

Decide less, communicate more: On the construct validity of end-to-end fact-checking in medicine

Anonymous ACL submission

Abstract

Technological progress has led to concrete advancements in tasks that were regarded as challenging, such as automatic fact-checking. Interest in adopting these systems for public health and medicine has grown due to the high-stakes nature of medical decisions and challenges in critically appraising a vast and diverse medical literature. Evidence-based medicine connects to every individual, and yet the nature of it is highly technical, rendering the medical literacy of majority users inadequate to sufficiently navigate the domain. Such problems with medical communication ripens the ground for end-to-end fact-checking agents: check a claim against current medical literature and return with an evidence-backed verdict. And yet, such systems remain largely unused.

In this position paper, developed with expert input, we present the first study examining how clinical experts verify real claims from social media by synthesizing medical evidence. In searching for this upper-bound, we reveal fundamental challenges in end-to-end fact-checking when applied to medicine: Difficulties connecting claims in the wild to scientific evidence in the form of clinical trials; ambiguities in underspecified claims mixed with mismatched intentions; and inherently subjective veracity labels. We argue that fact-checking should be approached and evaluated as an interactive communication problem, rather than an end-to-end process.

We will release data and code upon publication.

1 Introduction

Decision making in medicine is personal, intimate, and high-stakes. Traditionally the patient—often a lay person unfamiliar with medicine—converses with their care providers about questions about their health. However, the reality is far from this picture: most Americans resort to the web when they have a health-related question (Fox and Duggan, 2013).

Today, social media and AI have made medical knowledge seemingly accessible. But claims made by others on the web (or by a chatbot) can be inaccurate or inapplicable. This, combined with eroding health literacy (Champlin et al., 2017), have led to challenges in public health (Hassan et al., 2015) as well as patient-provider sessions.

Meanwhile, evidence-based medicine has continuously evolved, with an evidence base growing too rapidly for physicians to keep up (Bastian et al., 2010; Marshall et al., 2021). This does provide a unique opportunity for AI fact-checking systems: Advances in retrieval systems and Large Language Models (LLMs) have increased interest in fact-checking systems that can classify medical claims as ‘True’ or ‘False’ with supporting evidence. However, despite technological advances, such systems remain underutilized as they struggle to address diverse *claims in the wild*—naturally occurring statements, usually made by laypeople, that pervade public discourse (Das et al., 2023; Chen et al., 2022); this further applies to claims made by AI agents that can be sometimes questionable. Existing fact-checking datasets often miss such claims because they were collected from already curated sources such as fact-checking websites and news articles (Kotonya and Toni, 2020; Vladika et al., 2024). Prior datasets also extracted claims from their context (Sarrouti et al., 2021), generated synthetic claims (Saakyan et al., 2021), or filtered claims based on lexical criteria (Mohr et al., 2022). As a result, systems trained on these heavily curated datasets are likely to fail to understand real medical claims made by the public.

This position paper highlights practical gaps in AI-driven fact-checking systems that cannot be addressed by “building a better system” alone. We focus on real-world medical claims from social media, preserving their original context, and contend that fact-checking systems should mirror how experts (e.g., physicians, care providers) evaluate and

Dataset	Domains	Source	Labels	Explanations	Evidence Type	Claim Example
PUBHEALTH (2020)	Public Health	Fact-Checking Websites, health news	True, Unproven, False, Mixture	✗	Sentences from same claim article source	Expired boxes of cake and pancake mix are dangerously toxic.
SCIFACT (2020)	Science	Expert-written	Supports, Refutes	✗	Scientific Articles	Rapamycin slows aging in fruit flies.
HEALTHVER (2021)	COVID-19	News Articles, blogs, social media	Supports, Refutes, Neutral	✗	Scientific Articles on COVID-19	Coronavirus may have originated in bats or pangolins
COVID-FACT (2021)	COVID-19	Reddit	Supported, Refuted	✗	Google Search Results	Baricitinib restrains the immune dysregulation in COVID-19 patients
COVERT (2022)	COVID-19	Twitter	Supports, Refutes, Not Enough Information	✗	Google Search Results	5G networks caused covid
REDHOT (2023)	Medical Conditions	Reddit	N/A	N/A	Randomized Controlled Trial Abstract	Link between RA and migraines
HEALTHFC (2024)	Health	Medizin Transparent	Supported, Refuted, Not Enough Information	✓	Systematic Review and Clinical Trial	Does cat's claw improve joint disease symptoms?
OUR CASE STUDY (2025)	Health	Reddit	No Relevant Abstracts, Refutes, Partially Refutes, Inconclusive, Partially Supports, Supports	✓	Randomized Controlled Trial Abstract	Contextualized Claim (3)

Table 1: Comparative overview of related work in medical and health fact-checking. *N/A* indicates components not applicable to the task (e.g., REDHOT does not perform claim verification).

respond to such claims. As part of our exploration, we designed an annotation study that examines how medical experts verify claims using retrieved medical evidence. Experts were asked to assess medical claims present on social media (i.e., Reddit forums about a particular medical condition) by synthesizing retrieved randomized controlled trial (RCT) abstracts and explaining their judgments. This provides an idealized upper bound for systems verifying “*in the wild*” claims. However, we highlight fundamental obstacles that challenge the construct validity of end-to-end automated systems for fact-checking: given a claim, provide a veracity judgment. We identify inherent difficulties in this setup for even domain experts, including: connecting claims with evidence; ambiguity from underspecified claims leading to valid yet contradictory interpretations; and challenges in achieving annotator consensus due to the inherent subjectivity of veracity labels. These issues suggest that the existing framing fact-checking of as an end-to-end classification task is inadequate for real-world settings, which may explain in part why such systems have not been put into wide use.

To correct the flawed construct validity of this task, we contend that **fact-checking should be an interactive dialogue agent rather than an end-to-end system**. We envision a human-centered **communication model** for medical fact-checking inspired by interactions between patients and physicians. We explain how this model can overcome ex-

isting challenges and empower experts and laypeople to engage in constructive medical discourse.

2 Background: medical claim-checking

Guo et al. (2022) outlines the conventional framework for automated fact-checking which comprises three stages: (1) **Claim Detection**, (2) **Evidence Retrieval**, and (3) **Claim Verification**. In the **Claim Detection** stage, the system identifies claims—statements asserting verifiable facts—and often ranks them based on check-worthiness factors such as public interest, popularity, timeliness, and impact (Das et al., 2023; Micallef et al., 2022). Complex claims may also be automatically decomposed into sub-claims for individual verification (Wanner et al., 2024; Pan et al., 2023; Min et al., 2023; Kamoi et al., 2023a; Jing et al., 2024). **Evidence Retrieval** entails retrieving supporting evidence to inform verification (Chen et al., 2024). Finally, **Claim Verification** requires determining the claim’s veracity and generating a justification grounded in the retrieved evidence. There is a growing interest in using LLMs to automate these stages of the fact-checking pipeline (Vykopal et al., 2024; Iqbal et al., 2024; Quelle and Bovet, 2024).

We present a comparative overview of prior work in medical fact-checking in Table 1. With few exceptions (e.g., (Wadhwa et al., 2023) which we use in this study), existing work largely views claims as statements that “stand alone” without context, and has treated fact-checking as an end-to-end pipeline

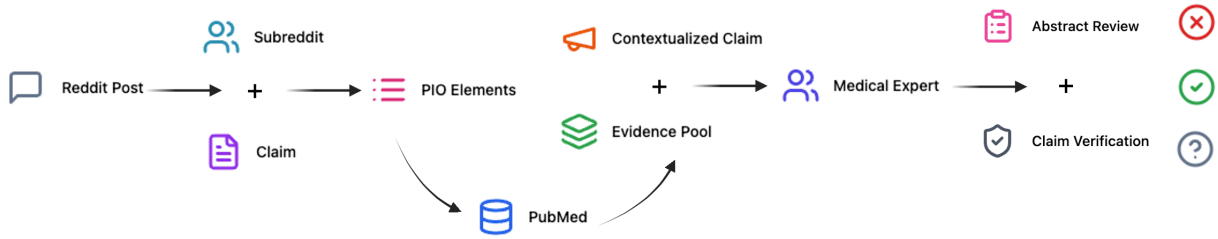


Figure 1: Overview of our AI-in-the-loop expert study pipeline. Given a claim from a subreddit, we extract PIO elements and automatically retrieve evidence. The evidence and its context are presented to a medical expert, who provides a veracity judgment and grounded rationale.

with the last step as a multi-label classification task (Sarrouti et al., 2021). Additionally, no work has yet examined **expert involvement in every stage of medical claim checking** – a fine-grained examination of claim interpretation, retrieved evidence, and veracity judgment, with natural language explanations. As a result, we do not yet have an understanding of an upper-bound that modern systems could achieve on evidence-based medical claim verification. Finally, while formal evidence synthesis has long been studied in the health literature (Thoma and Eaves, 2015; Sackett et al., 1996; Moberg et al., 2018; Cumpston and Thomas, 2019), it has not been integrated into LLM-based medical fact-checking systems.

3 Case Study Methods

To reveal the gaps in end-to-end fact-checking systems, we set out to establish an “idealized” scenario: given the task construction for automatic fact-checking systems discussed in Section 2, we asked domain experts to perform all subtasks that are delegable to humans. Thus, we formulate the following study for expert medical claim checking using the AI-in-the-loop pipeline shown in Figure 1.

Claim source We use claims from the Reddit Health Online Talk (RedHOT) corpus (Wadhwa et al., 2023), which contains 22,000 annotated posts from Reddit across 24 health conditions. RedHOT defines a claim as a statement indicating (often only implicitly) a causal relationship between an intervention and an outcome. To help experts contextualize these claims, we provide the full post and **Population, Intervention, Outcome (PIO)** descriptors. (See Appendix C, Appendix I for derivation details.) Note that we leave off the **Comparators** because in practice claims rarely mention the comparator explicitly (e.g., “Vitamin C cured my flu”).

Given a contextualized medical claim, a medical expert constructs a hierarchy of evidence (Guyatt et al., 2008a,b, 2011) based on relevance and quality, assesses the claim’s veracity, and provides a grounded explanation.

Annotation task setup The overall task, without any aid, places a high cognitive load on experts who must comb through and synthesize multiple pieces of evidence (Juneja and Mitra, 2022). With the aim of easing this, we automated several steps of the fact-checking process and integrated them into an intuitive web-based annotation interface (see Appendix G). These features enable experts to focus on critical aspects: evaluating evidence relevance and synthesizing it to support or refute claims. Detailed annotation guidelines are provided in Appendix E.

For each claim, we used its PIO elements to automatically retrieve ten published relevant RCT abstracts from Trialstreamer (Marshall et al., 2020), a continuously updated database of RCTs, as potential evidence. We provide experts with these RCT abstracts and their publication dates. We used a dense retrieval system using state-of-the-art embedding models. We provided a more detailed description of this retrieval methodology in Appendix K. We do not evaluate the feasibility of manual evidence search in this study. However, experts do provide judgments on how relevant retrieved abstracts are to the claim in question. For each RCT abstract, we collected annotations determining its **relevance** to the claim along four dimensions: **Population, Intervention, Outcome (PIO)**, and overall relevance. Claims are *contextualized* in their original Reddit posts during annotation. Annotators labeled each dimension as **(1) Relevant**, **(2) Somewhat Relevant**, or **(3) Irrelevant**. If an abstract was deemed *overall* relevant, annotators highlighted the most relevant text span and assessed whether the trial described in the abstract supports or refutes

the claim using the four labels: (1) Supports, (2) Partially Supports, (3) Partially Refutes, and (4) Refutes.

After annotating all ten abstracts, experts proceeded to synthesize the evidence. To help experts navigate through the evidence documents, we implemented a tiering step. Abstracts are initially tiered automatically based on their relevance annotations in the previous step, establishing a natural hierarchy. Annotators are free to further refine this hierarchy by considering evidence quality. Next, annotators verify the claim in two phases:

1. **Overall Support:** Verification based solely on the provided evidence.
2. **Expert Support:** Optional verification based on clinical expertise, particularly for claims unlikely to be studied in clinical trials.

This separation allows us to compare expert opinion with evidence-based conclusions and minimize bias. Annotators select from six labels for both phases: (1) No Relevant Abstracts/No Expert Opinion (for each of the above phases respectively), (2) Refutes, (3) Partially Refutes, (4) Inconclusive, (5) Partially Supports, and (6) Supports.

To justify their veracity label, experts write a paragraph-length explanation (see guideline in Appendix F). Annotators may optionally include a medical addendum detailing clinical practices typically used in response to the claim, providing practical context for users.

Annotation team Our annotation team consists of five clinical experts, one serving as the medical lead. All experts had experience reviewing medical articles and synthesizing them for biomedical research or patient care (annotator recruitment details are in Appendix H). To leverage their expertise effectively, we conducted the study in **three** rounds, with changes between rounds detailed in Appendix D. Following Klie et al. (2024), the two co-first authors held group meetings with the experts during each round to discuss disagreements and refine the annotation guidelines. In total, these meetings spanned four hours.

Across the three rounds, each of the five experts annotated 20 unique claims, yielding 1,000 abstract-level annotation instances (10 abstracts per claim) and 100 synthesis-level explanations.

Type	κ (\uparrow)
Population	0.416
Intervention	0.714
Outcome	0.200
Overall (Abstract-level)	0.155
Tab Support	0.170
Overall Support	0.124
Expert Support	-0.184

Table 2: Inter-annotator agreement for each portion of the fact-checking pipeline. **Blue**: abstract-level labels; **pink**: synthesis-level labels.

4 Challenges for end-to-end fact-checking

In this section, we present the fundamental challenges revealed by our expert annotations and inputs, gathered over multiple rounds of discussion.

4.1 Challenge 1: Connecting Medical Evidence with Claims

We present the inter-annotator agreement on the final round of annotations for five claims annotated by five experts (totaling 25 separate verifications) in Table 2. This round includes 50 abstracts per expert (250 total), with five labels per abstract, and two labels per claim. Despite multiple rounds of expert feedback to improve the annotation task, agreement remained low across all fields. A guideline for two annotators on two classes suggests that a Fleiss’ κ score of 0.21-0.40 represents fair agreement, 0.41-0.60 moderate agreement, and 0.61-0.80 substantial agreement (Landis and Koch, 1977). Since κ scores tend to be higher with fewer categories (Sim and Wright, 2005), we consider a reasonable κ score for a task with 4+ labels to be approximately 0.5.

Most instances—20 out of the 25 expert judgments—were labeled “No Relevant Abstracts”, indicating that the claims were unverifiable. For three of the five claims in this final round, all experts independently labeled them “No Relevant Abstracts” (see Section A for an illustrative example). This high rate of unverifiable claims underscores a broader challenge: even when claims are annotated to suggest a causal relationship between an Intervention and an Outcome (Wadhwa et al., 2023) and are paired with state-of-the-art evidence retrieval, they are often unverifiable in practice.

Our expert annotation study identifies four reasons why:

1. **No Intervention:** Claims lacking an intervention cannot be verified through an RCT.

Grapefruit**Post (r/Epilepsy):** Grapefruit and seizuresGuys, I am wondering if you have any issues or know about interactions between oxcarbazepine and/or levetiracetam and grapefruit? **I believe it may make those medications work differently, but I am not sure.****Population:** Epilepsy patients (implied by the subreddit r/Epilepsy and the mention of seizure medications)**Intervention:** Grapefruit consumption (in interaction with oxcarbazepine and/or levetiracetam)**Outcome:** Medication efficacy (i.e., how the medications work)

Table 3: An example of an unverifiable claim, as no RCTs have examined interactions between grapefruit, oxcarbazepine, and epilepsy, and such a study may be infeasible.

- 308 2. **Unethical Intervention:** Some interventions
309 are unethical to test via RCTs because they may
310 harm participants. For example, it is unethical
311 to study smoking as an intervention in an RCT.
- 312 3. **Lack of Feasibility:** Claims involving specific
313 PIO element combinations are often unverifi-
314 able due to the impracticality or improbability
315 of conducting such RCTs.
- 316 4. **Lack of Utility:** Some claims, while theoretic-
317 ally verifiable through an RCT, lack available
318 evidence as findings from such studies would
319 lack utility in the medical field.

320 These issues highlight the difficulty of collecting
321 evidence that is directly relevant to claims peo-
322 ple make (on social media). Medical evidence,
323 especially high-quality evidence, is bounded and
324 restricted by standards for feasibility and ethics in
325 a way that covers all the possible queries from the
326 public is impossible. Prior to annotation, we used
327 an automated method to filter out claims that could
328 not be verified by RCTs (detailed in Appendix J).
329 This issue of unverifiability remained despite this
330 effort. Part of addressing this challenge lies in
331 expanding the pool of evidence while ensuring it
332 remains trustworthy (see Section 5.4). Addition-
333 ally, systems should tackle how to best handle the
334 inevitable case in which a claim is unverifiable.
335 We discuss a guided retrieval approach to this in
336 Section 5.2.

337 4.2 Challenge 2: Variations in the 338 Interpretation of Claims

339 In our discussions with annotators, we noted consis-
340 tent disagreements in how they interpreted medical
341 claims, which in turn caused disagreement in the
342 claim verification task. To address this, we pro-
343 vided annotators with PIO descriptors (Table 3) of
344 claims to narrow the scope of possible interpre-
345 tations. However, even with this added context
346 reaching consensus among annotators remained a

challenge. This is because *claims in the wild* about
health tend to be *underspecified* and/or *misguided*,
causing annotators to deduce their own varying in-
terpretations on what the *patient*, the author of the
claim (usually a layperson), intended.

Underspecified Claims Naturally occurring
medical claims on social media are usually written
informally by laypeople, and tend therefore to be
underspecified (see Table 4). For example, a pa-
tient with ADHD claimed that “herbs are proven
to be helping with regulating the cycle.” It is un-
clear what “regulating” means here, and annota-
tors interpreted this in various ways, e.g., reducing
symptoms related to menstruation or skipping men-
struation altogether. Another patient claimed that
pineapple juice is “good for sinus issues” by reduc-
ing inflammation. However, it was unclear whether
they meant immediate or gradual improvement (no
time frame was offered). Resolving these under-
specifications requires understanding the patients’
intentions, which is challenging since intent cannot
always be inferred from the claim and its context
alone.

Misguided Claims Discussions with annotators
also identified another artifact of naturally occur-
ring medical claims: *misguided* claims formed
from incorrect premises. Annotators often dis-
agreed on how to handle such claims within our
task. The previously mentioned ADHD exam-
ple illustrates this issue. The patient, prescribed
methylphenidate, noticed no effect during the last
10 days of their menstrual cycle. Their doctor sug-
gested contraceptive pills as a potential solution.
Concerned about the side effects of contraceptives,
the patient considered using herbs as an alternative
to “regulate the cycle”. Annotators characterized
the underlying premise—that the medication’s ef-
ficacy is affected by the menstrual cycle—as in-
correct. Disagreement over whether to consider
the premise’s validity in veracity judgments led to

ADHD, Herbs, and Menstruation		
Post (r/ADHD):	Hello, menstruating people! How do your cycle and ADHD influence each other and how do you deal with it?	
EDIT:	After getting your responses I am reflecting again how medicine does not give a shit about women. It's truly insane. Thank you!	
	Hello! I have never paid too much attention to my menstrual cycle since it was never particularly bothersome. Now that I take methylo I feel big changes in how I function during the cycle. Like last 10 days of the cycle, my medication kind of stops working... That is like 1/3 of the time. I know it's still better than without meds nevertheless, it makes establishing a routine quite challenging. My doc suggested trying contraceptive pills, but I am not even sexually active ATM so taking more medication, with potential side effects, does not excite me.	
	I know there are herbs that are proven to be helping with regulating the cycle but I don't know if they would help with ADHD symptoms? Any tips?	
Population:	People with ADHD	Intervention: Herbs Outcome: Regulating the menstrual cycle
Expert Feedback:	<ul style="list-style-type: none"> - What does regulating the cycle mean? - ADHD has no bearing on one's menstruation cycle. It is a red herring. - Trials with these descriptors are unlikely to exist. - Is this claim really what the patient is concerned about in this post? 	
Pineapple Juice Reduces Inflammation		
Post (r/CysticFibrosis):	Anyone with sinus issues drinking pineapple juice?	
	It's a weird question, but I saw a post about pineapple juice being good for sinus issues (helps with the inflammation) and just wondered if anyone has done this? Some people were commenting about the high sugar content in pineapple juice not being good, but they get around that by taking a supplement instead of drinking the juice. Anyone?	
Population:	Patients With Cystic Fibrosis	Intervention: Pineapple Juice Outcome: Reduced Inflammation/Fewer Sinus Issues
Expert Feedback:	<ul style="list-style-type: none"> - How quickly is the poster expecting the intervention to produce results? - Just improving inflammation should not be the only criteria. - Trials with these exact descriptors are unlikely to exist. 	

Table 4: Examples of underspecified claims and corresponding expert feedback.

387 conflicting assessments among annotators. Prior
388 work in general-domain fact-checking used claim
389 decomposition to address false presuppositions in
390 claims (Chen et al., 2022; Kamoi et al., 2023b; Hu
391 et al., 2025), however this does not tackle *implicit*
392 premises.

393 **Mismatched Intent** Previous work on fact-
394 checking has addressed underspecified claims, of-
395 ten by decontextualizing them, i.e., removing con-
396 text and resolving underspecifications based on lo-
397 cal content (Deng et al., 2024; Gunjal and Durrett,
398 2024). This approach can clarify underspecifica-
399 tions, but it disconnects the claim from the patient's
400 original intent, as embedded in the global context.
401 This disconnect can result in verifications that do
402 not apply to the original claim, allowing subtle
403 falsehoods to slip through and potentially be ampli-
404 fied. For example, suppose the underspecification
405 of “regulating the cycle” were resolved and the
406 claim were deemed true without considering its
407 premise, this verification would fail to address the
408 patient's true goal of improving the efficacy of their
409 ADHD medication. Such verification would also
410 implicitly validate the misunderstood premise. It is
411 also the case that the patient's true intention is not
412 about this particular claim, but an overall desire to
413 communicate and discuss the underlying condition
414 to get better. To effectively address this, the focus
415 must be on meeting the patient's *information needs*,
416 and when needed, uncovering and assessing their
417 assumptions.

4.3 Challenge 3: Labeling the Severity of Inaccurate Statements is Inherently Subjective

421 Another factor that contributed to the disagree-
422 ments we observed is the subjectivity in labeling
423 the *veracity*, or the degree of truth of a medical
424 claim. This subjectivity, as we observe, is not
425 caused by differing interpretations of the claim.
426 Rather, experts are influenced by their backgrounds
427 and philosophies, applying different standards for
428 assessing a claim's “truthfulness” based on evi-
429 dence.

430 Consider the claim about pineapple juice in Ta-
431 ble 4. If the desired onset for resolving sinus issues
432 is clarified to mean within a few days, the claim
433 is technically false. However, the *severity* of this
434 falsehood is subjective. One expert might view it as
435 a minor inaccuracy, noting that pineapple juice may
436 help with sinus issues but works more slowly and
437 less effectively than targeted treatments. Another
438 expert might see it as a serious falsehood, arguing
439 that promoting pineapple juice as a quick fix is mis-
440 leading and potentially harmful. Both perspectives
441 are valid, highlighting the inherent subjectivity in
442 judging a claim's truthfulness. Experts' sensitivi-
443 ties can also vary depending on the topic, leading to
444 apparently inconsistent judgments. These findings
445 suggest that the existing practice—simply asking
446 for a veracity label—needs to be redefined to align
447 expert opinions in the first place.

5 How Can We Address These Challenges?

To address the challenges in Section 4, we argue that fact-checking alone is insufficient: A **communication model** is required to mirror a dialogue between a healthcare provider and a patient. This frames fact-checking as a dialogue aimed at addressing patient’s information needs. A vision of the user interaction with the communication model is shown in Figure 2.

5.1 Clarifying Intent Through Conversation

In Section 4.2, we describe how naturally occurring claims often contain underspecifications that require understanding the patient’s intentions. A communication model can address this through dialogue with the patient, similar to clinical interactions where physicians asking patients questions to understand their care needs. Prior work in dialogue systems has explored resolving ambiguity by generating clarification questions and modeling future conversations (Kim et al., 2023; Zhang et al., 2024a; Zhang and Choi, 2023).

To do this effectively, the system must identify underspecifications from the provided context. Similar work identifying ambiguities in user queries could provide a template for this task (Zhang et al., 2024b). However, as we discovered in our analysis, identifying these underspecifications also requires extensive expert knowledge of the claim’s subject matter, which could be encoded in a trained model or accessed via a retrieval-augmented system. The system must also identify “misguided” claims, as discussed in Section 4.2, which requires recognizing subtextual implications and commonly held misconceptions. An ideal system would proactively query the patient to uncover incorrect assumptions and address them with empathy, fostering constructive engagement.

Direct communication with the patient is not strictly necessary to achieve intent clarification. (Kim et al., 2023) proposed generating a tree of clarification questions, which, when fully answered, provides the context needed to resolve an ambiguous query. A similar approach could be applied here, where a tree of clarification questions resolves different interpretations of an underspecified claim, all of which must be verified and included in the system’s final output. However, for this approach to work research is needed to study how to align the final output with the patient’s original intent.

5.2 Guided Retrieval of Medical Evidence

In Section 4.1, we discussed the disconnect between what is practical and measurable in evidence-based medicine and what patients care about. The communication model enables providers to guide such claims toward verifiability while clarifying the patient’s intent. To support this process, evidence retrieval should inform the dialogue between the patient and the system. When no relevant evidence is found, the system should communicate this and guide the patient toward related, verifiable claims. When none is available, the system should *abstain*. This approach makes clear that the claim is unverifiable while offering a pathway for continued learning.

This *guided retrieval* could also help correct misguided claims, as this conversational approach digs into and exposes the patient’s thought process. Similar to physician-patient interactions, this process resembles a physician’s response to an unanswerable query. They might first gather more information about the patient’s question; if it remains unanswerable, they might recall related (answerable) questions. The field of Interactive Information Retrieval (IIR) studies the modeling and optimization of such back-and-forth interactions between users and retrieval systems (Zhai, 2020), e.g., for product retrieval (Wang et al., 2024; Aliannejadi et al., 2024). Similarly, this approach could be used to guide unverifiable claims—often misguided due to false assumptions—toward claims that satisfy the patient’s *information needs*.

5.3 Communicating Veracity Through Diverse Perspectives

As discussed in Section 4.3, our annotation study demonstrated that categorizing claims with fine-grained veracity labels is inherently subjective. While annotators often disagreed on labels, their reasoning in plain language explanations was often similar. During discussions, annotators often accepted each other’s explanations as valid despite disagreeing on the level of severity. This suggests the possibility of a wider range of ostensible “agreement” can be reached. We propose that an effective medical fact-checking agent should produce responses that reflect diverse expert perspectives, acknowledging the inherent heterogeneity of expert evaluation. The need for these diverse explanations is corroborated by fact-checking professionals, who acknowledge the need and complexity of

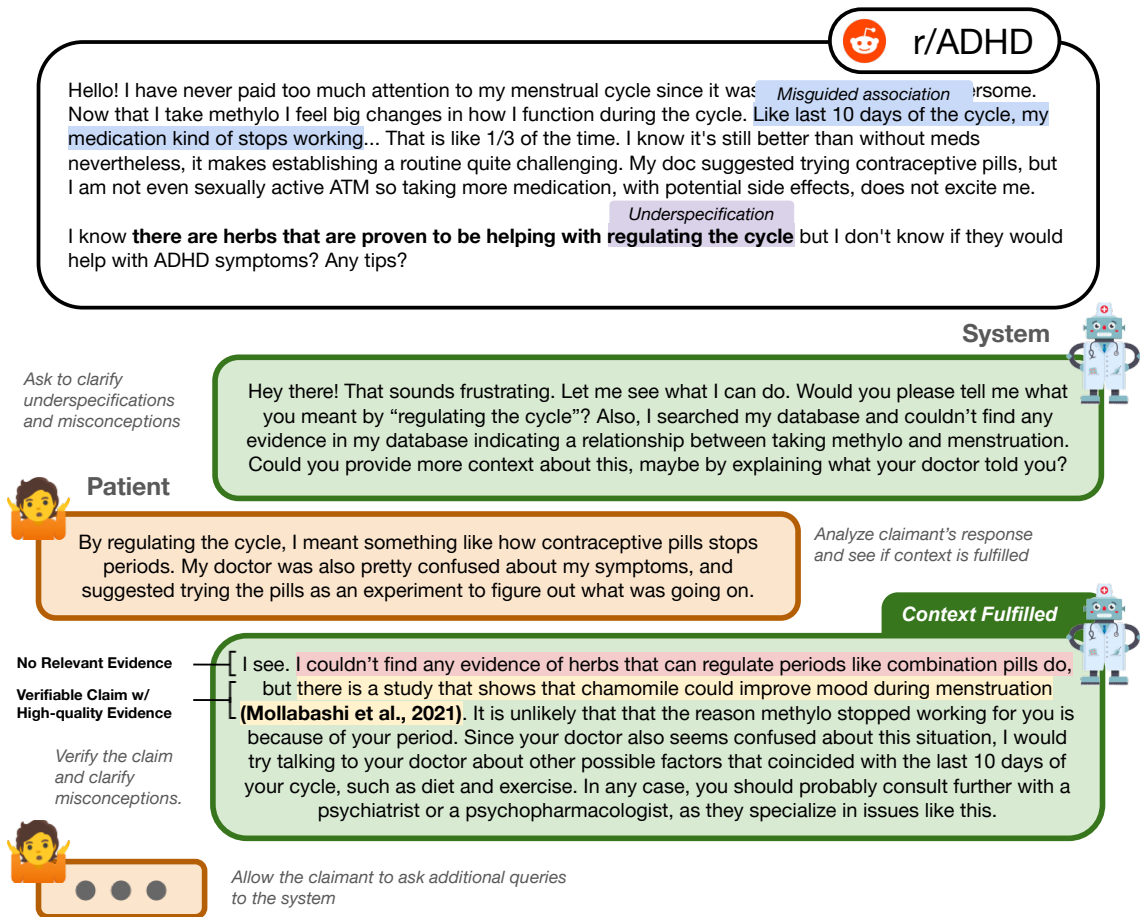


Figure 2: This figure illustrates the communication model for fact-checking, where the system engages the patient by asking clarifying questions, filling contextual gaps, and verifying claims while addressing misconceptions.

548 crafting thorough, nuanced explanations and calls
 549 for explanations to accommodate different audi-
 550 ence needs (Warren et al., 2025). Encouraging
 551 response diversity, rather than imposing artificial
 552 consensus via a single numerical value, is crucial
 553 for developing multi-agent medical fact-checking
 554 systems that integrate multiple expert viewpoints.

555 5.4 Medical Evidence Beyond RCTs

556 To bridge the gap between user questions and the
 557 limited pool of RCTs, future work could incor-
 558 porate other forms of medical evidence such as
 559 meta analyses, cohort studies, case-control stud-
 560 ies, case series, and case reports. Expanding the
 561 scope of medical evidence increases complexity, as
 562 systems must determine how to retrieve and synthe-
 563 size evidence of various types and strengths, how
 564 to appropriately communicate this to end users. To
 565 this end, systems should adequately communicate
 566 its uncertainty according to existing guidelines in
 567 medicine (Ratcliff et al., 2021; Simpkin and Arm-
 568 strong, 2019) and reference lower-grade evidence

only when the communication model fails to iden-
 569 tify a helpful, verifiable question and should clearly
 570 convey the quality of the source evidence and its
 571 limitations.
 572

573 6 Conclusion

574 In this position paper, we present an analysis of
 575 the construct validity of existing end-to-end auto-
 576 matic medical fact-checking systems, with expert
 577 engagement in all key aspects of the system. Our
 578 work highlights the unique challenges of automated
 579 medical fact-checking, showing that it should be
 580 approached with user interaction in mind, and not
 581 as an end-to-end system. We propose a commu-
 582 nication model to clarify underspecifications and
 583 guide unverifiable claims, aiming to improve user
 584 outcomes and real-world utility. We hope this work
 585 inspires further exploration of human-in-the-loop
 586 systems for medical fact-checking.

587 Limitations

588 Due to time and cost constraints, our annotation
589 study includes a limited number of claims. How-
590 ever, each claim was evaluated against 10 retrieved
591 abstracts, resulting in a substantial volume of
592 evidence-based annotations. These annotations and
593 expert discussions revealed fundamental challenges
594 in medical fact-checking that make large-scale an-
595 notation difficult to define in a principled way.

596 We used an automatic document retrieval system
597 instead of manual search by annotators to iden-
598 tify relevant abstracts for each claim. While this
599 approach may fail to retrieve the most relevant evi-
600 dence, retrieval performance was not the focus of
601 this study. Rather, our goal was to examine how
602 retrieved medical evidence is synthesized to ver-
603 ify claims, while avoiding undue burden on expert
604 annotators.

605 These limitations reflect broader challenges in
606 defining and evaluating medical fact-checking tasks
607 in practice.

608 Ethical Considerations

609 The posts found within the RedHOT dataset con-
610 tain health-related comments that are inherently
611 sensitive. To respect this sensitivity, the authors
612 of the RedHOT dataset notified all users of their
613 inclusion in this dataset and provided them with the
614 opportunity to opt-out. They also did not release
615 the data directly, but instead provided a script to
616 download content from Reddit so that individuals
617 can remove their post in the future. In this work,
618 we directly release a small subset of these posts
619 from RedHOT that we used in our annotation study.
620 In our released data, we do not reveal the username
621 of the author of the post. We only include the text
622 from the post and information about the subreddit
623 in which it was found. Considering the measures
624 taken by the authors of the RedHOT dataset and the
625 fact that these posts have been publicly available
626 on Reddit for at least more than a year, we believe
627 it is safe to publicly release this data.

628 We have consulted with our Institutional Review
629 Board (IRB) about the nature of our work and con-
630 firmed that the use of the RedHOT data and the
631 subsequent annotation study using these data do
632 not constitute research of human subjects. How-
633 ever, we acknowledge that certain uses of this data
634 may be considered sensitive. We strongly encour-
635 age researchers to obtain prior approval from their
636 own IRB regarding the intended use of the data

released by this work.

We also consider ethical aspects of the annota-
tion process. Annotators consented to participate
and were compensated at a competitive hourly rate
(see Appendix H for details).

References

- Mohammad Aliannejadi, Jacek Gwizdka, and Hamed Zamani. 2024. [Interactions with generative information retrieval systems](#). *Preprint*, arXiv:2407.11605.
- Hilda Bastian, Paul Glasziou, and Iain Chalmers. 2010. Seventy-five trials and eleven systematic reviews a day: how will we ever keep up? *PLoS medicine*, 7(9):e1000326.
- Sara Champlin, Michael Mackert, Elizabeth M Glowacki, and Erin E Donovan. 2017. Toward a better understanding of patient health literacy: A focus on the skills patients need to find health information. *Qualitative Health Research*, 27(8):1160–1176.
- Jifan Chen, Grace Kim, Aniruddh Sriram, Greg Durrett, and Eunsol Choi. 2024. [Complex claim verification with evidence retrieved in the wild](#). *Preprint*, arXiv:2305.11859.
- Jifan Chen, Aniruddh Sriram, Eunsol Choi, and Greg Durrett. 2022. [Generating literal and implied sub-questions to fact-check complex claims](#). In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 3495–3516, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.
- Li T Page MJ Chandler J Welch VA Higgins JPT Cumpston, M and J Thomas. 2019. [Updated guidance for trusted systematic reviews: a new edition of the cochrane handbook for systematic reviews of interventions](#). *Cochrane Database of Systematic Reviews*, (10).
- Anubrata Das, Houjiang Liu, Venelin Kovatchev, and Matthew Lease. 2023. [The state of human-centered nlp technology for fact-checking](#). *Information Processing & Management*, 60(2):103219.
- Zhenyun Deng, Michael Schlichtkrull, and Andreas Vlachos. 2024. [Document-level claim extraction and decontextualisation for fact-checking](#). *Preprint*, arXiv:2406.03239.
- Susannah Fox and Maeve Duggan. 2013. Health online 2013. pew research center. *National survey by the Pew Research Center’s Internet and American Life Project*.
- Anisha Gunjal and Greg Durrett. 2024. [Molecular facts: Desiderata for decontextualization in LLM fact verification](#). In *Findings of the Association for Computational Linguistics: EMNLP 2024*, pages 3751–3768, Miami, Florida, USA. Association for Computational Linguistics.

690	Zhijiang Guo, Michael Schlichtkrull, and Andreas Vlachos. 2022. A survey on automated fact-checking . <i>Transactions of the Association for Computational Linguistics</i> , 10:178–206.	
691		
692		
693		
694	Gordon Guyatt, Andrew D. Oxman, Elie A. Akl, Regina Kunz, Gunn Vist, Jan Brozek, Susan Norris, Yngve Falck-Ytter, Paul Glasziou, Hans deBeer, Roman Jaeschke, David Rind, Joerg Meerpohl, Philipp Dahm, and Holger J. Schünemann. 2011. Grade guidelines: 1. introduction—grade evidence profiles and summary of findings tables . <i>Journal of Clinical Epidemiology</i> , 64(4):383–394.	
695		
696		
697		
698		
699		
700		
701		
702	Gordon H Guyatt, Andrew D Oxman, Regina Kunz, Gunn E Vist, Yngve Falck-Ytter, and Holger J Schünemann. 2008a. What is “quality of evidence” and why is it important to clinicians? <i>BMJ</i> , 336(7651):995–998.	
703		
704		
705		
706		
707	Gordon H Guyatt, Andrew D Oxman, Gunn E Vist, Regina Kunz, Yngve Falck-Ytter, Pablo Alonso-Coello, and Holger J Schünemann. 2008b. Grade: an emerging consensus on rating quality of evidence and strength of recommendations . <i>BMJ</i> , 336(7650):924–926.	
708		
709		
710		
711		
712		
713	Naeemul Hassan, Bill Adair, James T. Hamilton, Chengkai Li, Mark Tremayne, Jun Yang, and Cong Yu. 2015. The quest to automate fact-checking .	
714		
715		
716	Qisheng Hu, Quanyu Long, and Wenya Wang. 2025. Decomposition dilemmas: Does claim decomposition boost or burden fact-checking performance? <i>Preprint</i> , arXiv:2411.02400.	
717		
718		
719		
720	Hasan Iqbal, Yuxia Wang, Minghan Wang, Georgi Nenkov Georgiev, Jiahui Geng, Iryna Gurevych, and Preslav Nakov. 2024. OpenFactCheck: A unified framework for factuality evaluation of LLMs . In <i>Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing: System Demonstrations</i> , pages 219–229, Miami, Florida, USA. Association for Computational Linguistics.	
721		
722		
723		
724		
725		
726		
727		
728		
729	Liqiang Jing, Ruosen Li, Yunmo Chen, and Xinya Du. 2024. FaithScore: Fine-grained evaluations of hallucinations in large vision-language models . In <i>Findings of the Association for Computational Linguistics: EMNLP 2024</i> , pages 5042–5063, Miami, Florida, USA. Association for Computational Linguistics.	
730		
731		
732		
733		
734		
735	Prerna Juneja and Tanushree Mitra. 2022. Human and technological infrastructures of fact-checking . <i>Proc. ACM Hum.-Comput. Interact.</i> , 6(CSCW2).	
736		
737		
738	Ryo Kamoi, Tanya Goyal, Juan Diego Rodriguez, and Greg Durrett. 2023a. WiCE: Real-world entailment for claims in Wikipedia . In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing</i> , pages 7561–7583, Singapore. Association for Computational Linguistics.	
739		
740		
741		
742		
743		
744	Ryo Kamoi, Tanya Goyal, Juan Diego Rodriguez, and Greg Durrett. 2023b. Wice: Real-world entailment for claims in wikipedia . <i>Preprint</i> , arXiv:2303.01432.	
745		
746		
	Gangwoo Kim, Sungdong Kim, Byeongguk Jeon, Joon-suk Park, and Jaewoo Kang. 2023. Tree of clarifications: Answering ambiguous questions with retrieval-augmented large language models . <i>Preprint</i> , arXiv:2310.14696.	747
		748
		749
		750
		751
	Jan-Christoph Klie, Richard Eckart de Castilho, and Iryna Gurevych. 2024. Analyzing Dataset Annotation Quality Management in the Wild . <i>Computational Linguistics</i> , pages 1–50.	752
		753
		754
		755
	Neema Kotonya and Francesca Toni. 2020. Explainable automated fact-checking for public health claims . <i>Preprint</i> , arXiv:2010.09926.	756
		757
		758
	J Richard Landis and Gary G. Koch. 1977. The measurement of observer agreement for categorical data . <i>Biometrics</i> , 33 1:159–74.	759
		760
		761
	Iain J Marshall, Benjamin Nye, Joël Kuiper, Anna Noel-Storr, Rachel Marshall, Rory Maclean, Frank Soboczenski, Ani Nenkova, James Thomas, and Byron C Wallace. 2020. Trialstreamer: A living, automatically updated database of clinical trial reports . <i>Journal of the American Medical Informatics Association</i> , 27(12):1903–1912.	762
		763
		764
		765
		766
		767
		768
	Iain James Marshall, Veline L’Esperance, Rachel Marshall, James Thomas, Anna Noel-Storr, Frank Soboczenski, Benjamin Nye, Ani Nenkova, and Byron C Wallace. 2021. State of the evidence: a survey of global disparities in clinical trials . <i>BMJ Global Health</i> , 6(1):e004145. PMID: PMC7786802.	769
		770
		771
		772
		773
		774
	Nicholas Micallef, Vivienne Armacost, Nasir Memon, and Sameer Patil. 2022. True or false: Studying the work practices of professional fact-checkers . <i>Proceedings of the ACM on Human-Computer Interaction</i> , 6(CSCW1):1–44.	775
		776
		777
		778
		779
	Sewon Min, Kalpesh Krishna, Xinxu Lyu, Mike Lewis, Wen-tau Yih, Pang Koh, Mohit Iyyer, Luke Zettlemoyer, and Hannaneh Hajishirzi. 2023. FActScore: Fine-grained atomic evaluation of factual precision in long form text generation . In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing</i> , pages 12076–12100, Singapore. Association for Computational Linguistics.	780
		781
		782
		783
		784
		785
		786
		787
	Jenny Moberg, Andrew D Oxman, Sarah Rosenbaum, Holger J Schünemann, Gordon Guyatt, Signe Flottorp, Claire Glenton, Simon Lewin, Angela Morelli, Gabriel Rada, Pablo Alonso-Coello, and GRADE Working Group. 2018. The GRADE evidence to decision (EtD) framework for health system and public health decisions . <i>Health Res. Policy Syst.</i> , 16(1):45.	788
		789
		790
		791
		792
		793
		794
	Isabelle Mohr, Amelie Wüthl, and Roman Klinger. 2022. CoVERT: A corpus of fact-checked biomedical COVID-19 tweets . In <i>Proceedings of the Thirteenth Language Resources and Evaluation Conference</i> , pages 244–257, Marseille, France. European Language Resources Association.	795
		796
		797
		798
		799
		800
	Liangming Pan, Xiaobao Wu, Xinyuan Lu, Anh Tuan Luu, William Yang Wang, Min-Yen Kan, and Preslav	801
		802

803	Nakov. 2023. Fact-checking complex claims with program-guided reasoning . In <i>Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 6981–7004, Toronto, Canada. Association for Computational Linguistics.	860
804		861
805		862
806		863
807		864
808		865
809	Dorian Quelle and Alexandre Bovet. 2024. The perils and promises of fact-checking with large language models . <i>Frontiers in Artificial Intelligence</i> , 7.	866
810		867
811		868
812	Chelsea L Ratcliff, Bob Wong, Jakob D Jensen, and Kimberly A Kaphingst. 2021. The impact of communicating uncertainty on public responses to precision medicine research. <i>Annals of Behavioral Medicine</i> , 55(11):1048–1061.	869
813		870
814		871
815		872
816		873
817	Arkadiy Saakyan, Tuhin Chakrabarty, and Smaranda Muresan. 2021. COVID-fact: Fact extraction and verification of real-world claims on COVID-19 pandemic . In <i>Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)</i> , pages 2116–2129, Online. Association for Computational Linguistics.	874
818		875
819		876
820		877
821		878
822		879
823		880
824		881
825		882
826	David L Sackett, William M C Rosenberg, J A Muir Gray, R Brian Haynes, and W Scott Richardson. 1996. Evidence based medicine: what it is and what it isn't . <i>BMJ</i> , 312(7023):71–72.	883
827		884
828		885
829		886
830	Mourad Sarrouti, Asma Ben Abacha, Yassine Mrabet, and Dina Demner-Fushman. 2021. Evidence-based fact-checking of health-related claims . In <i>Findings of the Association for Computational Linguistics: EMNLP 2021</i> , pages 3499–3512, Punta Cana, Dominican Republic. Association for Computational Linguistics.	887
831		888
832		889
833		890
834		891
835		892
836		893
837	Julius Sim and Chris C Wright. 2005. The kappa statistic in reliability studies: Use, interpretation, and sample size requirements . <i>Physical Therapy</i> , 85(3):257–268.	894
838		895
839		896
840		897
841	Arabella L Simpkin and Katrina A Armstrong. 2019. Communicating uncertainty: a narrative review and framework for future research . <i>Journal of general internal medicine</i> , 34:2586–2591.	898
842		899
843		900
844		901
845	Achilleas Thoma and III Eaves, Felmont F. 2015. A brief history of evidence-based medicine (ebm) and the contributions of dr david sackett . <i>Aesthetic Surgery Journal</i> , 35(8):NP261–NP263.	902
846		903
847		
848		
849	Juraj Vladika, Phillip Schneider, and Florian Matthes. 2024. HealthFC: Verifying health claims with evidence-based medical fact-checking . In <i>Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024)</i> , pages 8095–8107, Torino, Italia. ELRA and ICCL.	
850		
851		
852		
853		
854		
855		
856	Ivan Vykopal, Matúš Pikuliak, Simon Ostermann, and Marián Šimko. 2024. Generative large language models in automated fact-checking: A survey . <i>Preprint</i> , arXiv:2407.02351.	
857		
858		
859		
	Somin Wadhwa, Vivek Khetan, Silvio Amir, and Byron Wallace. 2023. RedHOT: A corpus of annotated medical questions, experiences, and claims on social media . In <i>Findings of the Association for Computational Linguistics: EACL 2023</i> , pages 809–827, Dubrovnik, Croatia. Association for Computational Linguistics.	
	Mengzhao Wang, Haotian Wu, Xiangyu Ke, Yunjun Gao, Xiaoliang Xu, and Lu Chen. 2024. An interactive multi-modal query answering system with retrieval-augmented large language models . <i>Preprint</i> , arXiv:2407.04217.	
	Miriam Wanner, Seth Ebner, Zhengping Jiang, Mark Dredze, and Benjamin Van Durme. 2024. A closer look at claim decomposition . In <i>Proceedings of the 13th Joint Conference on Lexical and Computational Semantics (*SEM 2024)</i> , pages 153–175, Mexico City, Mexico. Association for Computational Linguistics.	
	Greta Warren, Irina Shklovski, and Isabelle Augenstein. 2025. Show me the work: Fact-checkers' requirements for explainable automated fact-checking . In <i>Proceedings of the 2025 CHI Conference on Human Factors in Computing Systems</i> , CHI '25, New York, NY, USA. Association for Computing Machinery.	
	ChengXiang Zhai. 2020. Interactive information retrieval: Models, algorithms, and evaluation . In <i>Proceedings of the 43rd International ACM SIGIR conference on research and development in Information Retrieval, SIGIR 2020, Virtual Event, China, July 25-30, 2020</i> , pages 2444–2447. ACM.	
	Michael J. Q. Zhang and Eunsol Choi. 2023. Clarify when necessary: Resolving ambiguity through interaction with lms . <i>Preprint</i> , arXiv:2311.09469.	
	Michael J. Q. Zhang, W. Bradley Knox, and Eunsol Choi. 2024a. Modeling future conversation turns to teach llms to ask clarifying questions . <i>Preprint</i> , arXiv:2410.13788.	
	Tong Zhang, Peixin Qin, Yang Deng, Chen Huang, Wenqiang Lei, Junhong Liu, Dingnan Jin, Hongru Liang, and Tat-Seng Chua. 2024b. Clamber: A benchmark of identifying and clarifying ambiguous information needs in large language models . <i>Preprint</i> , arXiv:2405.12063.	

A Trikafta Claim Example

Dandruff and Trikafta		
Post (r/CysticFibrosis): Anyone had really bad flaky scalp or dandruff lately ? Think it could be due to trikafta . Could it be anything else		
Population:	Patients with Cystic Fibrosis (implied by the subreddit r/CysticFibrosis)	
Intervention:	Trikafta (a medication)	
Outcome:	Flaky scalp or dandruff	
Expert 1: None of the abstracts directly addressed patients with cystic fibrosis experiencing dandruff/skin-related adverse effects secondary to trikafta use. All of the abstracts include some form of scalp flaking whether it be dandruff in general or specific conditions such as seborrheic dermatitis. However, they cannot be considered relevant as none of them address patients with cystic fibrosis or experiencing dandruff secondary to medication side effect.	Expert 2: A1, a2, a3, a4, a6 include the treatment of dandruff in a population with dandruff. The only relevant element in these abstracts is the outcome measures. A5, a7 includes a population with psoriasis and therefore, even the outcome measure here is irrelevant (i.e., all PIOs irrelevant). Similarly, a8 and a9 included patients with seborrheic dermatitis, whereby the population and intervention were irrelevant, but scaling was an outcome measure (I would suggest somewhat relevant outcome). A10 also involved psoriasis population and intervention, but outcomes included scaling, therefore partially relevant outcome. Overall, the results are entirely inconclusive, since no abstract was relevant.	Expert 3: Overall, no relevant abstracts were available to analyze if trikafta may cause bad flaky scalp or dandruff. This is a very specific claim, and it is usually verified in the side effects results of RCTs. If not asked, participants may ignore the symptom if it is not significant.
Expert 4: There are no relevant abstracts to determine overall support. None of them include cystic fibrosis patients or the medication (Trikafta) from the claim. The outcome is mentioned in abstracts a1, a3, a4, a6, but has no relation to the claim.	Expert 5: The available studies included individuals with dandruff however none were diagnosed with cystic fibrosis, none were given Trikafta (a1, a2, a3, a4, a5, a6, a7, a8, a9, a10). Expert opinion: Dandruff can be caused by the underlying condition (cystic fibrosis) rather than as an effect of the medication itself (Trikafta)	

Table 5: All experts found the claim unverifiable based on the available RCTs. They attributed this to the claim's high specificity, noting it is unlikely—and potentially unethical—for a trial to match the described scenario.

ADHD, Herbs, and Menstruation

Post (r/ADHD): Hello, menstruating people! How do your cycle and ADHD influence each other and how do you deal with it?
 EDIT: After getting your responses I am reflecting again how medicine does not give a shit about women. It's truly insane. Thank you!
 Hello! I have never paid too much attention to my menstrual cycle since it was never particularly bothersome. Now that I take methylo I feel big changes in how I function during the cycle. Like last 10 days of the cycle, my medication kind of stops working... That is like 1/3 of the time. I know it's still better than without meds nevertheless, it makes establishing a routine quite challenging. My doc suggested trying contraceptive pills, but I am not even sexually **active ATM so taking** more medication, with potential side effects, does not excite me.
I know there are herbs that are proven to be helping with regulating the cycle but I don't know if they would help with ADHD symptoms? Any tips?

Population: People with ADHD
Intervention: Herbs
Outcome: Regulating the menstrual cycle

Stimulants and Sodium

Post (r/ADHD): Stimulants vs. Sodium
 Im wondering if anyone else has experienced this. **I find that my stimulant medications (Adderall IR and Vyvanse) make me very sensitive to salt.** If I have a higher sodium meal (eg ramen or canned soup, or even just mustard on my sandwich), I get very bloated. Its uncomfortable and lasts for a few days. Whenever I take a break from my meds, this doesnt happen. Ive had labs done for it in the past and it doesnt seem like anything medically problematic, but its uncomfortable and it really stresses me out.

Population: People with ADHD
Intervention: Stimulant medications (Adderall IR and Vyvanse)
Outcome: Sensitivity to salt (resulting in bloating)

Dandruff and Trikafta

Post (r/CysticFibrosis): Anyone had really **bad flaky scalp** or **dandruff** lately ? **Think it could be due to trikafta.** Could it be anything else

Population: Patients with Cystic Fibrosis (implied by the subreddit r/CysticFibrosis)
Intervention: Trikafta (a medication)
Outcome: Flaky scalp or dandruff

Pineapple Juice Reduces Inflammation

Post (r/CysticFibrosis): Anyone with sinus issues drinking pineapple juice?
 It's a weird question, but **I saw a post about pineapple juice being good for sinus issues (helps with the inflammation)** and just wondered if anyone has done this? Some people were commenting about the high sugar content in pineapple juice not being good, but they get around that by taking a supplement instead of drinking the juice. Anyone?

Population: Patients With Cystic Fibrosis
Intervention: Pineapple Juice
Outcome: Reduced Inflammation/Fewer Sinus Issues

Trikafta and PMDD

Post (r/CysticFibrosis): Trikafta & PMDD
 So, **I believe trikafta has given me PMDD premenstrual dysphagia disorder.** Every month, the week before my period I have extreme **anxiety** in a running dialogue in my head that is constantly negative. I've never been this way before. I also have horrible **hormonal acne** on my back & forehead which are very new to me as well.
 My question is: any one else having this problem? My Dr said they are noticing a "negative interaction with estrogen and trikafta". Anyone find anything that helps??

Population: Patients with cystic fibrosis (implied by the Reddit thread r/CysticFibrosis), specifically females of reproductive age
Intervention: Trikafta
Outcome: Development of PMDD (premenstrual dysphoric disorder, not dysphagia disorder) symptoms, including extreme anxiety and hormonal acne.

Table 6: All claims from the final annotation split. The extracted claim span is highlighted in yellow; Population annotations are highlighted in blue, Intervention in pink, and Outcome in green.

C Rationale for Re-extracting PIO Elements

During the pilot and refinement rounds, we identified a key source of expert disagreement arising from underspecified focus when multiple Population, Intervention, and Outcome (PIO) elements were present in a claim. In such cases, experts differed on whether to consider all PIO elements or prioritize a subset, leading to inconsistent judgments. An illustrative example is shown in Table 7.

Although RedHOT provides PIO highlights for each claim, we observed occasional errors and inconsistencies in these annotations (e.g., mislabeling a Population as an Intervention), which could further contribute to ambiguity during expert review. To reduce this source of disagreement and standardize expert attention, we re-extracted explicit PIO elements from each post using an LLM-based pipeline, with prompt design and expert validation detailed in Appendix I.

Prednisone

Post (r/lupus): Cytoxan and prednisone

Rheumatologist says **celcept** failed to protect my kidneys and now I have developed **lupus nephritis**. Im so upset. **Prednisone** messed up my hips so badly that they both need to be replaced I dont want to get back on it but **rheumatologist says its to bring the inflammation down in my kidneys**. Ive never been on **Cytoxan** but the side effects sound identical to a lupus flare. How am I supposed to be positive with news like this? I feel so defeated I dont know what to do.

Expert 1: Overall consensus of relevant abstracts is that the combination of prednisone with cyclophosphamide (Cytoxan) is effective in treating kidney inflammation due to lupus nephritis (a3, a6, a8, a9, 10). This is in support of the original claim. However, one abstract found that kidney function continues to gradually deteriorate even with treatment (a2). Due to the majority of abstracts supporting the original claim however, the conclusion can be made that cyclophosphamide and prednisone combination therapy is effective for decreasing renal inflammation in lupus nephritis.

Expert 2: The overall conclusion is that the claim is partially refuted. Prednisone alone appears to lead to renal deterioration (a2, a3, a6, a7, a8), but in combination with immunosuppressants, can have beneficial effects (a9). Abstracts a4 and a5 were somewhat relevant (did not specifically test prednisone effects). Irrelevant abstracts included a1 and a10. All relevant abstracts included lupus nephritis as the population. The person who made the claim appears to have arthritis plus lupus nephritis, therefore, none of the abstracts reported this exact population.

Expert 3: None of the given abstracts were relevant to verify the claim.

None of the studies had a control group for prednisone treatment in lupus nephritis. All groups in all of the given studies received prednisone as a base treatment and compared this to the effects of an additional immunosuppressive drug. Since lupus (-nephritis) is an autoimmune disease, it is usually (depending on the severity) treated with immunosuppressive glucocorticoids such as prednisone to inhibit the autodestruction of tissues and organs.

Expert 4: Overall, the abstracts partially support the use of prednisone to reduce kidney inflammation. With the exception of study a7, every other study included either prednisone or glucocorticoid in both the treatment and control groups. The difference usually is between glucocorticoid only or low-dose. Even a7, the only one that does not show a low dose of glucocorticoid use, might not appear to do so because methods may not be totally revealed in the abstract. Therefore, the abstracts suggest that this patient might receive a prescription for at least low-dose prednisone.

Expert 5: Core: The evidence supports the benefit of immunosuppressive medications such as Cyclophosphamide in addition to steroids and oral maintenance medications for those with lupus nephritis. Addendum: There is no study to support the superiority of Cyclophosphamide over other immunosuppressive medications especially if the patient has already had a poor response to other immunosuppressive medications such as Mycophenolate mofetil.

Expert 6: Abstracts a9, a3 a10, a6, a8 conclude that combination of cytoxan and glucocorticoids shows better outcomes in patients with nephritis related to systemic lupus, which are two of the meds mentioned in the claim. The abstracts a7, a5, a2 are somewhat relevant. Since the majority of the abstracts did mention better outcomes using the medications from the claim, but cytoxan wasn't the only immunosuppressive drug compared in the studies, I'd say overall the abstracts partially support.

Table 7: Example illustrating expert disagreement when multiple PIO elements are present and no explicit guidance is provided on which elements to prioritize.

926 **D Changes between annotation rounds**

927 In the first round, ten claims were annotated with-
928 out explicit PIO contextualization or the expert
929 support field. In the second round, we refined
930 the annotation guidelines to clarify label defini-
931 tions; however, agreement on the refinement set of
932 five claims remained low. Based on expert feed-
933 back, we substantially revised the setup for the third
934 round by filtering out non-RCT-verifiable claims,
935 re-extracting and providing PIO elements to guide
936 expert focus, improving the retrieval system, and
937 further clarifying the annotation guidelines. Ex-
938 perts then annotated five claims under this updated
939 setting.

940 **E Annotation Guidelines**

941 We present the guideline given to our expert annota-
942 tors below. These instructions were given in slide-
943 deck format to annotators with images from the
944 annotation interface spliced in-between to clearly
945 indicate how to annotate.

946 **The Task**

947 This annotation task involves verifying medical claims made on Reddit posts using retrieved
948 evidence. You will be looking at the provided abstracts to determine whether, when
949 considering all the evidence, you can support or refute the claim.

950 **Post**

951 We will give you a Reddit post, which is annotated to contain the following.

- 953 • What subreddit the post is from.
- 954 • Spans indicating PIO (Population, Intervention, Outcome) elements.
 - 955 – Population indicates the affected subjects (ex: COVID patients, diabetics).
 - 956 – Intervention indicates any treatments applied to the subjects (ex: remdesivir,
957 Ozempic).
 - 958 – Outcome indicates how the effects of the intervention are evaluated (ex: pain,
959 weight, 30 day mortality).
- 960 • Claim Span: Part of the post that makes the medical claim that we analyze.

961 **Post & Derived Claim**

- 962 • You will NOT be directly evaluating the information in the post. It is presented to you
963 as to inform you of the context in which the claim is made.
- 964 • What you will be directly evaluating is the claim derived from the post. We present
965 you with a (P, I, O) tuple extracted from the post that we use to make the claim as clear
966 as possible.
- 967 • In some cases, the claim in the post may be ambiguous. In this case, we will present a
968 disambiguated claim for you to evaluate.
969

970 **RCT-Verifiability**

- 971 • A claim is RCT-Verifiable if there exists (or should exist) a reasonable RCT that will
972 be able other either support or refute it.
 - 973 – A reasonable RCT is one that can be practically and ethically conducted.
- 974 • Most of the claims we give should be RCT-Verifiable. However, this may not always
975 be the case.
- 976 • In the case when a claim is not RCT-Verifiable, you should indicate as such.
 - 977 – You will be forced to write a 10 word explanation for why the claim is not
978 RCT-Verifiable. Please ensure that the explanation is for a legitimate reason,
979 for example, an unethical intervention, as we will review. You can only
980 continue to annotate once you are done.

981 **Retrieved Abstracts**

- 982 • For each claim, you are given 10 abstracts that are retrieved automatically based on
983 information in the claim.

- Each abstract has the following: 984
 - Title 985
 - Published Date 986
 - Informative Highlights: PIO Spans and Abstract Punchline (Span describing
the core of the abstract's findings) 987
988

- We provide you a way to flag abstracts that you believe to be of poor quality in the
interface. Be sure to keep in mind the quality of the RCT experiment described in the
abstract when annotating them. 989
990
991

992 **Relevance Annotations**

- For these annotations, you will be analyzing the relevance of the abstract to the claim
being evaluated. 993
994
- You will analyzing the relevance of the following four components: 995
 - Population: Is the population being studied in the abstract relevant to the
population the claim is addressing? 996
997
 - Intervention: Is the intervention being studied in the abstract relevant to the
population the claim is addressing? 998
999
 - Outcome: Are any of the outcome measures used in the RCT described in
the abstract relevant to the population the claim is addressing? 1000
1001
 - Overall: Is the abstract relevant enough to the claim for it to be used to verify
the claim? 1002
1003

- As mentioned before, you may flag the abstract if you think it is of concerning quality. 1004

1005 **Relevance Labels**

- For each relevance component, you are given 4 labels to choose from. They are as
follows: 1006
1007
 - Select (Default) 1008
 - * For PIO: The element is missing. 1009
 - * For Overall: The abstract does not describe an RCT. 1010
 - Irrelevant 1011
 - * For PIO: The element in the abstract has no relation at all to the
corresponding element in the post. 1012
1013
 - Ex for Population - Claim: Patients with LPR - Abstract: Healthy
Patients 1014
1015
 - * For Overall: No part of this abstract can be used to make even an
inference on whether the claim can be supported or refuted. 1016
1017

- Somewhat Relevant 1018

- For PIO: Indicates that the element has some relation to the corresponding
element in the post, but is not close enough for it to be used to directly verify
the claim even if all other elements are 100% relevant. 1019
1020
1021
 - * Ex for Population - Claim: Patients with LPR - Abstract: Patients with
gastro-oesophageal reflux disease 1022
1023
- For Overall: Some parts of this abstract can be used to make an inference on
whether the claim can be supported or refuted. However, this abstract still
cannot be used as direct evidence to support or refute the claim. 1024
1025
1026

- Relevant 1027

- For PIO: Indicates that the element in the abstract is close enough to the
corresponding element in the post for the purpose of verifying the claim. 1028
1029
 - * Ex for Population - Claim: Patients with LPR - Abstract: Patients with
laryngopharyngeal reflux 1030
1031
- For Overall: The abstract can be used as direct evidence to support or refute
the claim. 1032
1033

1034 **Abstract Support**

- If overall, the abstract is relevant to the claim. You will be given the opportunity to
annotate for whether the abstract supports/refutes the claim. 1035
1036

- There are four labels you can choose from: 1037

- Refutes: This abstract fully refutes the claim in the post. 1038
- Partially Refutes: This abstract refutes the claim given some condition or
caveat. 1039
1040
- Partially Supports: This abstract supports the claim given some condition or
caveat. 1041
1042
- Supports: This abstract supports the claim in the post. 1043

- A partial support or refute indicates that there is some nuance in the RCT result that
prevents the abstract from fully supporting or refuting the claim. 1044
1045

- A sub-group of the population experienced different results from the rest. 1046
- The results can only be reproduced under specific conditions that cannot be
generalized. 1047
1048

1049 **Relevant Span**

1050 When you are done determining the support label for abstract. You must determine which
span of text in the abstract is most relevant in indicating whether the abstract can be
supported or refuted. 1051

1052 **Tiering**

1053
1054

- 1055
1056
- After you are done with annotating all the abstracts you can start the tiering and synthesis.
- 1057
- In the tiering phase, you are organizing the abstracts into tiers.
 - Abstracts are automatically tiered according to the relevance annotations.
 - You should attempt to further categorize the abstract according to their quality or their importance regarding the claim, as well as temporal relevance (up your own medical expertise).
- 1058
1059
1060
1061
- Synthesis**
- 1062
- For this task, you must pick the label determining whether the claim is supported or refuted according to two criteria.
 - Overall Support (OS): Determine whether the claim is supported or refuted using only the provided evidence.
 - Expert Opinion (EO): Determine whether the claim is supported or refuted using your expert knowledge.
- 1063
1064
- 1065
1066
1067
1068
- You are given 6 labels to choose from:
 - No Relevant Abstracts/No Expert Opinion:
 - * OS: There are no relevant abstracts to determine overall support.
 - * EO: You don't have the expert knowledge in the field to make this decision.
 - Refutes:
 - * OS: Overall, considering all the abstracts, there is strong evidence that the claim can be refuted.
 - * EO: According to your expert knowledge, this claim can be strongly refuted
 - Partially Refutes:
 - * OS: Overall, considering all the abstracts, there is evidence that the claim can be refuted depending on some general condition or caveat.
 - * EO: According to your expert knowledge, this claim can be refuted depending on some general condition or caveat.
 - Inconclusive:
 - * OS: This should rarely happen. Only pick this in cases, where there is true deadlock within the evidence as to whether the claim can be supported or refuted.
 - * EO: According to your expert knowledge, there is no scientific consensus that points to the claim being supported or refuted.
 - Partially Supports:
 - * OS: Overall, considering all the abstracts, there is evidence that the claim can be supported depending on some general condition or caveat.
 - * EO: According to your expert knowledge, this claim can be supported depending on some general condition or caveat.
 - Supports:
 - * OS: Overall, considering all the abstracts, there is strong evidence that the claim can be supported.
 - * EO: According to your expert knowledge, this claim can be strongly supported.

1070
1071
1072
1073
1074
1075
1076
1077
1078
1079
1080
1081
1082
1083
1084
1085
1086
1087
1088
1089
1090
1091
1092
1093
1094
1095
1096
1097
1098
1099
1100
1101

Synthesis Explanation

- 1102
- Afterwards, write an explanation of why you picked that option. Use the tiers you created earlier to help develop this explanation (quality, temporal, relevance). Cite the abstracts in your explanation.
 - For example: There were a few abstracts that refuted the claim (a2, a5). . . .
- 1103
1104
1105
- 1106
- Give sufficient context and details where someone can follow the reasoning without looking at your annotations. Think of the perspective of you explaining to a patient.
- 1107
1108
- 1109
- Try and keep the grade level at around Middle School, try to avoid complex jargon
- 1110
- We would like for you to include the following in your explanation:
 - Main statement explaining why you selected the label.
 - Rundown of how relevant evidence (abstracts) supports/conditionally supports/refutes claim.
 - (optional) Addendum with relevant clinical experience regarding claim
 - Use the terminology of a(abstract number) like a9 to refer to abstract 9 in your explanations.
 - Aim for the length (without addendum) to be around 100 words. If you don't need that much explanation, 50 words is fine. If you really need to explain something in more detail, please keep it under 150 words.
- 1111
1112
1113
1114
1115
1116
1117
1118
1119
- 1120
- Please do not include:
 - Any direct references to the tiers (Ex: The abstracts in tiers 1).
- 1121

- F Plain Language Explanation Guideline** 1122
- Include an overall sentence either at the beginning or end of your synthesis explanation. 1123
1124
 - Target to aim the explanations at 100 words or less, 150 words if there are details that must be elaborated on. 1125
1126
1127
 - Include details of abstracts identified as relevant and explanations of how it supports the ultimate label, including some nuance. 1128
1129
1130
 - (Optional) Medical Addendum at end. 1131

G Annotation Interface

RedHOT Post r/Sinusitis

After your balloon sinuplasty, how long did it take you to get rid of your face pressure and shallow breathing?
Its been 5 days since my procedure and I still have facial pressure. This effects me in negative ways like not being able to focus with my eyes. I also have trouble automatically deep breathing. Im just curious how long does it typically take to recover from these symptoms? I was told I should experiencing some relief after 2 days but still nothing. Im starting to wonder if this is dental related.

Selected Claim Not RCT-Verifiable

Original Claim: I was told I should experiencing some relief after 2 days but still nothing.
Population: Patients who have undergone balloon sinuplasty (similar to the author of the post, who has sinusitis)
Intervention: Balloon sinuplasty
Outcome: Relief from facial pressure and breathing symptoms (specifically within 2 days after the procedure)

Retrieved Abstracts Close

a1 a2 a3 a4 a5 **a6** a7 a8 a9 a10

Relative importance of antibiotic and improved clearance in topical treatment of chronic mucopurulent rhinosinusitis. A controlled study.

1986 Aug 16

50 patients with chronic mucopurulent rhinosinusitis were randomly allocated to treatment with nasal sprays of dexamethasone, tramazoline, and neomycin, dexamethasone and tramazoline with no antibiotic, or matched placebo (propellant alone) four times daily to both nostrils for 2 weeks. The patients were assessed in a double-blind manner for symptomatic response and improvement in nasal mucociliary clearance, nasal airway resistance, sinus radiographs, and intranasal bacteriology and appearance. Both active preparations (with antibiotic 14 of 20 patients responded; without antibiotic 12 of 20 patients responded) were more effective than the placebo (2 of 10 patients responded). There was no significant difference in response between the active preparations with and without antibiotic. Thus, in treatment of chronic mucopurulent rhinosinusitis, reduction of the inflammatory response and decongestion make topical antibiotic unnecessary, probably by allowing host clearance mechanisms to recover.

Annotations: F

Relevance

Population:

Intervention:

Outcome:

Overall:

Optional Comments

Comments

Figure 3: Presentation of claims, PIO elements, and abstracts in the annotation interface.

The interface displays four tiers of abstracts:

- T1:** Contains abstracts a5, a6, and a10. Label: All Relevant Abstracts.
- T2:** Contains abstracts a2, a4, and a8. Label: All Somewhat Relevant Abstracts.
- T3:** Contains abstract a3. Label: Irrelevant Abstracts containing at least 1 non-irrelevant PIO element.
- T4:** Contains abstracts a1, a7, and a9. Label: Completely Irrelevant Abstracts.

Below the tiers are two buttons: "Add Tier" and "Stop Sorting".

Synthesis Annotations:

- Overall Support: N/A
- Expert Opinion: No Expert Opinion
- Overall Explanation: (Empty text area)
- Word Count: 0

Figure 4: Presentation of the tiering and synthesis annotations interface.

We used a web-based annotation interface to collect annotations from our expert medical annotators. Figure 3 shows how we present a claim with its surrounding context, extracted PIO elements for that claim, and retrieved abstracts corresponding to that claim. In this setup, we highlight the claim within the post in which its found along with any PIO spans as determined by data found in the RedHOT dataset. We also present extracted PIO elements (see Appendix I) in a separate box with a rewritten claim created by inputting these elements in a template. All of this information was provided for the benefit of the annotator to clearly understand the claim in question.

Reading each abstract can be a cumbersome task. Therefore, we also provide for each abstract, information highlights and the published date of the paper associated with that abstract. These informational highlights covered the PIO elements in the abstract as well as the punchline of the abstract. This information is provided along with the abstracts in the TrialStreamer dataset.

Figure 4 shows the interface in which the expert annotators would tier abstracts and then provide the overall synthesis annotations with explanations. After the annotator is done with their relevance annotations, the interface will automatically tier abstracts according to these annotations. These automatic tiers are:

- **All Relevant Abstracts:** This tier contains all abstracts that were determined to be overall relevant to the claim.
- **All Somewhat Relevant Abstracts:** This tier contains all abstracts that were determined to be somewhat relevant to the claim.
- **Irrelevant Abstracts containing at least 1 non-irrelevant PIO element:** As the name suggests, this tier contains all irrelevant abstracts with at least 1 non-irrelevant PIO element in relation to the claim.
- **Completely Irrelevant Abstracts:** This tier contains all abstracts with the overall and all the PIO elements labeled as irrelevant in relation to the claim.

Expert annotators, when presented with these tiers, should try to further categories the collection of abstracts if possible. They are able to add tiers, manipulate their order, and change their names. They can also double click on an abstract tag, and

the interface will display the abstract corresponding to that tag. All of these features serve to make the process of tiering abstracts as streamlined as possible for the expert annotators.

H Annotator Recruitment

We recruited five medical experts via Upwork over a four-week period. We received 117 proposals, screened 19 candidates using a sample annotation task (Table 8) to assess instruction-following, medical knowledge, and explanation quality, and short-listed seven candidates for interviews, from which five were selected. Experts worked between 3 and 20 hours per week and required approximately 20 minutes to annotate each claim end-to-end. Compensation ranged from \$22 to \$35 per hour. The total cost of the annotation study was \$1,432.20.

I PIO Extraction

We used an automatic PIO extraction approach based on Llama-3.1-405B-Instruct. After preliminary prompt testing and consultation with a medical expert, we incorporated PIO element definitions and guiding questions adapted from Duke University’s PICO evidence-based medicine guidelines¹.

PATIENT OR PROBLEM
How would you describe a group of patients similar to yours? What are the most important characteristics of the patient? Example: COVID patients, diabetics

INTERVENTION, EXPOSURE, PROGNOSTIC FACTOR
What main intervention, exposure, or prognostic factor are you considering? What do you want to do with this patient? Example: Remdesivir, Ozempic

OUTCOME
What are you trying to accomplish, measure, improve or affect? Example: pain, weight, 30 day mortality

Extract the Population, Intervention, and Outcome elements from the following claim from the following text. Write “None” if the element does not exist in the text.

Text posted by someone in Reddit thread r/[sub][sub_description]:
[post]
Highlighted claim:

¹<https://guides.mclibrary.duke.edu/ebm/pico>

Gaviscon Advance

Post (r/GERD): Can I buy liquid alginate suspension (Gaviscon Advance) in the U.S.? Hi everyone. I'm newly diagnosed with LPR and doing a lot of research on the best treatments. I've read that a liquid alginate suspension (Gaviscon Advance) is quite effective at treating LPR but it looks like it's not sold in the U.S. Does anyone know how I can find it here?

Table 8

1229	<i>[claim]</i>		
1230	On a validation set of 55 samples (five claims	model on a small development set, with relevance	1272
1231	randomly sampled from each of 11 conditions),	judgments provided by a medical expert. Aggre-	1273
1232	a medical expert reviewed the extracted PIO el-	gate results are shown in Table 9. At the time of ex-	1274
1233	ements. The expert reported over 90% accuracy	perimentation (09/11/2024), stella_en_400M_v5	1275
1234	for each element, noting only minor issues such	was among the top-performing open-source embed-	1276
1235	as overlapping Population and Intervention spans	ding models on the HuggingFace MTEB bench-	1277
1236	and occasionally implied elements. Based on this	mark while remaining computationally feasible for	1278
1237	review, the extraction pipeline was deemed suffi-	indexing approximately 800,000 RCT abstracts.	1279
1238	ciently reliable for use in our study.		
1239	J Prompt for RCT Verifiability	L Implementation Considerations for the	1280
1240	<i>You are a potential clinical trialist. I will give</i>	Communication Model	1281
1241	<i>you a claim and post. The claim is part of the</i>	We emphasize that the communication model is	1282
1242	<i>post, and the post can give you context. I want</i>	an abstraction, and that the actual system could	1283
1243	<i>you to tell me if the claim can be studied in</i>	be implemented in many ways. The model could	1284
1244	<i>a randomized controlled trial (RCT). An RCT</i>	be a Reddit or Chrome extension tool for users	1285
1245	<i>can test an intervention to measure a benefit or</i>	to interact with while they write their Reddit post.	1286
1246	<i>non-inferiority. However, the RCT must be ethical:</i>	Alternatively, the model could also publicly post	1287
1247	<i>Ethical Guidelines: The intervention should not</i>	follow-up questions as comments as a public reply	1288
1248	<i>cause harm or have a significant risk of toxicity. It</i>	to the users' post. The benefit of this approach is	1289
1249	<i>should not test exposures known to be potentially</i>	that there could be an expert in the loop that veri-	1290
1250	<i>harmful, such as food-drug interactions that might</i>	fies the models' response, as well as these public	1291
1251	<i>cause adverse effects. The safety of participants</i>	responses could be helpful to other users reading	1292
1252	<i>is the primary concern, and interventions that</i>	the thread. We encourage the research community	1293
1253	<i>pose significant health risks should not be tested</i>	to explore various implementations of this system,	1294
1254	<i>in an RCT. Design Requirements: The trial needs</i>	as well as focus on extensive human and expert	1295
1255	<i>to have a control group, with the only difference</i>	evaluation and systematic HITL methods.	1296
1256	<i>being the intervention. There must be a feasible</i>		
1257	<i>and ethical way to measure outcomes without</i>		
1258	<i>exposing participants to undue risk. Wait for my</i>		
1259	<i>text to classify whether the claim can be ethically</i>		
1260	<i>studied in an RCT.</i>		
1261			
1262	<i>Claim: [claim]</i>		
1263			
1264	<i>Text from r/[subreddit]: [post]</i>		
1265			
1266	<i>Format your response starting with Classifica-</i>		
1267	<i>tion: Can be ethically studied in a RCT or Classifi-</i>		
1268	<i>cation: Cannot be ethically studied in a RCT</i>		
1269	K Retrieval Configurations		
1270	We evaluated several retrieval configurations using		
1271	the dunzhang/stella_en_400M_v5 embedding		

Strategy	Query	Document	Pop.	Inter.	Out.	Overall
S2P	Question PIO	PIO	2.2	1.6	2	1.5
S2P	PIO	Abstract	2.3	1.7	2.1	1.7
S2P	PIO	PIO	2.3	1.7	2.1	1.6
S2S	PIO	Abstract	2.4	1.7	2.2	1.7
S2S	PIO	PIO & Abstract	2.2	1.6	2.0	1.6
S2S	PIO	PIO	2.3	1.6	1.9	1.5
S2S	PIO	PIO & Title & Abstract	2.2	1.6	2.0	1.6
S2S	PIO	Title & Abstract	2.3	1.7	2.1	1.6

Table 9: Retrieval strategy results on a test set of five claims. Scores report average relevance per claim, computed over 10 retrieved abstracts (1 = irrelevant, 2 = somewhat relevant, 3 = relevant). The highlighted row indicates the configuration used in the main study.

Post

microclots

Post (r/Diabetes): Could **microclots** help explain the mystery of long Covid? **Acute Covid-19 is not only a lung disease, but actually significantly affects the vascular (blood flow) and coagulation (blood clotting) systems. A connection to the damage done by diabetes might be possible.**

Diabetes

Post (r/Diabetes): Affording Medication

so im on a family plan with a 3k/6k out of pocket expense, I think it's hdhp with a hsa that my husband employer contributes a bit too. I know when I had coverage with my job i had a ppo plan. he's the one that chooses the plans at his job so im not the best when it comes to explaining the details for it..

was originally taking metformin but it's horrible and over the past 2 months it's been making me sick as a dog so I asked my endocrinologist can I go on something else. she recommended **ozempic since alot of patients responded well to it, lost weight, and had a good effect on their sugar.** plus it's taken only weekly in which sounds great for someone like me since I'm not the best with keeping up with medications. back in December since we had hit our 6k deductible I had paid nothing when I recieved the medication so I had no clue what the actual price would be but I nearly had a heart attack when I tried picking it up in the store recently...with my plan I'm at 800 bucks for the thing and optum informed me it's 2300 (1981 with the discount card) for a 90 day supply. that's ALOT of money...I was going to purchase farxiga today with optum (1500 dollars) but literally don't have the money to afford to do so...my car needs a new catalytic converter so financially I had to make the cut to my medication (my cardiologist put me on that to prevent heart failure since I have "resistant hypertension" that's not responding well medication)

i made a joke to my husband and said I may have to divorce him just so I qualify for that government health insurance. hell looking at it now I may be serious! as a diabetic or anyone in America on any type of medication how are ppl able to afford their insulin/pills/machines/ whatever. our household income is around 85k so there's not much assistance we can get that im aware of

Hallucinations

Post (r/narcolepsy): Do people with IH experience hallucinations?

I am so confused! My MSLT showed IH but my doctor gave me a clinical diagnosis of **narcolepsy** because I experience hypnopompic **hallucinations** and **sleep paralysis**. **She told me people with IH dont experience those things which is why she switched the diagnosis.** Im confused because I've read articles that say they are symptoms of IH. I know it doesnt really matter because treatment is the same, but I have this thing in me where I just need to know.

COVID

Post (r/Epilepsy): **Epilepsy Patients Much More Likely to Die of COVID**

Long Covid

Post (r/CFS): **I Had Never Felt Worse: Long Covid Sufferers Are Struggling With Exercise And experts have some theories as to why.** - The New York Times

Glycemic

Post (r/Diabetes): **Dietary carbohydrate restriction augments weight loss-induced improvements in glycaemic control and liver fat in individuals with type 2 diabetes:** a randomised controlled trial. (Pub Date: 2022-01-07)

Pfizer vaccine

Post (r/CysticFibrosis): Pfizer vaccine

My son, non cf, is having his second pfizer **vaccine**. He is 25 yrs old. For some reason I'm really nervous about it as **he has been told not to exercise for 48hrs afterwards due to heart inflammable young people are getting**...obvs this is rare...but my son is extremely active & I'm in a tizz. He's having now as i write this. I'm extremely proud he is having it as alot of youngsters are refusing it atm but the anxiety over it is making me feel sick.

mold

Post (r/rheumatoidarthritis): Mold and RA

I'm having a bit of a weird issue with **mold**. I'm currently in the process of being diagnosed with **RA**. I've got **achy joints**, **swelling** whole nine yards. I transferred job locations earlier this month and was starting to feel better and my hand swelling finally went down. I then signed up for some overtime in my old job location and after about 2 hours my elbows and hands started to **ache** and swell. Every time it rains at this building water runs through the walls. I'm certain theirs mold in the walls. **Google says long-term toxic mold exposure can mimic RA**. Had anyone else had an experience with RA symptoms not ending up being RA or having one large trigger to RA symptoms. After going home and sleeping on things my hands started to feel better but not completely.

Continued on next page

Post

Copaxone

Post (r/MultipleSclerosis): copaxone vs aubagio?

My gf is about to switch from once a day copaxone injections to aubagio at the advice of her new neurologist.

After doing some research before starting the treatment, she is a bit worried about the liver function concerns with the drug.

My gf is bipolar, has high anxiety, and is on several meds for her mental health. I just pulled up a site that compared these 2 drugs and was really angry to see that copaxone patients reported it caused depression, anxiety, and other things the doctor never mentioned. So I am cautiously optimistic that the change is in her best interest.

Any thoughts or experiences would be greatly appreciated. My research seems to lean towards the new medication, but we are obviously concerned at least about the liver function monitoring.

tyvm

Calcium

Post (r/thyroidcancer): Calcium supplements (Citracal slow release) and total thyroidectomy Hey all, I'm almost 6 years post total thyroidectomy, and since my providers at the time of my TT didn't really share any of this info/it was hard to track down, I wanted to put it out there for others.

First- you'll probably want to get on a calcium supplement. That part I was told. How much calcium I wasn't told, but eventually found out from a pharmacist to go a bit above the recommended for your age/assigned sex at birth. Normally my recommended would be 1000mg, but because of the TT, it's 1200mg.

Second- wait at least 4 hours after taking your thyroid hormone replacement before taking a calcium supplement. Also was told that, also something everyone here probably already knows.

Third- our bodies can only absorb around 500mg of calcium in one go.

Fourth- if you've had a TT, your body will absorb calcium citrate more effectively than calcium carbonate. I learned this literally a month ago from a PCP who doesn't specialize in thyroid health, and I'd love to know why my endocrinologist never told me.

Now, all of the above led me to be interested in the Citracal slow release, as it's 1200mg, but released slowly so you can take it once a day and still get all of it. My only issue was that I couldn't find anywhere that said how long it took to fully release. I was worried that if it took too long, it would prevent my levothyroxine from absorbing the next day. I couldn't find the answer online, but finally called their questions line today and found out it's 8 hours.

Obviously I'm not a medical provider, I just want more of us to have access to this info, especially since a lot of us have worked with medical providers that don't give us all of the information we need. Also not here to advertise for that brand specifically, I just wanted something convenient enough to take once a day, and figured others might have had the same question!

Table 10: Pilot claims.

N Refinement Claims

Post

Ivabradine

Post (r/POTS): Anyone here with low bp take Ivabradine?

Im just wanting to do a bit of research on different meds before my doctors appointment. Last time they told me they cant medicate me because my bp dropped pretty low during the TTT. However, at rest my bp is normal and even standing up it doesnt drop noticeably low unless Im standing still for a longer period of time.

So I just wanted to know if any of yall are in a similar situation and have good (or bad) experiences with this drug. I hear midodrine is good for low bp but its more expensive and the side effects sound kind of iffy to me.

Fluoxetine

Post (r/Dysthymia): Long-lasting apathetic tendencies, anhedonia etc.

I'm just apathetic in general, and am unable to do even the smallest things. Fluoxetine might have had some positive effects, and I'm supposed to be taking it now, but I can't even be bothered to get a refill.

I can't tell whether or not my asocial tendencies are a personality trait. I currently have no interest in maintaining a relationship with family or friends.

I've never been diagnosed with dysthymia - only depression - but a lot of the symptoms seem relevant, and my doctor did mention it at one point.

Psychosis and Antidepressants

Post (r/Psychosis): Psychosis and antidepressants

Hey everyone!

So some crazy stuff happened to me over the last week. I am on abilify for my psychosis and I have been suffering from depression.

My doctor decided to prescribe me Wellbutrin 150mg first. Took it for about five days, started having extreme anxiety and dry mouth. I mentioned this to my doctor and he switched me to Lexapro 5mg. Extreme anxiety and dry mouth but something new happened this time -my fucking delusions and hallucinations came back. I had to legit tell myself my thoughts werent based on reality. But holy crap was it difficult. I didnt take any antidepressants today and already feel better.

This is crazy, has anyone experienced anything like this? I didnt think anti depressants would bring out my psychosis. Guess I might have to go the natural route for my depression :(

Prednisone

Post (r/lupus): Cytoxan and prednisone

Rheumatologist says cellcept failed to protect my kidneys and now I have developed lupus nephritis.Im so upset. Prednisone messed up my hips so badly that they both need to be replacedI dont want to get back on it but rheumatologist says its to bring the inflammation down in my kidneys. Ive never been on Cytoxan but the side effects sound identical to a lupus flare. How am I supposed to be positive with news like this? I feel so defeatedI dont know what to do.

Metformin replace Insulin

Post (r/Diabetes): Can metformin replace insulin?

I realize this is definitely case-by-case but Im curious to know if anyone has been able to get off of insulin and take just metformin? When I was diagnosed with type 2 I was automatically put on insulin and generally take quite a bit of it, but now after some research Im considering asking my doctor to try treatment through metformin in February.

Table 11

1299

O License Information

1300

We predominantly used the RedHOT dataset and abstracts from the TrialStreamer database in our work. Both of these works are licensed on a Creative Commons Attribution 4.0 International License.

1301

1302

1303

1304

1305

1306

We will also release our work under a Creative Commons Attribution 4.0 International License.