی لزوماً شایعه نیستند، نویسنده‌گان از مدل GPT-3.5-Turbo برای تولید نسخه «شایعه» از روی پاسخ‌های صحیح استفاده کردند.
- **ساختار هر نمونه:** شامل پرسش اصلی، عنوان شایعه، متن شایعه (کوتاه/بلند)، پاسخ رد شایعه (Debunking)، پاسخ صحیح پزشکی و کلیدواژه‌ها.

2-۳) پایگاه دانش مرجع (Knowledge Base for Retrieval)

این پایگاه محلی است که سیستم در زمان پاسخ‌دهی (Inference) در آن به جستجوی سند می‌پردازد.

- **منابع:** ۱۹ وبسایت و سازمان معتبر سلامت (مانند سایت‌های دولتی بهداشت و بیمارستان‌های تخصصی).
- **اندازه :**
 - بخش Elasticsearch (اسناد کامل): ۳۵۸,۱۱۴ سند.
 - بخش Milvus (تکه‌های متن): ۱,۳۴۷,۹۷۶ تکه (chunk).
 - Chunking: به دلیل محدودیت مدل امبدینگ m3e ، متن مقاله‌ها (به ۵۱۲ توکن) تکه‌تکه می‌شود و سپس embedding تولید می‌شود.
- **فرآیند بهروزرسانی:** این پایگاه به صورتی آپدیت می‌شود که اطلاعات جدید پزشکی در دسترس مدل باشد.

3-3) دیتاست ارزیابی (Evaluation Set)

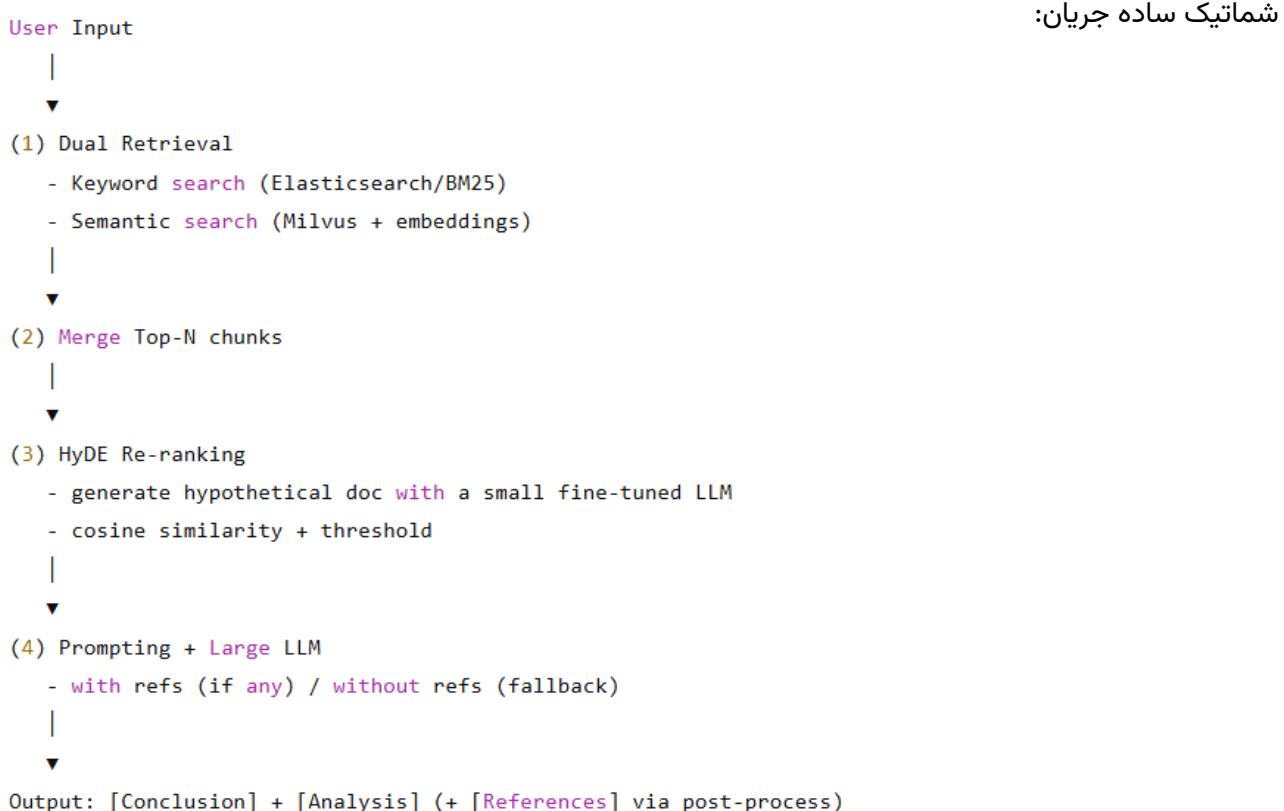
- **اندازه:** 2500 نمونه
- **برچسب‌ها و توزیع:** Not related to health information: 148 ;Not rumor: 855 ;Rumor: 1497
- **ترکیبی از نمونه‌های غیرمرتبه با سلامت (دستی ساخته شده) و بدون همپوشانی با داده فاینتیون.**

۴) روش پیشنهادی مقاله به زبان ساده + شماتیک یا شبکه.

۱-۴) ایده کلی HRDE (چهار بخش اصلی)

HRDE چهار مراحل دارد:

- (۱) جمعآوری و ذخیره اسناد مرجع
- (۲) بازیابی اسناد مرجع
- (۳) رتبه‌بندی دوباره (Re-ranking)
- (۴) تولید پاسخ تشخیص شایعه توسط LLM



۴-۲) بازیابی دوگانه (Dual Retrieval)

- استخراج کلیدواژه‌ها برای بازیابی مبتنی بر متن (BM25) از پایگاه داده Elasticsearch
- ورودی کاربر (توسط m3e) برای بازیابی برداری با embedding Milvus از پایگاه داده
- هدف این است که هم «دقیق کلیدواژه‌ای» و هم «پوشش معنایی» را داشته باشند.

به دلیل اینکه خروجی Elasticsearch مقاله کامل است، chunk می‌شود و با chunk می‌گردد تا Top-N ساخته شود.

در تنظیمات آزمایش، مقاله می‌گوید از Elasticsearch تعداد ۵ سند کامل و از Milvus تعداد ۲۵ chunk بازیابی می‌شود.

(4-3) Re-ranking با HyDE (فیلتر کردن اسناد واقعاً مفید)

برای اینکه تعداد chunk‌ها به ورودی LLM اصلی قابل مدیریت باشد و نویز کم شود:

- یک «سند فرضی» (hypothetical document) با یک مدل کوچکتر (Qwen1.5-4B-Chat) فاین‌تیون شده تولید می‌کنند.

- شباهت معنایی (cosine similarity) بین هر chunk و سند فرضی محاسبه می‌شود و رتبه‌بندی انجام می‌گیرد؛ سپس K-Top انتخاب می‌شود و یک آستانه شباهت هم برای حذف موارد خیلی نامرتب اعمال می‌شود.
- در تنظیمات پایه، Top-K نهایی = ۵ قطعه و آستانه شباهت برابر 0.5 است.

(4-4) تولید پاسخ (Generation) با پaramپت‌های دو حالت

- دو نوع پaramپت طراحی می‌شود: با رفرنس و بدون رفرنس. هر پaramپت ۴ بخش دارد: توضیح کار، ورودی‌ها، الزامات دقیق، و فرمت خروجی.
- اگر هیچ سند معتبری بازیابی نشود یا همه در Re-ranking حذف شوند، سیستم با پaramپت بدون رفرنس پاسخ می‌دهد (fallback روی دانش داخلی مدل).

(5) نتایج اصلی، محدودیت‌ها و ایده‌های ادامه

(5-1) نتایج اصلی (کمی + کیفی)

- معیارهای ارزیابی: F1 و Accuracy برای تشخیص برچسب؛ و سه معیار کیفی پاسخ (Relevance, Reliability, Richness). همچنین نرخ پاسخ معتبر (Valid Answer Rate) هم ثبت می‌شود.
- برای نمره‌دهی معیارهای کیفی، مقاله از GPT-4-1106-Preview برای ارزیابی با پaramپت‌های مشخص استفاده کرده است.

:HRDE عملکرد

- average accuracy: 91.04%
- F1: 91.58%

طبق آزمایش‌های انجام شده، این مدل تخصصی‌شده حدود ۱۲٪ عملکرد بهتری نسبت به GPT-4 (نسخه Preview-1106) در تشخیص شایعات چینی داشته است.

• (RAG و SFT اثر Ablation):

- جدول Ablation نشان می‌دهد ترکیب SFT+RAG بهترین عملکرد را می‌دهد و استفاده تنها از یکی از آن‌ها ایده‌آل نیست (مثلاً RAG بدون فاینتیون می‌تواند مدل را گمراه کند).

• نمونهٔ کیفی (Case Study):

- در یک مثال، سیستم درست تشخیص می‌دهد «اطلاعات شکر و دیابت» شایعه نیست و برای تحلیل، چند سند مرجع می‌آورد و حتی به دیدگاه WHO درباره «قند افزوده و سرطان» اشاره می‌کند تا نگرانی کاربر را کاهش دهد.

• (5-2) محدودیت‌های اعلام شده در مقاله

مقاله دو محدودیت مهم را صریح می‌گوید:

- خطا در موضوعات غیر سلامت: به این دلیل که داده‌های آموزش (Fine-tuning) عمدتاً پزشکی است، مدل گاهی ورودی‌های غیر پزشکی کاربران را به اشتباہ تفسیر می‌کند.
- گرچه مدل‌های زبانی سرعت استنتاج بالایی دارند (مدل ۴۶۹۸ توکن/ثانیه و مدل ۱۴۱ میلیاردی حدود ۱۴۳ توکن/ثانیه پردازش می‌کنند)، اما گلوگاه اصلی سیستم، زمان بر بودن فرایند بازیابی اسناد (Retrieval) است که کندی کلی سیستم را ایجاد می‌کند.

برنامه آینده نویسندگان: بهبود دیتاست دستورالعملی و بهینه‌سازی فرایند RAG برای ارتقای عملکرد کلی.

• (5-3) ایده‌های ادامه (تحلیلی و پیشنهادی بر پایه همین چارچوب)

با توجه به محدودیت‌های خود مقاله و طراحی HRDE، چند مسیر توسعه منطقی است:

- بالانس‌کردن کلاس "Not related to health information": چون خود نویسندگان می‌گویند این نوع نمونه‌ها در آموزش کم بوده، اضافه کردن نمونه‌های متنوع غیرسلامت (و حتی ورودی‌های مرزی مثل تغذیه عمومی/ورزش/سبک زندگی) می‌تواند خطای تفسیر را کم کند.
- کاهش زمان بازیابی: چون سرعت پاسخ‌دهی تحت تأثیر retrieval است، می‌توان روی کاهش تعداد candidates، کش‌کردن نتایج پر تکرار، یا ساده‌سازی rerank در بعضی سناریوهای کار کرد (مقاله خودش «زمان بر بودن retrieval را محدودیت می‌داند»).