# Clinical Text Summarization: Adapting Large Language Models Can Outperform Human Experts

Dave Van Veen, Cara Van Uden, Louis Blankemeier, Jean-Benoit Delbrouck, Asad Aali,
Christian Bluethgen, Anuj Pareek, Malgorzata Polacin, Eduardo Pontes Reis,
Anna Seehofnerová, Nidhi Rohatgi, Poonam Hosamani, William Collins, Neera Ahuja,
Curtis P. Langlotz, Jason Hom, Sergios Gatidis, John Pauly, Akshay S. Chaudhari
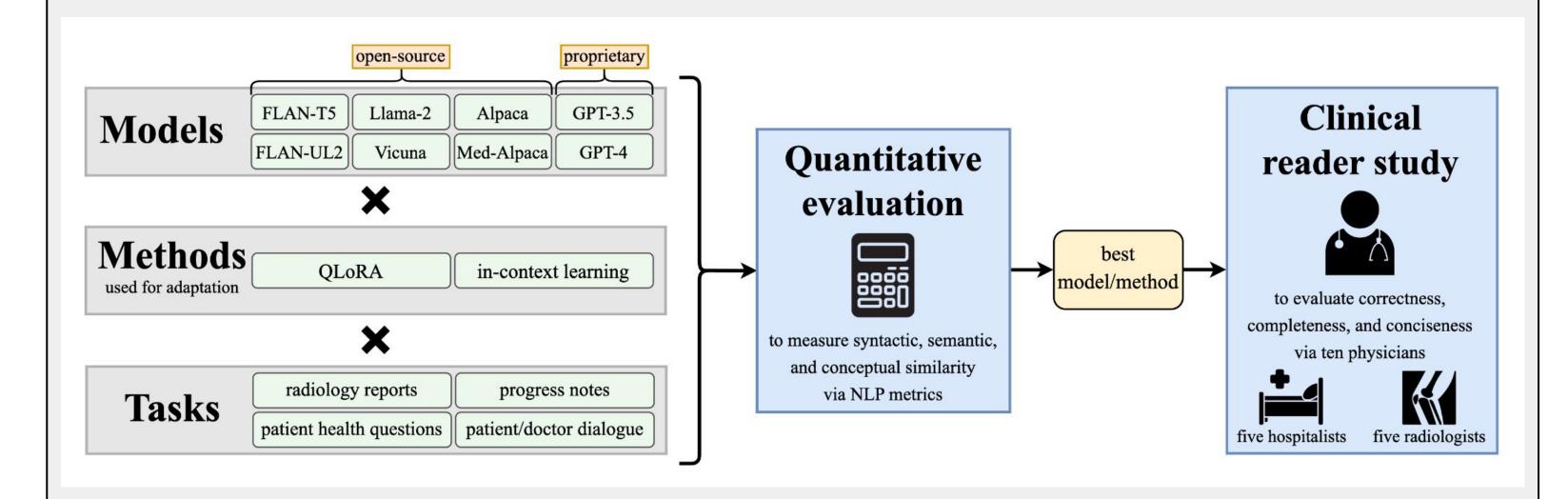
## Abstract

Motivation
- Summarizing key information from electronic health records (EHR) imposes a substantial burden on how clinicians allocate their time.
- Large language models (LLMs) do well on general natural language processing (NLP) tasks, but their efficacy on summarizing clinical text has not been demonstrated.

Outcome
- Our research marks the first evidence of LLMs outperforming human experts for clinical text summarization.
- This implies that integrating LLMs into clinical workflows could alleviate documentation burden, enabling clinicians to focus more directly on patient care.

## Overview

First, we quantitatively evaluate each valid combination (×) of LLM and adaptation method across four distinct summarization tasks comprising six datasets. We then conduct a clinical reader study in which ten physicians compare summaries of the best model/method against those of a human expert.
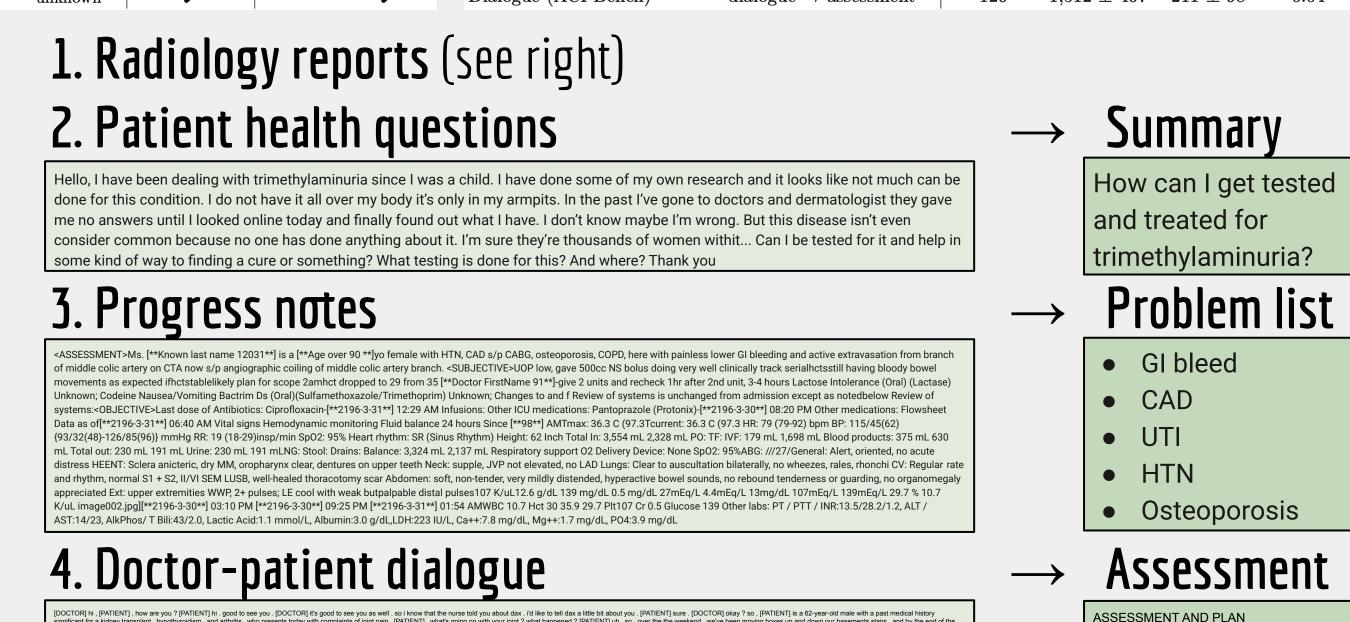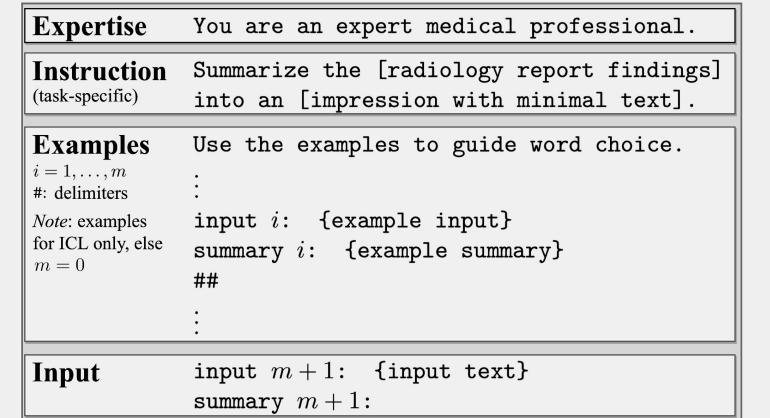


## Approach

We compare two methods—quantized low-rank adaptation (QLoRA) and in-context learning (ICL)—for adapting eight models (**left**) to four summarization tasks (**right**).

| Model | Context | Parameters | Proprietary? | Seq2seq | Autoreg. |
|---|---|---|---|---|---|
| FLAN-T5 | 512 | 2.7B | - | ✔ | - |
| FLAN-UL2 | 2,048 | 20B | - | ✔ | - |
| Alpaca | 2,048 | 7B | - | - | ✔ |
| Med-Alpaca | 2,048 | 7B | - | - | ✔ |
| Vicuna | 2,048 | 7B | - | - | ✔ |
| Llama-2 | 4,096 | 7B, 13B | - | - | ✔ |
| GPT-3.5 | 16,384 | 175B | ✔ | - | ✔ |
| GPT-4 | 32,768 | unknown | ✔ | - | ✔ |

| Task (Dataset) | Task description | Number of samples | Avg. number of tokens Input | Avg. number of tokens Target | Lexical variance |
|---|---|---|---|---|---|
| Radiol. reports (Open-i) | findings → impression | 3.4K | 52 ± 22 | 14 ± 12 | 0.11 |
| Radiol. reports (MIMIC-CXR) | findings → impression | 128K | 75 ± 31 | 22 ± 17 | 0.08 |
| Radiol. reports (MIMIC-III) | findings → impression | 67K | 83 ± 42 | 61 ± 45 | 0.09 |
| Patient questions (MeQSum) | verbose → short question | 1.2K | 83 ± 67 | 14 ± 6 | 0.21 |
| Progress notes (ProbSum) | notes → problem list | 755 | 1,013 ± 299 | 23 ± 16 | 0.15 |
| Dialogue (ACI-Bench) | dialogue → assessment | 126 | 1,512 ± 467 | 211 ± 98 | 0.04 |

### Task examples

**1. Radiology reports** (see right)

**2. Patient health questions**
> Hello, I have been dealing with trimethylaminuria since I was a child. I have done some of my own research and it looks like not much can be done for this condition. I do not have it all over my body it's only in my armpits. In the past I've gone to doctors and dermatologist they gave me no answers until I looked online today and finally found out what I have. I don't know maybe I'm wrong. But this disease isn't even consider common because no one has done anything about it. I'm sure they're thousands of women within... Can I be tested for it and help in some kind of way to finding a cure or something? What testing is done for this? And where? Thank you

→ **Summary**
> How can I get tested and treated for trimethylaminuria?

**3. Progress notes**

→ **Problem list**
- GI bleed
- CAD
- UTI
- HTN
- Osteoporosis

**4. Doctor-patient dialogue**

→ **Assessment**

### Prompt Anatomy

| Expertise | You are an expert medical professional. |
|---|---|
| Instruction (task-specific) | Summarize the [radiology report findings] into an [impression with minimal text]. |
| Examples | Use the examples to guide word choice. |
| | $i = 1, \ldots, m$ |
| | #: delimiters |
| | *Note:* examples for ICL only, else $m = 0$ |
| | input $i$: {example input} |
| | summary $i$: {example summary} |
| | ## |
| Input | input $m+1$: {input text} |
| | summary $m+1$: |

## Quantitative evaluation

**Metrics** ____ measure ____ similarity between generated and reference texts.
**BLEU, ROUGE-L:** syntactic   **BERTScore:** semantic   **MEDCON:** medical conceptual
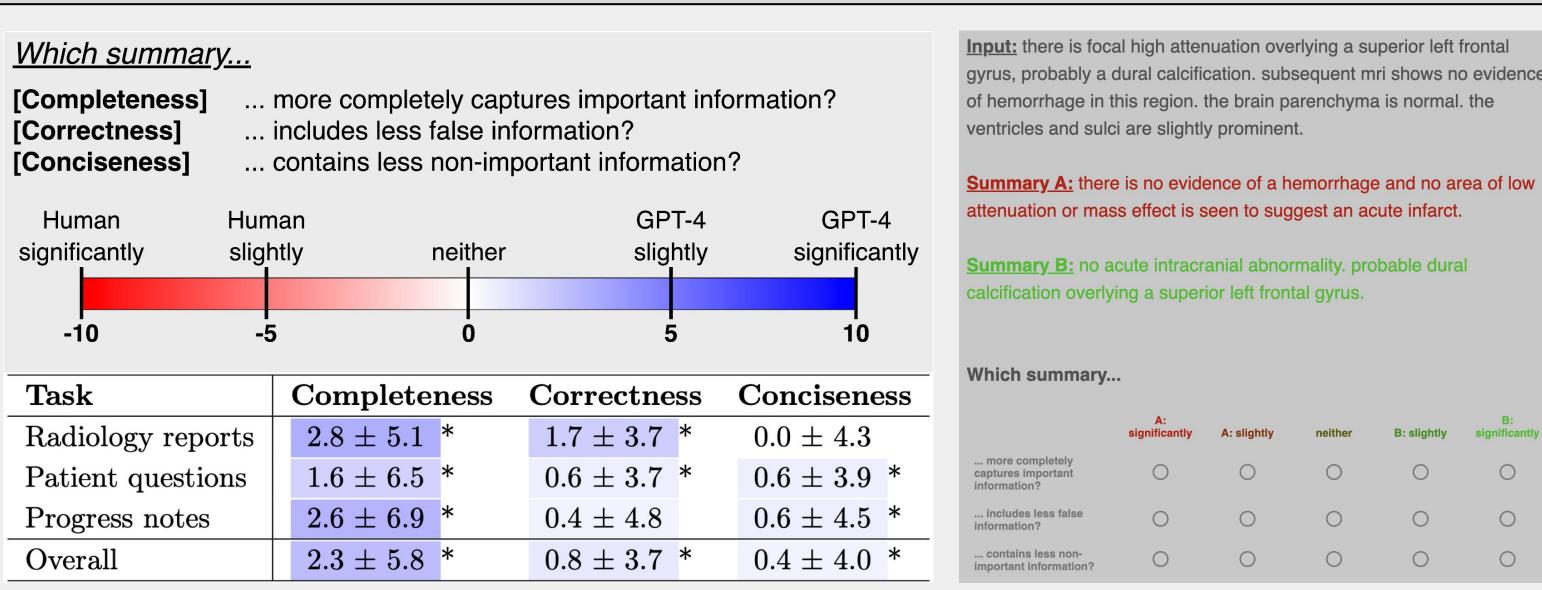


**Top left:** FLAN-T5 emerges as the best open-source model fine-tuned with QLoRA.
**Bottom left:** Given enough examples, ICL surpasses QLoRA (dashed line). Including one example drastically improves performance compared to zero-shot prompting.
**Right:** Head-to-head win rates of each model combination. GPT-4 generally performs best. Seq2seq models outperform open-source autoregressive models.

→ *we conclude the best configuration is GPT-4 with maximal\* in-context examples*
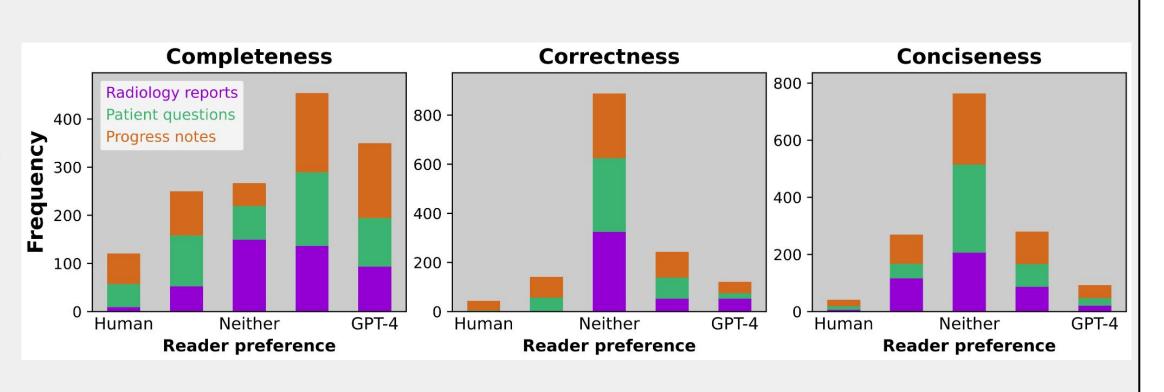
*\*determined by example size and model context length*

## Clinical Reader Study



Which summary...
- **[Completeness]** ... more completely captures important information?
- **[Correctness]** ... includes less false information?
- **[Conciseness]** ... contains less non-important information?

| Task | Completeness | Correctness | Conciseness |
|---|---|---|---|
| Radiology reports | 2.8 ± 5.1 * | 1.7 ± 3.7 * | 0.0 ± 4.3 |
| Patient questions | 1.6 ± 6.5 * | 0.6 ± 3.7 * | 0.6 ± 3.9 * |
| Progress notes | 2.6 ± 6.9 * | 0.4 ± 4.8 | 0.6 ± 4.5 * |
| Overall | 2.3 ± 5.8 * | 0.8 ± 3.7 * | 0.4 ± 4.0 * |

**Left:** Summaries generated via GPT-4 with in-context examples are rated higher with statistical significance (*) on all three attributes. **Right:** Example of study user interface.

Distribution of reader scores across five-point Likert scale.

Annotations of radiology report examples to illustrate the strengths and weaknesses of GPT-4 and humans.

**Contact:** Dave Van Veen | vanveen@stanford.edu | davevanveen.com