# What Would it Take to get Biomedical QA Systems into Practice?

**Gregory Kell** and **Iain J. Marshall**
King's College London
{name.surname}@kcl.ac.uk

**Byron C. Wallace**
Northeastern University
b.wallace@northeastern.edu

**André Jaun**
Metadvice
ajaun@metadvice.com

## Abstract

Medical question answering (QA) systems have the potential to answer clinicians' uncertainties about treatment and diagnosis on-demand, informed by the latest evidence. However, despite the significant progress in general QA made by the NLP community, medical QA systems are still not widely used in clinical environments. One likely reason for this is that clinicians may not readily trust QA system outputs, in part because transparency, trustworthiness, and provenance have not been key considerations in the design of such models. In this paper we discuss a set of criteria that, if met, we argue would likely increase the utility of biomedical QA systems, which may in turn lead to adoption of such systems in practice. We assess existing models, tasks, and datasets with respect to these criteria, highlighting shortcomings of previously proposed approaches and pointing toward what might be more usable QA systems.

## 1 Introduction

During consultations in primary care, clinicians generate at least one question for every two patients (Del Fiol et al., 2014). Nonetheless, clinicians look for answers to only half of the questions due to time constraints and the belief that answers to certain questions do not exist (Del Fiol et al., 2014), despite the plethora of available evidence (Bastian et al., 2010). When clinicians do search for answers, they usually spend fewer than three minutes per question doing so (Del Fiol et al., 2014; Hoogendam et al., 2008).

Our focus in this paper is on questions pertaining to patient care decisions, for example seeking guidance about diagnosis or treatment. Ideally, clinicians would search for answers to such questions with reference to high-quality studies and up-to-date evidence syntheses, typically indexed in medical databases such as PubMed[1] and the Cochrane Library.[2] This practice of emphasizing use of rigorous empirical evidence is known as *evidence-based medicine* (EBM). Under this framework, evidence compiled from *all relevant high-quality research* (in the form of, e.g., systematic reviews and rigorously produced clinical guidelines) is preferred to individual studies or expert opinion (Ebell et al., 2004; Guyatt et al., 2008; Alper and Haynes, 2016a). Unfortunately, searching existing sources for relevant, high-quality information is onerous. Due to the time constraints imposed on clinicians, this leads to widespread reliance on general information sources such as Google (Hider et al., 2009). However, while simple to use, general-purpose search engines rank results according to criteria not directly aligned with EBM principles such as rigour, comprehensiveness, and reliability (Hider et al., 2009).

Aside from internet search, clinicians often engage in informal discussions about decisions with colleagues in what are sometimes referred to as "curbside consultations" (Papermaster and Champion, 2017, 2020; O'leary and Mhaolrúnaigh, 2012). It is common for practitioners to engage in at least one such discussion per week for practical reasons, including convenience, or an urgent need for information (Smith, 1996). These inform the "mindlines" that clinicians acquire over their careers (i.e., mental models of medicine) and that are also based on other sources including guideline documents, training, background reading, and experience (Gabbay and le May, 2016). However, the information exchanged in informal consultations may be inaccurate, incomplete, and lead to practice influenced more by expert opinion than the scientific literature (Papermaster and Champion, 2017).

Medical question answering (QA) systems have the potential to address these issues by answering clinicians' questions in real-time on the ba-

---

[1] https://pubmed.ncbi.nlm.nih.gov

[2] https://www.cochranelibrary.com

sis of the latest evidence. This has motivated development of QA systems and associated medical QA datasets used to train them. For example, BioASQ (Tsatsaronis et al., 2015) and PubMedQA (Jin et al., 2019) have been created to train and evaluate systems that answer clinicians' questions based on medical research literature, while emrQA (Pampari et al., 2018), emrKBQA (Raghavan et al., 2021) and why-QA (Fan, 2019) were constructed using queries concerning patient data from electronic health records (EHRs). MEDIQA-QA (Ben Abacha et al., 2019) and LiveQA-Medical (Abacha et al., 2017) are datasets designed for systems that answer consumer (patient) queries. MEDIQA-AnS (Savery et al., 2020) accompanies the answers from MEDIQA-QA with summaries that consumers would understand more easily. Systems for QA over EHRs aim to answer questions about the medical history or prior care of individual patients. By contrast, our focus here is on systems that can provide general evidence-based guidance in response to queries; we therefore omit emrQA, emrKBQA and why-QA from our discussion.

Existing biomedical QA systems that answer questions with reference to the medical literature typically provide answers in the form of yes/no, factoids, lists, and/or definitions (Sarrouti and Ouatik El Alaoui, 2020; Ben Abacha and Zweigenbaum, 2015; Cao et al., 2011; Zahid et al., 2018; Yu et al., 2007) without supplying justifications, e.g., source journals, extracted text snippets, and/or associated statistics. However, this answer format does not readily translate into clinical practice.

Take, for example, the question "Which antibiotic should I use for urinary tract infections?". A factoid-based QA system might (reasonably) return the answer "trimethoprim 200mg". However, a "correct" answer is not sufficient to translate into clinical use. An answer here is only as reliable as the source from which it was extracted. The source therefore needs to be judiciously chosen, and presented transparently. Furthermore, in this example, the knowledge of the best treatment requires information about the patients' age and any additional health problems (for instance, dosing may vary in children, or where someone has impaired kidney function). The optimal treatment might vary by location, reflecting local or individual bacterial resistance patterns (which frequently change over time), or vary depending on the cost of drug acquisition or availability. A factoid answer

does not allow the possibility of changing practice, or providing critical information which is not a direct response to the narrow question asked (perhaps an antibiotic is not always needed). These issues both need to be considered in producing an answer, and need to be *seen* to have been considered by the clinician before s/he can feel confident in following the recommendation.

In this context, reliability is multi-faceted. For an answer to be reliable it must have been extracted from a trustworthy source, accurately transcribed, and relevant to the clinical context (was the dosing information extracted for the correct clinical condition?). It should also be locally applicable, and recent. In this example, a methodologically sound national clinical guideline is likely to be highly dependable, whereas a journal editorial or case study giving one expert's idiosyncratic opinion might be safely ignored. A question-answering system which does not understand the difference is not likely to be useful.

We argue that the deployment of EBM-guided QA systems—by which we mean those intended to provide answers to clinical questions based on published evidence—in clinical practices is contingent on the outputs being reliable and actionable. Clinicians should be able to trust that the most robust evidence was retrieved, and that conflicting evidence was handled appropriately. Uncertainties associated with answers should be communicated to the clinician.



Figure 1: Yellow box shows text snippet used to answer "what dose of flucloxacillin should I prescribe for a 5 year old child?"[3].
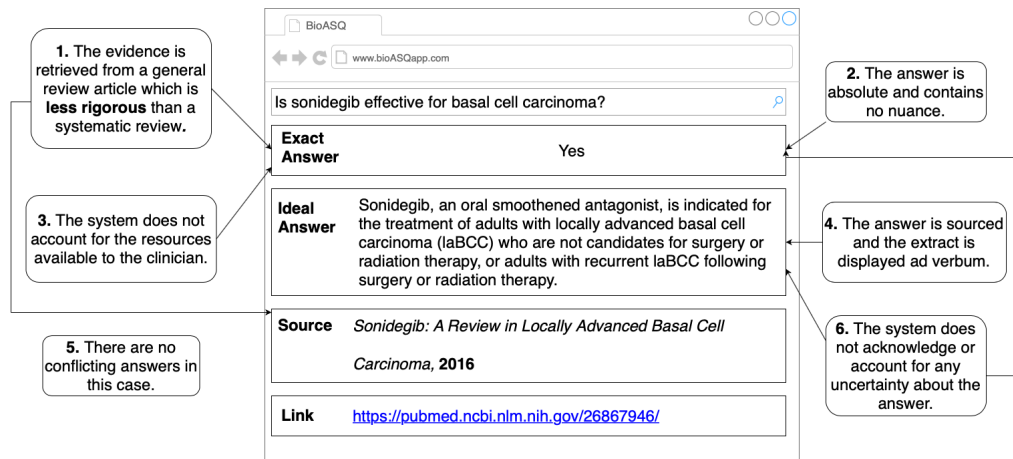
**BioASQ**

www.bioASQapp.com

Is sonidegib effective for basal cell carcinoma?

| Exact Answer | Yes |
|---|---|
| Ideal Answer | Sonidegib, an oral smoothened antagonist, is indicated for the treatment of adults with locally advanced basal cell carcinoma (laBCC) who are not candidates for surgery or radiation therapy, or adults with recurrent laBCC following surgery or radiation therapy. |
| Source | *Sonidegib: A Review in Locally Advanced Basal Cell Carcinoma*, **2016** |
| Link | https://pubmed.ncbi.nlm.nih.gov/26867946/ |

**1.** The evidence is retrieved from a general review article which is **less rigorous** than a systematic review.

**3.** The system does not account for the resources available to the clinician.

**5.** There are no conflicting answers in this case.

**2.** The answer is absolute and contains no nuance.

**4.** The answer is sourced and the extract is displayed ad verbum.

**6.** The system does not acknowledge or account for any uncertainty about the answer.

Figure 2: Web interface for QA system developed using BioASQ.



**PubMedQA**

https://www.PubMedQAapp.com

Do preoperative statins reduce atrial fibrillation after coronary artery bypass grafting?

| Answer | Yes |
|---|---|
| Long Answer | *(Conclusion)* Our study indicated that preoperative statin therapy seems to reduce AF development after CABG. |
| Context | *(Objective)* Recent studies have demonstrated that statins have pleiotropic effects, including anti-inflammatory effects and atrial fibrillation (AF) preventive effects [...]<br><br>*(Methods)* 221 patients underwent CABG in our hospital from 2004 to 2007. 14 patients with preoperative AF and 4 patients with concomitant valve surgery [...]<br><br>*(Results)* The overall incidence of postoperative AF was 26%. *Postoperative AF was significantly lower in the Statin group compared with the Non-statin group (16% versus 33%, p=0.005).* Multivariate analysis demonstrated that independent predictors of AF [...] |

**1.** The source of the answer is unclear.

**3.** The system does not account for the resources available to the clinician.

**5.** There are no conflicting answers in this case.

**2.** The answer is absolute and contains no nuance.

**4.** The long answer is sourced and the extract is displayed ad verbum.

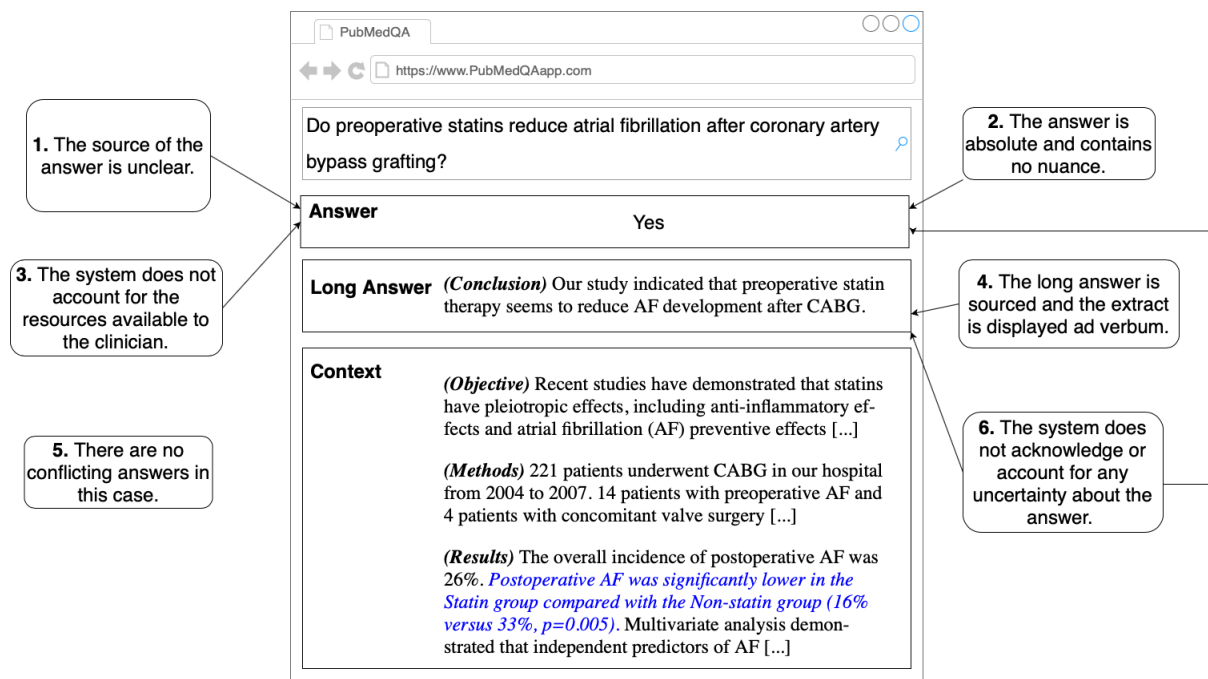**6.** The system does not acknowledge or account for any uncertainty about the answer.

Figure 3: Web interface for QA system developed using PubMedQA.

## 2 Desiderata for Medical QA

What would be needed for clinicians to trust, and actually act upon answers provided by a QA system? In our view, the necessary criteria include: Provenance of the evidence and its reliability; Faithfulness of the evidence to the source, and; Transparency with respect to how answers are chosen, and how conflicting evidence is resolved. In accordance with these criteria, we suggest the following questions to assess the transparency of QA systems:

1. **Do the answers come from reliable sources for health information?** All research articles are not equal, and there exist mature approaches to help clinicians identify the most reliable advice from the health literature. Evidence-Based Medicine is one such framework in which the findings of the most rigorous study designs (typically high quality clinical guidelines, and *systematic reviews* of the primary literature) are preferred to case studies and observational research (Alper and Haynes, 2016b; Sackett et al., 1985).

More sophisticated approaches (e.g., *risk of bias* assessment tools and the *GRADE* framework; Higgins et al. 2011; Guyatt et al. 2008) go further by estimating how confident one should be in a research finding, taking into account aspects such as study type, the precision

**Figure 4 (diagram):**

Left callouts:
- **1.** The evidence is retrieved from a clinical guideline which is one of the most **reliable** sources of evidence**.**
- **3.** The clinician would be able to perform a DEXA scan in their practice.
- **5.** There are no conflicting answers in this case.

Browser window — Medical QA — https://www.medicalqa.com

Should men receiving long-term GnRH analogues for prostate cancer be offered regular DEXA scans to monitor potential loss of bone density?

**Result** No relevant local or national guidelines are available, but here is one from Wirral Community Teaching Hospital:

*A DEXA scan should be performed by secondary care physician in patients with major risk factors for decreased bone mineral content and treatment should not be initiated if result is below normal levels.*

**Source** *Gonadorelin analogue and gonadotrophin-releasing hormone (GnRH) antagonist depots for treatment of prostate cancer, 2020*

**Link** https://mm.wirral.nhs.uk/document_uploads/shared-care/GonadorelinandGnRHantagonistsfortreatmentofprostatecancerSCGV1.pdf

Right callouts:
- **2.** The answer is framed as a suggestion.
- **4.** The answer is sourced and the extract is displayed ad verbum.
- **6.** The system acknowledges uncertainty with regards to the applicability of the Wirral guidelines.
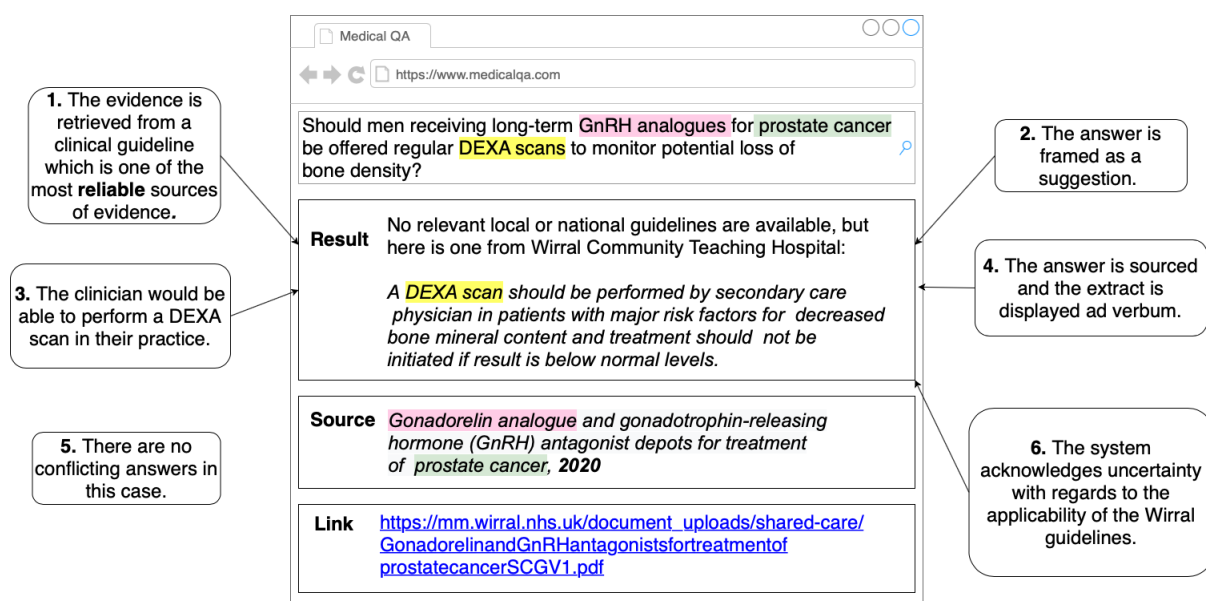
Figure 4: Example of a medical QA system output that meets the criteria in §2. The assumption is that the clinician is based in Nottingham while the most relevant guideline is for Wirral Community Teaching hospital (which is in a different region). The corresponding text spans in the question and response are highlighted with the same color.

**Figure 5 (diagram):**

Left callouts:
- **1.** The evidence is retrieved from systematic reviews which are among the most **reliable** sources of evidence.
- **3.** The clinician would need to devote time to investigate the listed reviews.
- **5.** Given the conflicting evidence, the system refrained from providing a recommendation. Instead, the relevant reviews were listed.

Browser window — Medical QA — https://www.medicalqa.com

Should spinal manipulations be used to treat headaches?

**Result** 3 systematic reviews were found, but their conclusions are contradictory.

**Source** Manual therapies for migraines: a systematic review, **2011**
Spinal manipulations for cervicogenic headaches: a systematic review of randomised clinical trials, **2011**
Evidence-based guidelines for the chiropractic treatment of adults with headache, **2011**
Conflicting Snippets

**Link** https://www.ncbi.nlm.nih.gov/pmc/articles/PMC3072494/
https://pubmed.ncbi.nlm.nih.gov/21649656/
https://pubmed.ncbi.nlm.nih.gov/21640251/

Right callouts:
- **2.** The answer is framed as an explanation and a recommendation.
- **4.** No specific guidance about spinal manipulations could be given due to conflicting evidence.
- **6.** The system acknowledges uncertainty with regards to the correct answer.
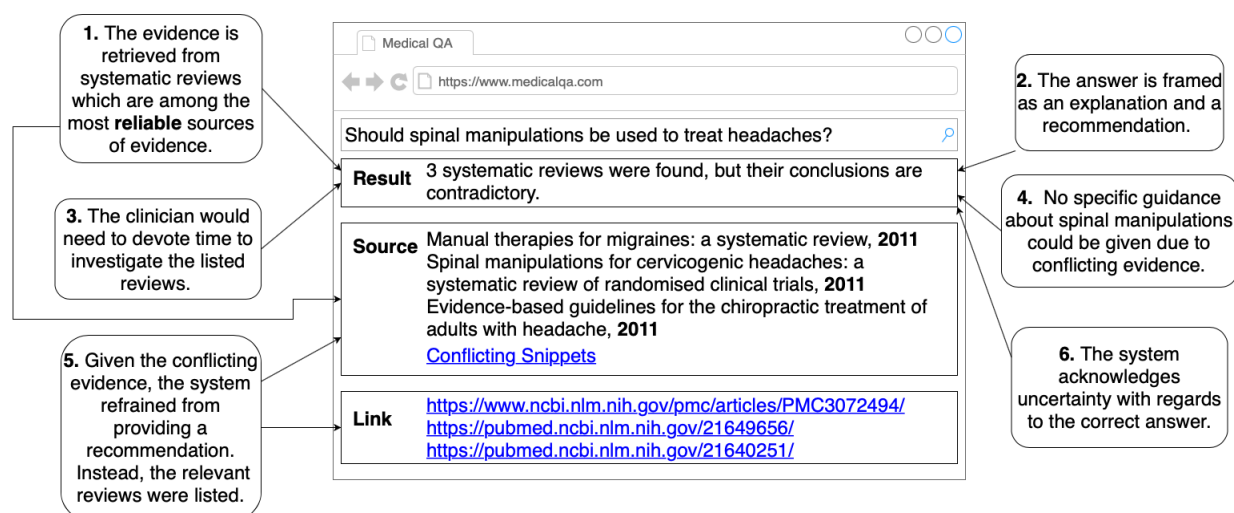
Figure 5: Example of a medical QA system output that handles the conflicting conclusions of 3 systematic reviews.

of the statistical results, and whether problems in study design were likely to have led to bias. QA systems which take a naive approach to evidence extraction—for example, selecting an answer from an undifferentiated corpus of scientific literature, treating all studies as equally reliable—are likely to be considerably less useful to clinicians. This is particularly true because there is often no definitive "correct" answer to a query; an overview of the best available evidence is what is sought. We suggest that QA systems should aim to explicitly use more rigorous, theoretically informed approaches to sorting the literature, mirroring the best current practice of manual question answering and evidence synthesis.

2. **Does the system provide guidance?** When searching for answers clinicians are looking for *guidance*, not just information. Guidance consists of recommendations of what to do in various clinical situations, while Boolean or factoid answers appear more absolute. The demand for guidance is reflected by the fact that many questions are of the form "Should I ...?" (Del Fiol et al., 2014; Ely et al., 2000; Papermaster and Champion, 2017). Therefore, the system could respond with "study/review

**Spinal manipulations for cervicogenic headaches: a systematic review of randomized clinical trials, 2011**

The results are mixed and the only trial accounting for placebo effects fails to be positive. Therefore, the therapeutic value of this approach remains uncertain.

**Manual therapies for migraine: a systematic review, 2011**

Therefore, any firm conclusion will require future, well-conducted RCTs on manual therapies for migraine.

**Evidence-based guidelines for the chiropractic treatment of adults with headache, 2011**

Evidence suggests that chiropractic care, including spinal manipulation, improves migraine and cervicogenic headaches. [...] Evidence for the use of spinal manipulation as an isolated intervention for patients with tension-type headache remains equivocal.

Figure 6: Contradictory source snippets leading to the response presented in figure 5.

X suggests the following action... ". This response could encourage the clinician to engage with the guidance and think critically about how to apply it in practice.

In the aforementioned example on urinary tract infections (UTI), the NICE[4] guideline (NICE, 2019) recommends Nitrofurantoin under specific conditions: If the estimated glomerular filtration rate (eGFR) $\geq$ 45 ml/minute then 100 mg modified-release twice a day (or if unavailable, 50 mg four times a day) for 3 days.

3. **Are the answers useful in the context in which the provider is practicing?** The usefulness of a QA system could be limited by factors such as drug availability, antibiotic resistance, and local or national funding/resources. Therefore, QA systems should account for the resources that are available to

[4]The National Institute for Health and Care Excellence: the UK national health guideline producer

clinicians when providing guidance. In addition, what is deemed as "best practice" may vary by location (i.e., region or country).

If a clinician were to consult a QA system on whether "men receiving long-term GnRH analogues for prostate cancer should be offered regular DEXA scans to monitor potential loss of bone density", guidelines from Wirral Community Teaching Hospital might be retrieved. The clinician would need to decide whether the guidelines apply to their locality (e.g., Nottingham) where DEXA scans may or may not be readily available.

4. **Is there sufficient "rationale" for the answer provided?** Prior work has shown that users of QA systems prefer answers to consist of paragraph-sized chunks of text as opposed to concise phrases (Lin et al., 2003). Lengthier "answers" provide context, and allow users to ensure that the information in the source text is consistent with the final answer. As answers should be faithful to the source, any generated summaries should probably be extractive rather than abstractive.

For example, the answer to "what dose of flucloxacillin should I prescribe for a 5 year old child?" could consist of the snippet highlighted in Figure 1. However, in cases where the answer is derived from multiple sources it may be necessary to generate a summary.

5. **Does the system resolve conflicting evidence appropriately?** Higher quality information should be prioritized using frameworks for rating the quality of evidence (Ebell et al., 2004; Guyatt et al., 2008; Alper and Haynes, 2016a). If there are conflicts between equally relevant and reliable sources, the system should refrain from providing oversimplified guidance and inform the clinician of the conflicting sources. This could form the basis for further investigation by the clinician.

The query "Should spinal manipulations be used to treat headaches?" could return three conflicting systematic reviews: one concluding that they should (Bryans et al., 2011) and two others that judge the evidence to be inconclusive (Chaibi et al., 2011; Posadzki and Ernst, 2011). A QA system should inform the clinician of these contradictions. An ideal system would assess the relative methodological

quality of the reviews, and present the most rigorous and reliable first.

6. **Does the system handle and communicate uncertainties adequately?** When providing guidance, the system should communicate any sources of uncertainty. If appropriate, the system should abstain from providing explicit guidance (e.g., where information conflicts or where supporting evidence is either absent or of low quality).

   In the case of the regular DEXA scans for men recieving long-term GnRH analogues, the system should communicate its uncertainty on whether the guidelines from Wirral Community Teaching hospital are applicable to the clinician's region.

   Additionally, the question "Does speech and language therapy help dysarthria after a brain injury?" could return no relevant studies (Sellars et al., 2002). It is important that the system explain that the question is unanswerable using the available literature.

There are several research challenges associated with the above criteria. Reframing the QA task will require new datasets which include answers (with accompanying rationales) from trusted sources; rankings by evidence quality; locality and patient contextualizing information; and which incorporate real-world conflicting answers and questions which lack answers. Quantitative measures would need to be created to assess how well the datasets and systems meet each criterion.

While we expect that an improved system using these criteria might be more trustworthy (and hence potentially help to translate health research more effectively info clinical practice), we note that our criteria need to be empirically tested. To achieve this, we need to move beyond dataset evaluation, and consider user-centred design methodology. Ultimately, we should aim to improve and evaluate systems through research conducted in real-world clinical practice.

We next review prior work on Biomedical QA with respect to the above criteria. We display typical responses of these systems in a hypothetical web interface, and assess how well these responses meet the criteria.

# 3 Existing Medical QA Datasets and Systems

The primary focus of prior medical QA work has been on developing systems that answer the following types of questions: boolean (yes/no), factoid, list (of factoids), and definitional, e.g. (Sarrouti and Ouatik El Alaoui, 2020; Ben Abacha and Zweigenbaum, 2015; Cao et al., 2011; Zahid et al., 2018; Yu et al., 2007). Several datasets have been created to train and evaluate systems that handle the aforementioned question types, including BioASQ (Tsatsaronis et al., 2015), emrQA (Pampari et al., 2018), emrKBQA (Raghavan et al., 2021), PubMedQA (Jin et al., 2019), why-QA (Fan, 2019), MEDIQA-QA (Ben Abacha et al., 2019), LiveQA-Medical (Abacha et al., 2017) and MEDIQA-AnS (Savery et al., 2020). BioASQ, PubMedQA, MEDIQA-QA, MEDIQA-AnS and LiveQA-Medical derive answers from a corpus of biomedical literature, whereas emrQA, emrKBQA and why-QA are based on patient notes within EHRs. As stated above, our focus here is on systems that can answer general questions (independent of individual patients) based on the latest evidence, so we do not discuss emrQA, emrKBQA and why-QA. A comparison of the systems and datasets is provided in Table 1.

While BioASQ, MEDIQA-QA, MEDIQA-AnS and LiveQA-Medical are large-scale information retrieval (IR) and question answering (QA) datasets, PubMedQA is designed for "reading comprehension" question answering (RCQA) based on scientific abstracts. Each question of PubMedQA is accompanied by the abstract containing the answer.

The BioASQ Phase B challenge comprises the following question types (Tsatsaronis et al., 2015):

- Exact: "yes" or "no", e.g., "Is the protein Papilin secreted?";

- Factoid: named entities, e.g., "Name synonym of Acrokeratosis paraneoplastica.";

- List: list of named entities, e.g., "Which miR-NAs could be used as potential biomarkers for epithelial ovarian cancer?";

- Ideal: paragraph-sized summaries (text spans), e.g., "What is the effect of TRH on myocardial contractility?"

While BioASQ has been instrumental to the progress of the field (Nentidis et al., 2017, 2018,

| QA System/ Dataset | D1 | D2 | D3 | D4 | D5 | D6 |
|---|---|---|---|---|---|---|
| BioASQ (Krallinger et al., 2020) | ✗ | ✗ | ✗ | ✓ | ✗ | ✗ |
| PubMedQA (Jin et al., 2019) | ✗ | ✗ | ✗ | ✓ | ✗ | ✗ |
| MEDIQA-QA (Ben Abacha et al., 2019) | ✓ | ✓ | ✗ | ✓ | ✗ | ✗ |
| MEDIQA-AnS (Savery et al., 2020) | ✓ | ✓ | ✗ | ✗ | ✗ | ✗ |
| LiveQA-Medical (Abacha et al., 2017) | ✓ | ✓ | ✗ | ✓ | ✗ | ✗ |
| MEANS (Ben Abacha and Zweigenbaum, 2015) | ✗ | ✗ | ✗ | ✓ | ✗ | ✗ |
| AskHERMES (Cao et al., 2011) | ✓ | ✗ | ✗ | ✓ | ✗ | ✗ |
| CLINIQA (Zahid et al., 2018) | ✗ | ✓ | ✗ | ✓ | ✗ | ✗ |
| MedQA (Yu et al., 2007) | ✗ | ✓ | ✗ | ✓ | ✗ | ✗ |

Table 1: Comparision of how well QA systems and datasets meet the desiderata outlined in §2.

2020; Krallinger et al., 2020), it satisfies only one of the criteria we have enumerated above, namely 4. Figure 2 shows the expected output of a system developed using BioASQ. In this example, the extract is provided verbatim (criterion 4).

However, the answer is sourced from a general review; these reviews are less reliable than guidelines or systematic reviews (criterion 1). Furthermore, the system outputs absolute answers rather than guidance (criterion 2) which limits their usefulness to clinicians. A more suitable answer would be "the following guidance is provided in X...". It is unclear what resources are available to the clinician and the BioASQ dataset does not account for this (criterion 3). There is no contradictory evidence in the example and BioASQ has been preprocessed to ensure there are no conflicting papers (criterion 5). Unless the trained model is acting on a curated knowledge base, it would not be robust to conflicts. Finally, the absolute nature of the answer does not allow the system to recognise and account for uncertainty (criterion 6).

In contrast to BioASQ, PubMedQA provides answers to only Boolean (yes/no) questions, e.g. "Do preoperative statins reduce atrial fibrillation after coronary artery bypass grafting?". Accompanying these responses is a "long answer", supplied in the form of the conclusions of the source abstracts. As per Figure 3, the outputs of systems trained on PubMedQA can only satisfy criterion 4. The conclusion is given verbatim to support the short answer. Nevertheless, the source of the answer is not specified (criterion 1), the answer is absolute (criterion 2) and it does not account for any uncertainty (criterion 6). Systems developed using PubMedQA cannot ensure that the answer is useful to the clinician (criterion 3). Given the task is framed as "reading comprehension", there is only one abstract per question. This prevents systems from being trained to handle conflicts (criterion 5).

MEDIQA-QA is a consumer QA dataset whose answers consist of exact snippets from MedlinePlus. Consumer questions are more focused on general information, symptom or person/organization questions (Roberts and Demner-Fushman, 2016). The answers that are required by consumers are less complex and more easily understandable than those given to clinicians (Savery et al., 2020). This has motivated the development of MEDIQA-AnS which summarises the answers of MEDIQA-QA. As shown in Figures 7 and 8, MEDIQA-QA satisfies desiderata 1,2 and 4 while MEDIQA-AnS satisfies only 4.

Although the LiveQA-Medical dataset uses the same answers and sources as MEDIQA-QA and MEDIQA-Ans, it differs by providing answers to each subquestion of the query. Additionally, verbatim extracts of MedlinePlus are used in the responses (Figure 9). Hence criteria 1, 2 and 4 are fulfilled.

MEANS returns only an extract of the original source, without any contextualizing information (Figure 10), i.e. the provenance of the answer. Therefore, only condition 4 is satisfied.

On the other hand, AskHERMES provides a list of answers which are labelled with topics from the question (Figure 11). The extracts shown are from the original sources and are accompanied by links, authors, and dates. Thus, AskHERMES satisfies desiderata 1 and 4.

CLINIQA responds to queries with original abstracts that are accompanied with the PMID and the title of the source paper (Figure 12). However, the results are not rank according to reliability, so

only criteria 2 and 4 are met.

Finally, MedQA's answers comprise sourced extracts from Medline and Google:Definition (Figure 13). Answers are not ranked according to reliability, so the system only satisfies criteria 2 and 4.

None of the aforementioned datasets or systems address conflicts (criterion 5) or communicate uncertainty to clinicians (criterion 6). What might QA systems that satisfy all desiderata look like?

## 4 Presentation of Answers

We have seen that systems trained on BioASQ and PubMedQA do not satisfy all the criteria defined in §2. In this section we present illustrative outputs of hypothetical systems that meet the full set of criteria we have put forth.

Figure 4 presents an example output which satisfies the criteria but where no conflicts occur (criterion 5). The answer is sourced from a systematic review (criterion 1) and is in the form of guidance (criterion 2). While the guidance is actionable given the resources available (criterion 3) and the source extract is reproduced directly (criterion 4), the uncertainty in the answer is acknowledged (criterion 6) by stating the absence of relevant local and national guidelines. The corresponding words and phrases in the question, answer and title used to extract the text snippet are highlighted.

A demonstration of how conflicting evidence could be addressed is provided in Figure 5. In this scenario, the question "Should spinal manipulations be used to treat headaches?" returned three contradictory systematic reviews (Bryans et al., 2011; Chaibi et al., 2011; Posadzki and Ernst, 2011) (criterion 1). Therefore, the system refrains from providing explicit guidance (criterion 6) and instead provides the clinician with the names and links of conflicting reviews (criterion 5). In addition, the clinician is able to investigate the contradictory snippets further by clicking on "Conflicting Snippets" which would show the snippets in Figure 6. Criterion 4 is inapplicable in this case as no answer was retrieved from the documents.

One promising direction which may permit improved handling of contradictory evidence involves use of argumentation-based logic to "reason" about multiple potentially conflicting inputs (Chapman et al., 2019; Cyras et al., 2018), perhaps after explicitly inferring the reported findings concerning treatment efficacies (Lehman et al., 2019; Nye et al., 2020). An alternative (more audacious) direction would be to generate comparative summaries for clinicians that compose narrative summaries of the evidence on a given topic from primary sources, including discussion of conflicting evidence (Wallace et al., 2020; Shah et al., 2021).

Developing and assessing systems according to the criteria outlined in §2 would ensure the output is useful, actionable and reliable to clinicians. It would additionally improve the accountability of both the clinician and the system as the form of the output would be conducive to debugging and root cause analysis.

## 5 Conclusions

We have introduced criteria for assessing the transparency of medical question answering systems. These have been guided by the following question: What would be needed for clinicians to trust, and act upon answers from a QA system? In part we have argued that these systems should be explicitly informed by principles of EBM. The adequacy of existing medical systems and datasets, including BioASQ, PubMedQA, MEDIQA-QA, MEDIQA-AnS, LiveQA-Medical, MEANS, AskHERMES, CLINIQA and MedQA, was assessed using the transparency criteria that we proposed. We found that they met some, but not all, of the conditions.

We presented hypothetical examples of system outputs that satisfy all of the criteria and explained how they could be useful to clinicians. These included conflicts between sources of similar reliability. In these cases, the best course of action was to refrain from giving guidance and instead return the sources to the clinicians for further examination. The examples could form the basis of new datasets and systems that provide actionable answers to clinicians.

We believe that these avenues of investigation would assist with the *deployment* of medical QA systems, ultimately furthering the practice of EBM.

## 6 Acknowledgements

# References

Asma Ben Abacha, Eugene Agichtein, Yuval Pinter, and Dina Demner-Fushman. 2017. Overview of the medical question answering task at trec 2017 liveqa. In *TREC*.

Brian S Alper and R Brian Haynes. 2016a. EBHC pyramid 5.0 for accessing preappraised evidence and guidance. *Evidence Based Medicine*, 21(4):123.

Brian S Alper and R Brian Haynes. 2016b. Ebhc pyramid 5.0 for accessing preappraised evidence and guidance. *BMJ evidence-based medicine*, 21(4):123–125.

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. Publisher: Public Library of Science.

Asma Ben Abacha, Chaitanya Shivade, and Dina Demner-Fushman. 2019. Overview of the MEDIQA 2019 shared task on textual inference, question entailment and question answering. In *Proceedings of the 18th BioNLP Workshop and Shared Task*, pages 370–379, Florence, Italy. Association for Computational Linguistics.

Asma Ben Abacha and Pierre Zweigenbaum. 2015. MEANS: A medical question-answering system combining NLP techniques and semantic Web technologies. *Information Processing & Management*, 51(5):570–594.

Roland Bryans, Martin Descarreaux, Mireille Duranleau, Henri Marcoux, Brock Potter, Rick Ruegg, Lynn Shaw, Robert Watkin, and Eleanor White. 2011. Evidence-Based Guidelines for the Chiropractic Treatment of Adults With Headache. *Journal of Manipulative and Physiological Therapeutics*, 34(5):274–289.

YongGang Cao, Feifan Liu, Pippa Simpson, Lamont Antieau, Andrew Bennett, James Cimino, John Ely, and Hong Yu. 2011. AskHERMES: An online question answering system for complex clinical questions. *Journal of biomedical informatics*, 44:277–88.

Aleksander Chaibi, Peter J Tuchin, and Michael Bjørn Russell. 2011. Manual therapies for migraine: a systematic review. *The journal of headache and pain*, 12(2):127–133. Edition: 2011/02/05 Publisher: Springer Milan.

Martin Chapman, Panagiotis Balatsoukas, Mark Ashworth, Vasa Curcin, Nadin Kökciyan, Kai Essers, Isabel Sassoon, Sanjay Modgil, Simon Parsons, and Elizabeth I. Sklar. 2019. Computational argumentation-based clinical decision support. In *Proceedings of the 18th International Conference on Autonomous Agents and MultiAgent Systems*, AAMAS '19, page 2345–2347, Richland, SC. International Foundation for Autonomous Agents and Multiagent Systems.

K. Cyras, B. Delaney, Denys Prociuk, Francesca Toni, M. Chapman, Jesús Domínguez, and V. Curcin. 2018. Argumentation for explainable reasoning with conflicting medical recommendations. In *MedRACER+WOMoCoE@KR*.

Guilherme Del Fiol, T. Elizabeth Workman, and Paul N. Gorman. 2014. Clinical Questions Raised by Clinicians at the Point of Care: A Systematic Review. *JAMA Internal Medicine*, 174(5):710–718.

Mark Ebell, Jay Siwek, Barry Weiss, Steven Woolf, Jeffrey Susman, Bernard Ewigman, and Marjorie Bowman. 2004. Strength of Recommendation Taxonomy (SORT): A Patient-Centered Approach to Grading Evidence in the Medical Literature. *The Journal of the American Board of Family Practice / American Board of Family Practice*, 17:59–67.

John Ely, Jerome Osheroff, Paul Gorman, Mark Ebell, M Chambliss, Eric Pifer, and P Stavri. 2000. A taxonomy of generic clinical questions: Classification study. *BMJ (Clinical research ed.)*, 321:429–32.

Jungwei Fan. 2019. Annotating and characterizing clinical sentences with explicit why-QA cues. In *Proceedings of the 2nd Clinical Natural Language Processing Workshop*, pages 101–106, Minneapolis, Minnesota, USA. Association for Computational Linguistics.

John Gabbay and Andrée le May. 2016. Mindlines: making sense of evidence in practice. *British Journal of General Practice*, 66(649):402.

Gordon Guyatt, Andrew Oxman, Gunn Vist, Regina Kunz, Yngve Falck-Ytter, Pablo Alonso, and Holger Schünemann. 2008. GRADE: An emerging consensus on rating quality of evidence and strength of recommendations. *BMJ (Clinical research ed.)*, 336:924–6.

Phil Hider, Gemma Griffin, Marg Walker, and Edward Coughlan. 2009. The information-seeking behavior of clinical staff in a large health care organization. *Journal of the Medical Library Association : JMLA*, 97:47–50.

Julian PT Higgins, Douglas G Altman, Peter C Gøtzsche, Peter Jüni, David Moher, Andrew D Oxman, Jelena Savović, Kenneth F Schulz, Laura Weeks, and Jonathan AC Sterne. 2011. The cochrane collaboration's tool for assessing risk of bias in randomised trials. *Bmj*, 343.

Arjen Hoogendam, Anton F H Stalenhoef, Pieter F de Vries Robbé, and A John P M Overbeke. 2008. Answers to questions posed during daily patient care are more likely to be answered by UpToDate than PubMed. *Journal of medical Internet research*, 10(4):e29–e29. Publisher: Gunther Eysenbach.

Qiao Jin, Bhuwan Dhingra, Zhengping Liu, William Cohen, and Xinghua Lu. 2019. *PubMedQA: A Dataset for Biomedical Research Question Answering*. Pages: 2577.

Martin Krallinger, Anastasia Krithara, A. Nentidis, G. Paliouras, and Marta Villegas. 2020. Bioasq at clef2020: Large-scale biomedical semantic indexing and question answering. *Advances in Information Retrieval*, 12036:550 – 556.

Eric Lehman, Jay DeYoung, Regina Barzilay, and Byron C. Wallace. 2019. Inferring Which Medical Treatments Work from Reports of Clinical Trials. In *Proceedings of the Conference of the North American Chapter of the Association for Computational Linguistics (NAACL)*, pages 3705–3717.

Jimmy Lin, Dennis Quan, Vineet Sinha, Karun Bakshi, David Huynh, Boris Katz, and David R. Karger. 2003. What makes a good answer? the role of context in question answering. In *Proceedings of INTERACT 2003*, pages 25–32.

Anastasios Nentidis, Konstantinos Bougiatiotis, Anastasia Krithara, and Georgios Paliouras. 2020. Results of the seventh edition of the BioASQ challenge. In *Machine Learning and Knowledge Discovery in Databases*, pages 553–568. Springer International Publishing.

Anastasios Nentidis, Konstantinos Bougiatiotis, Anastasia Krithara, Georgios Paliouras, and Ioannis Kakadiaris. 2017. Results of the fifth edition of the BioASQ challenge. In *BioNLP 2017*, pages 48–57, Vancouver, Canada,. Association for Computational Linguistics.

Anastasios Nentidis, Anastasia Krithara, Konstantinos Bougiatiotis, Georgios Paliouras, and Ioannis Kakadiaris. 2018. Results of the sixth edition of the BioASQ challenge. In *Proceedings of the 6th BioASQ Workshop A challenge on large-scale biomedical semantic indexing and question answering*, pages 1–10, Brussels, Belgium. Association for Computational Linguistics.

NICE. 2019. UTI (lower): antimicrobial prescribing. https://www.nice.org.uk/guidance/ng109/resources/visual-summary-pdf-6544021069.

Benjamin E Nye, Jay DeYoung, Eric Lehman, Ani Nenkova, Iain J Marshall, and Byron C Wallace. 2020. Understanding clinical trial reports: Extracting medical entities and their relations. *arXiv preprint arXiv:2010.03550*.

Denise Fiona O'leary and Siobhán Ni Mhaolrúnaigh. 2012. Information-seeking behaviour of nurses: where is information sought and what processes are followed? *Journal of Advanced Nursing*, 68(2):379–390. Publisher: John Wiley & Sons, Ltd.

Anusri Pampari, Preethi Raghavan, Jennifer Liang, and Jian Peng. 2018. *emrQA: A Large Corpus for Question Answering on Electronic Medical Records*.

Amy Papermaster and Jane Dimmitt Champion. 2017. The common practice of "curbside consultation": A systematic review. *Journal of the American Association of Nurse Practitioners*, 29(10).

Amy E Papermaster and Jane Dimmitt Champion. 2020. Exploring the use of curbside consultations for interprofessional collaboration and clinical decision-making. *Journal of Interprofessional Care*, pages 1–8. Publisher: Taylor & Francis.

Paul Posadzki and Edzard Ernst. 2011. Spinal manipulations for cervicogenic headaches: a systematic review of randomized clinical trials. *Headache*, 51(7):1132–1139.

Preethi Raghavan, Jennifer J. Liang, Diwakar Mahajan, Rachita Chandra, and Peter Szolovits. 2021. emrkbqa: A clinical knowledge-base question answering dataset. In *BIONLP*.

Kirk Roberts and Dina Demner-Fushman. 2016. Interactive use of online health resources: a comparison of consumer and professional questions. *Journal of the American Medical Informatics Association*, 23(4):802–811.

David L Sackett, R Brian Haynes, Peter Tugwell, et al. 1985. *Clinical epidemiology: a basic science for clinical medicine.* Little, Brown and Company.

Mourad Sarrouti and Said Ouatik El Alaoui. 2020. SemBioNLQA: A semantic biomedical question answering system for retrieving exact and ideal answers to natural language questions. *Artificial Intelligence in Medicine*, 102:101767.

Max Savery, Asma Ben Abacha, Soumya Gayen, and Dina Demner-Fushman. 2020. Question-driven summarization of answers to consumer health questions. *Scientific Data*, 7(1):322.

C. Sellars, T. Hughes, and P. Langhorne. 2002. Speech and language therapy for dysarthria due to non-progressive brain damage. *Cochrane Database of Systematic Reviews*, (3). Publisher: John Wiley & Sons, Ltd.

Darsh J. Shah, Lili Yu, Tao Lei, and R. Barzilay. 2021. Nutribullets hybrid: Multi-document health summarization. *ArXiv*, abs/2104.03465.

Richard Smith. 1996. What clinical information do doctors need? *BMJ*, 313(7064):1062.

George Tsatsaronis, Georgios Balikas, Prodromos Malakasiotis, Ioannis Partalas, Matthias Zschunke, Michael R. Alvers, Dirk Weissenborn, Anastasia Krithara, Sergios Petridis, Dimitris Polychronopoulos, Yannis Almirantis, John Pavlopoulos, Nicolas Baskiotis, Patrick Gallinari, Thierry Artiéres, Axel-Cyrille Ngonga Ngomo, Norman Heino, Eric Gaussier, Liliana Barrio-Alvers, Michael Schroeder, Ion Androutsopoulos, and Georgios Paliouras. 2015. An overview of the BIOASQ large-scale biomedical semantic indexing and question answering competition. *BMC Bioinformatics*, 16(1):138.

Byron C. Wallace, Sayantan Saha, Frank Soboczenski, and I. Marshall. 2020. Generating (factual?) narrative summaries of rcts: Experiments

with neural multi-document summarization. *ArXiv*, abs/2008.11293.

Hong Yu, Minsuk Lee, David Kaufman, John Ely, Jerome A. Osheroff, George Hripcsak, and James Cimino. 2007. Development, implementation, and a cognitive evaluation of a definitional question answering system for physicians. *J. of Biomedical Informatics*, 40(3):236–251.

M Zahid, Ankush Mittal, R. Joshi, and G. Atluri. 2018. *CLINIQA: A Machine Intelligence Based Clinical Question Answering System.*
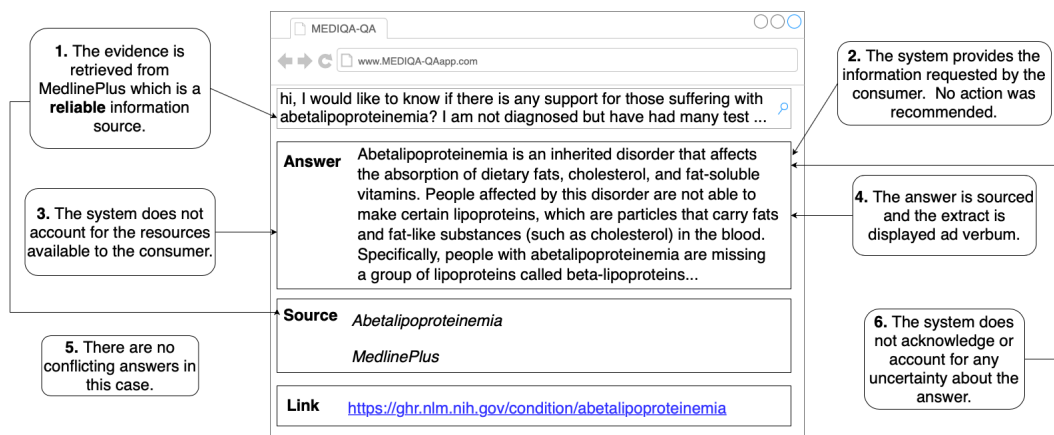
## A   Additional figures of QA interfaces

**1.** The evidence is retrieved from MedlinePlus which is a **reliable** information source.

**3.** The system does not account for the resources available to the consumer.

**5.** There are no conflicting answers in this case.

MEDIQA-QA

www.MEDIQA-QAapp.com

hi, I would like to know if there is any support for those suffering with abetalipoproteinemia? I am not diagnosed but have had many test ...

**Answer** Abetalipoproteinemia is an inherited disorder that affects the absorption of dietary fats, cholesterol, and fat-soluble vitamins. People affected by this disorder are not able to make certain lipoproteins, which are particles that carry fats and fat-like substances (such as cholesterol) in the blood. Specifically, people with abetalipoproteinemia are missing a group of lipoproteins called beta-lipoproteins...

**Source** *Abetalipoproteinemia*

*MedlinePlus*

**Link** https://ghr.nlm.nih.gov/condition/abetalipoproteinemia

**2.** The system provides the information requested by the consumer. No action was recommended.

**4.** The answer is sourced and the extract is displayed ad verbum.

**6.** The system does not acknowledge or account for any uncertainty about the answer.

Figure 7: Web interface for QA system developed using MEDIQA-QA.



**1.** The evidence is retrieved from MedlinePlus which is a **reliable** information source.

**3.** The system does not account for the resources available to the consumer.

**5.** There are no conflicting answers in this case.

MEDIQA-AnS

www.MEDIQA-AnSapp.com

hi, I would like to know if there is any support for those suffering with abetalipoproteinemia? I am not diagnosed but have had many test ...

**Answer** Test that can help diagnose abetalipoproteinemia, also known as Bassen-Kornzweig syndrome, include blood tests to look for vitamin deficiencies, red blood cell malformations, complete blood count (CBC) and cholesterol. Other tests may be needed: electromyography, eye exam, nerve conduction velocity, stool sample analysis and genetic testing. Treatment for abetalipoproteinemia includes large doses of fat-soluble vitamins (A, D, E, and K) and linoleic acid supplements...

**Source** *Abetalipoproteinemia*

*MedlinePlus*

**Link** https://ghr.nlm.nih.gov/condition/abetalipoproteinemia

**2.** The system provides the information requested by the consumer. No action was recommended.

**4.** The answer is sourced and an abstractive summary is shown.

**6.** The system does not acknowledge or account for any uncertainty about the answer.

Figure 8: Web interface for QA system developed using MEDIQA-AnS.



**1.** The evidence is retrieved from National institutes for Health websites which are **reliable** information sources.

**3.** The system does not account for the resources available to the consumer.

**5.** There are no conflicting answers in this case.

MEDIQA-AnS

www.MEDIQA-AnSapp.com

hi, I would like to know if there is any support for those suffering with abetalipoproteinemia? I am not diagnosed but have had many test ...

**Answer to 1st subquestion** There may be damage to the retina of the eye (retinitis pigmentosa). Tests that may be done to help diagnose thiscondition include: - Apolipoprotein B blood test - Blood tests to look for vitamin deficiencies (fat-soluble vitamins A, D, E, and K)...

**Answer to 2nd subquestion** Treatment involves large doses of vitamin supplements containing fat-soluble vitamins (vitamin A, vitamin D, vitamin E, and vitamin K). Linoleic acid supplements are also recommended. People with this condition should talk to a nutritionist. Diet changes are needed to prevent stomach problems...

**2.** The system provides the information requested by the consumer. No action was recommended.

**4.** The answers are sourced and the extract is displayed ad verbum.

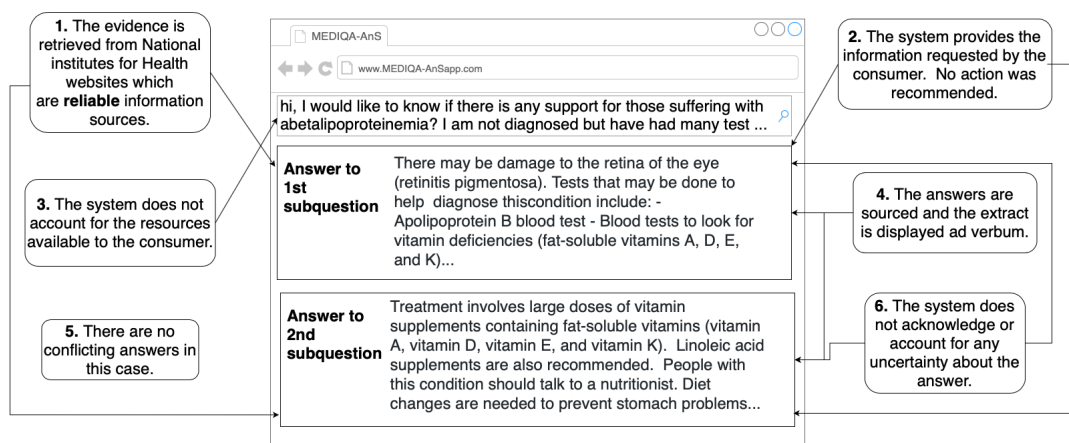**6.** The system does not acknowledge or account for any uncertainty about the answer.

Figure 9: Web interface for QA system developed using LiveQA-Medical.
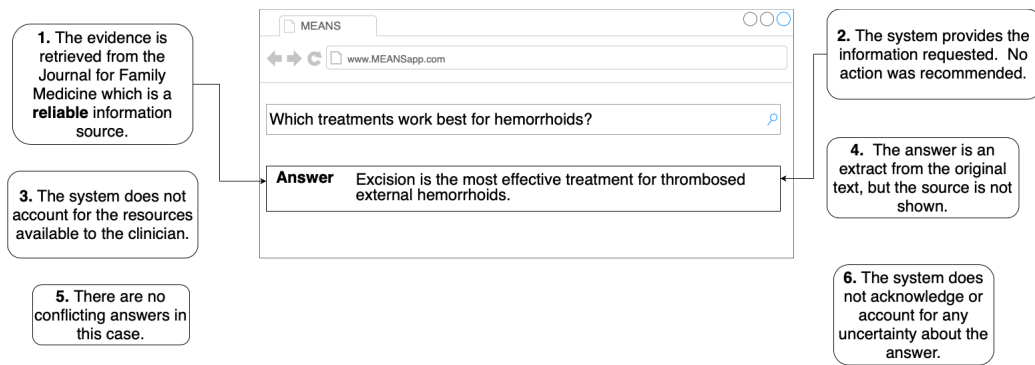
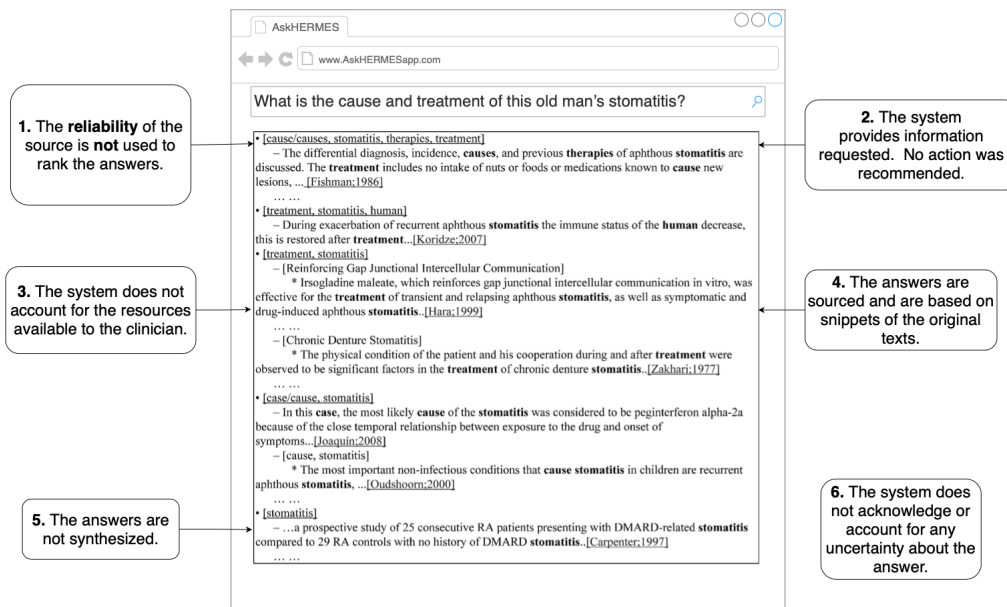Figure 10: Web interface for MEANS.
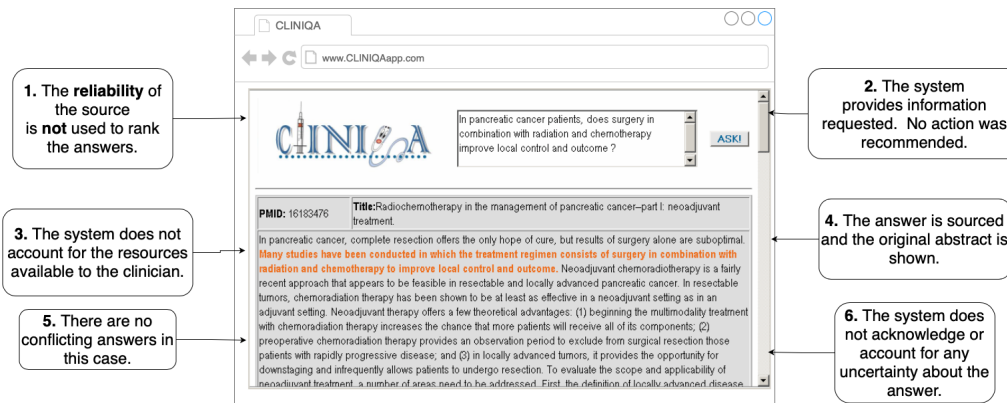


Figure 11: Web interface for AskHERMES.



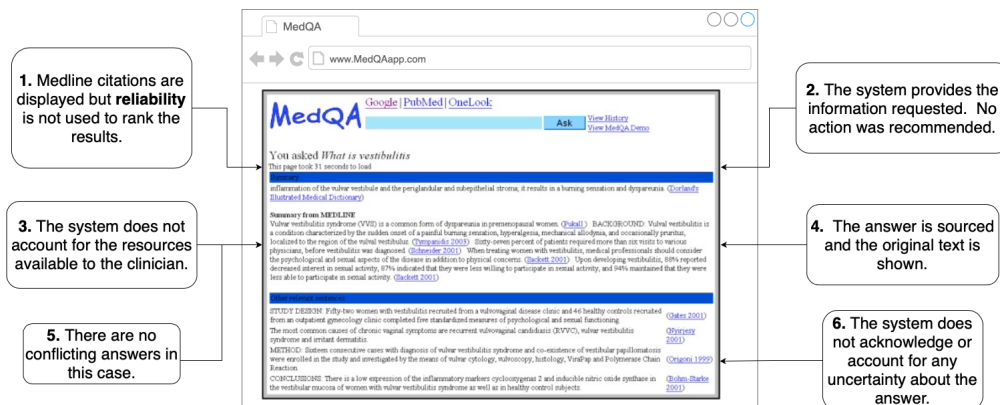Figure 12: Web interface for CLINIQA which includes figure 5 from (Zahid et al., 2018).

Figure 13: Web interface for MedQA which includes figure 3 from (Yu et al., 2007).