# MiBi at BioASQ 2024:
# Retrieval-Augmented Generation for Answering Biomedical Questions

September 9, 2024

Jan Heinrich
Merker

Alexander
Bondarenko

Matthias
Hagen

Adrian
Viehweger

Friedrich-Schiller-Universität Jena
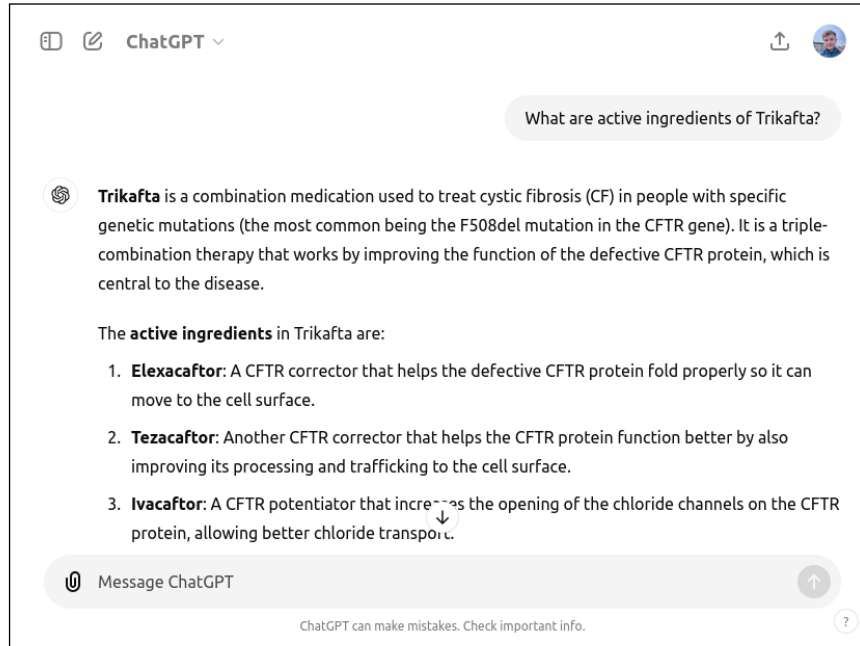Leipzig University

https://webis.de

CLEF 2024
GRENOBLE

15th International Conference of the CLEF Association (CLEF 2024)

## Medical Q&A

Example: What are active ingredients of Trikafta?

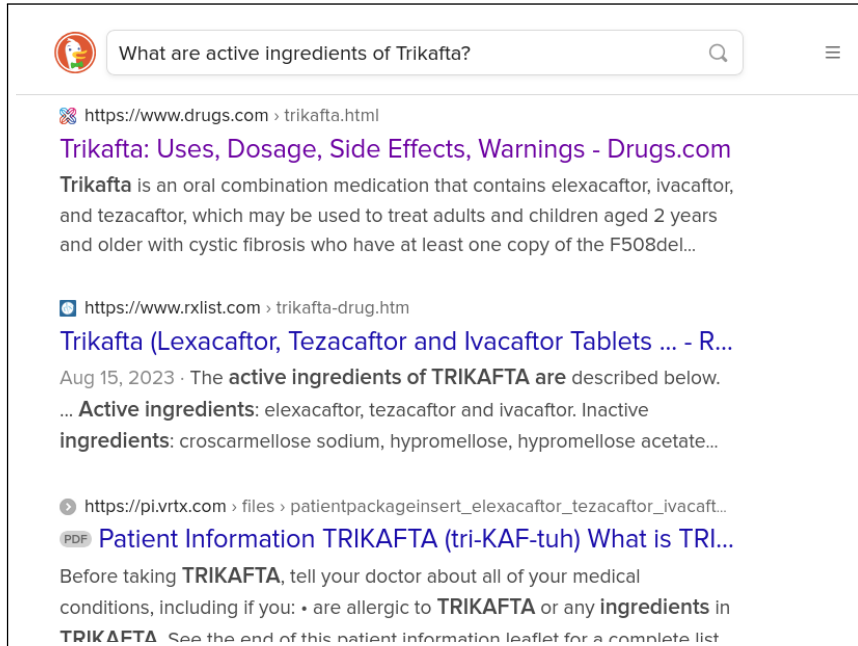# MiBi at BioASQ 2024
## Baselines

### Example: What are active ingredients of Trikafta?



💬 Why not just use GPT . . . ?
(correct ingredients, no dosage, no sources)
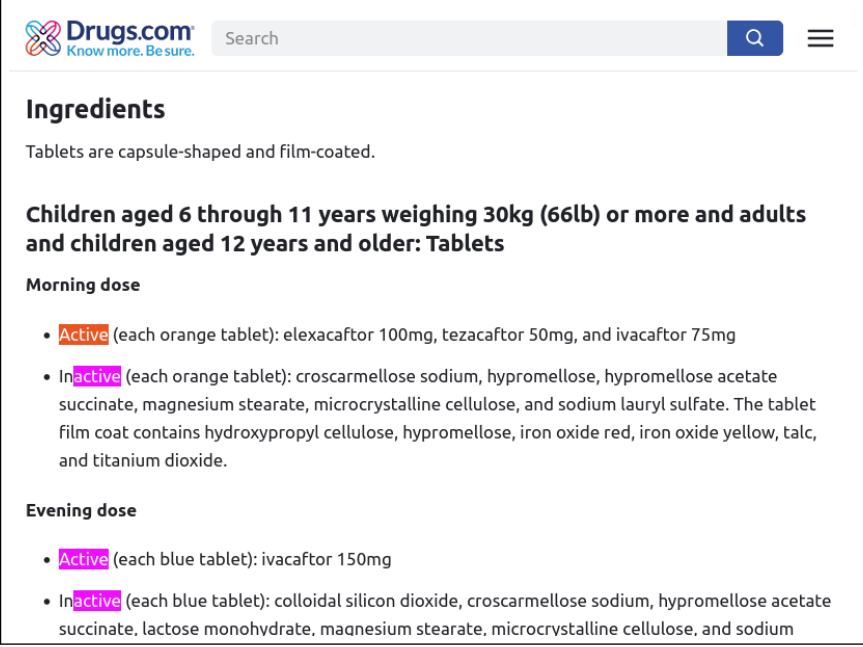
Example: What are active ingredients of Trikafta?



🔍 ... or a quick web search ... ?

# MiBi at BioASQ 2024

## Baselines

Example: What are active ingredients of Trikafta?



🌐    … with Ctrl+F on the first result?

(correct ingredients and dosage, good source, but takes longer)

# MiBi at BioASQ 2024
## Baselines

Example: What are active ingredients of Trikafta?



😳 . . . and what about this?

(correct ingredients and dosage, good source, fastest?!)

# MiBi at BioASQ 2024

## Goal: RAG for Medical Questions

| Document retrieval | → | Snippet extraction and re-ranking | → | Exact answer generation | → | Ideal answer generation |
|---|---|---|---|---|---|---|

## Stages

- Document retrieval ➔ Find relevant medical articles (from PubMed).

- Snippet extraction and re-ranking ➔ Extract snippets and rank by relevance.

- Answer generation ➔ Generate exact answer and "ideal" summary answer.

- RAG ➔ Combine retrieval- and generation-focused components.

# MiBi at BioASQ 2024

## Approaches: Document Retrieval

```
┌──────────────────┐   ┌──────────────────┐   ┌──────────────────┐   ┌──────────────────┐
│ Document retrieval│──▶│ Snippet extraction│──▶│  Exact answer    │──▶│  Ideal answer    │
│                  │   │  and re-ranking  │   │   generation     │   │   generation     │
└──────────────────┘   └──────────────────┘   └──────────────────┘   └──────────────────┘
```

Goal: Find relevant medical articles (from PubMed).

❑ PubMed search API

   – Re-rank with BM25, MiniLM, and MPNet

❑ Custom BM25 index with metadata (Elasticsearch)

   – Match abstract, title, and MeSH terms

   – Disallow non-peer-reviewed publication types

➔ Do we need to wory about indexing?

# MiBi at BioASQ 2024

## Approaches: Snippet Extraction and Re-Ranking

| Document retrieval | → | Snippet extraction and re-ranking | → | Exact answer generation | → | Ideal answer generation |

Goal: Extract concise snippets from the article's abstract (or title).
Rank extracted snippets by relevance to the question.

❏ Using LLMs

– Chain-of-thought 3-shot prompt (GPT-3.5-turbo)
– No re-ranking

❏ Rule-based

– Split abstract in sentences
– Candidates: full title + sentence $n$-grams (up to 3 sent.) from abstract
– Re-rank with TAS-B and duoT5

➔ Are LLMs better at snippet extraction?

# MiBi at BioASQ 2024

## Approaches: Answer Generation with LLMs

| Document retrieval | → | Snippet extraction and re-ranking | → | Exact answer generation | → | Ideal answer generation |

Goal: Generate exact (e.g., yes-no) answer and "ideal" summary answer.

❑ Few-shot prompting with function calling

  – Manual prompts per question/answer type (e.g., yes-no / exact)

  – Context: top-3 abstracts or all (top-10) snippets

  – GPT-3.5-turbo and GPT-4

❑ Modular "programming" with DSPy

  – Automatic prompts via DSPy (signature of in-/outputs are Python classes)

  – Context: abstracts, snippets, previous answer

  – Mixtral-7B

→ Do we need manual prompts?

# MiBi at BioASQ 2024
## Approaches: RAG paradigms

| Document retrieval | ? | Snippet extraction and re-ranking | ? | Exact answer generation | ? | Ideal answer generation |
|---|---|---|---|---|---|---|

Goal: Combine retrieval- and generation-focused components.

❑ Retrieve-then-generate (exact → ideal → documents → snippets)

❑ Generate-then-retrieve (documents → snippets → exact → ideal)

❑ GtRtG / RtGtR (e.g., exact → ideal → documents → snippets → exact → ideal)

❑ Let the LLM decide (DSPy, Mixtral-8x7B)

→ Which paradigm to use when?

# MiBi at BioASQ 2024
## Results

- 42 submitted runs, different systems per phase and batch

- Retrieval:
  - PubMed search API struggles with question-like queries
  - Enhancing index with metadata pays off for domain-specific retrieval

- Snippet extraction and re-ranking:
  - Neither GPT- nor rule-based snippet extraction competitive

- Answer generation:
  - Snippets instead of abstracts as context less "confusing" for LLMs
  - Models: GPT-4 $>>$ GPT-3.5 $>$ Mixtral-7B
  - No difference by prompting strategy (manual vs. DSPy)

- RAG paradigm:
  - With ground truth: RtG, GtRtG
  - Without ground truth: GtR, GtRtG

# MiBi at BioASQ 2024

## Summary

❏ Mixed results, but some recommendations:
  – Put work into first-stage retrieval
  – Show snippets to LLMs, not long texts
  – Use the latest LLMs (e.g., GPT-4)
  – GtRtG seems to work well with/without ground-truth evidence

❏ Limitation: comparability across test batches

❏ Future work: systematic evaluation of grounded RAG paradigms

## Code and Data

○ github.com/webis-de/CLEF-24

# MiBi at BioASQ 2024
## Summary

- ❏ Mixed results, but some recommendations:
  - – Put work into first-stage retrieval
  - – Show snippets to LLMs, not long texts
  - – Use the latest LLMs (e.g., GPT-4)
  - – GtRtG seems to work well with/without ground-truth evidence

- ❏ Limitation: comparability across test batches

- ❏ Future work: systematic evaluation of grounded RAG paradigms

## Code and Data

 github.com/webis-de/CLEF-24

*Thank you & merci!*