

# How Much Would a Clinician Edit This Draft? Evaluating LLM Alignment for Patient Message Response Drafting

Parker Seegmiller<sup>1</sup>, Joseph Gatto<sup>1</sup>, Sarah E. Greer<sup>1</sup>,  
Ganza Belise Isingizwe<sup>1</sup>, Rohan Ray<sup>1</sup>, Timothy Burdick<sup>2,3</sup>, Sarah M. Preum<sup>1</sup>

<sup>1</sup> Department of Computer Science, Dartmouth College

<sup>2</sup> Department of Community and Family Medicine, Dartmouth Health

<sup>3</sup> The Dartmouth Institute, Dartmouth College

{pkseeg.gr, sarah.masud.preum}@dartmouth.edu

## Abstract

Large language models (LLMs) show promise in drafting responses to patient portal messages, yet their integration into clinical workflows raises various concerns, including whether they would actually save clinicians time and effort in their portal workload. We investigate LLM alignment with individual clinicians through a comprehensive evaluation of the patient message response drafting task. We develop a novel taxonomy of thematic elements in clinician responses and propose a novel evaluation framework for assessing clinician editing load of LLM-drafted responses at both content and theme levels. We release an expert-annotated dataset and conduct large-scale evaluations of local and commercial LLMs using various adaptation techniques including thematic prompting, retrieval-augmented generation, supervised fine-tuning, and direct preference optimization. Our results reveal substantial epistemic uncertainty in aligning LLM drafts with clinician responses. While LLMs demonstrate capability in drafting certain thematic elements, they struggle with clinician-aligned generation in other themes, particularly question asking to elicit further information from patients. Theme-driven adaptation strategies yield improvements across most themes. Our findings underscore the necessity of adapting LLMs to individual clinician preferences to enable reliable and responsible use in patient-clinician communication workflows.

## 1 Introduction

The use of large language models (LLMs) for **drafting responses to asynchronous patient messages** has garnered significant interest in the medical community (Hu et al., 2025). This would involve the integration of LLMs in the patient-clinician communication loop by drafting an initial clinician response to an incoming patient message, which the clinician would then edit and send to the patient.

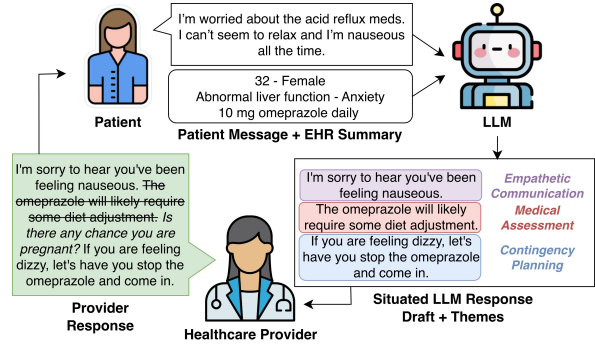


Figure 1: Patient message response drafting. LLMs draft responses to patient messages, then clinicians edit the draft by ~~deleting~~ and *adding* content as needed. We evaluate content-level and theme-level alignment between clinicians and LLMs.

Figure 1 shows an example of this task: generating a response draft to a patient-initiated message, given the message and a summary of the patient’s relevant electronic health record (EHR) data. Responding to patient portal messages places a heavy burden on clinicians due to increasing use of the patient portal and significant clinical workforce constraints (Budd, 2023; Underdahl et al., 2024; Martinez et al., 2023; Yan et al., 2021). As such, there is growing interest in developing AI-mediated support for improving efficiency and engagement in patient portal messaging (Gatto et al., 2025; Biro et al., 2025). Thus, patient portal messaging is a high-stakes, real-world setting for evaluating LLMs on the task of drafting responses.

Prior work has gathered clinician feedback on LLM response drafts to patient portal messages with mixed results. Some studies report that these responses can be useful (Hu et al., 2025; Garcia et al., 2024; Bootsma-Robroeks et al., 2025; English et al., 2024b). However, there is evidence that LLM responses often diverge from clinician responses in style and content, and lack accuracy (Hu et al., 2025; Biro et al., 2025; Sharma et al.).

Divergence between LLM response drafts and

Theme	Example Frame	Example Response Element
Empathy	Encouragement of patient treatment effort	You’ve been doing a great job with your tapering.
Symptom Questions	Asking about location of symptoms	Has your pain only been in your lower back?
Medication Questions	Asking about intake of medications	Have you been taking your Amoxicillin regularly?
Medical Assessment	Explanation of test result	Your iron levels look normal.
Medical Planning	Confirmation of required testing	Let’s get you in for a bloodwork test.
Logistics	Confirmation of clinic policy	We can only offer telehealth in the state.
Care Coordination	Promise of future patient contact	We’ll reach out after we receive the results.
Contingency Planning	Symptom-related backup plan	If you’re feeling dizzy, please call triage

Table 1: Themes derived from clinician responses to patient portal messages, alongside representative frames and example response elements/utterances. For example, “explanation of test result” is a frame within the medical assessment theme, and “your iron levels look normal” is a clinician response component that falls under this frame. In total, we derive 8 clinician response themes comprised of 67 unique frames (examples in supplemental materials).

clinician responses may lead to either **unreliability**, if LLM response drafts must be significantly edited, contributing to clinicians’ workload in responding to patient messages, or **irresponsibility**, if unedited low-quality LLM response draft elements are sent to the patient. Reliability is important, as clinicians spending significant time editing/improving the drafted response defeats the purpose of using LLMs to improve efficiency (Tai-Seale et al., 2024; Bootsma-Robroeks et al., 2025). Clinician responsibility is critical, as LLM-generated drafts may contain clinically-significant errors and adversely impact the standards of care (Biro et al., 2025; Sharma et al.; Chen et al., 2025).

We investigate the use of LLM drafts in supporting clinician responses to patient messages, by evaluating alignment of LLMs to responses generated by real clinicians. Specifically, we aim to explore the content-level and theme-level alignment between clinician-written and LLM-generated responses, to inform responsible use of NLP in patient message response drafting. We answer three relevant research questions. **RQ1:** What constitutes a high-quality clinician response to a patient message? **RQ2:** How might we automate evaluation of LLM response draft quality, with respect to clinician editing workload? **RQ3:** How can we adapt LLMs to support clinicians in generating quality responses to patient messages?

In answering these research questions, we make **four key contributions**. **First**, we use a clinicians-in-the-loop, hybrid approach to develop a clinically-relevant set of “themes” and frames to systematically characterize clinician responses to patient messages. **Second**, we develop and validate a novel two-level evaluation framework for assessing clinician editing load given LLM-drafted responses to patient messages. **Third**, we annotate and release an expert-clinician-annotated dataset for evaluating performance on the patient message

response drafting task<sup>1</sup>. **Finally**, we conduct a rigorous evaluation of three local and three commercial LLMs on this task, using five LLM adaptation techniques varying in degree of supervision, finding that theme-driven adaptation of LLMs improves response drafting performance by 33% over 0-shot models.

## 2 Overview of Data

The patient-clinician conversations used in our experiments are collected from a large academic hospital in the United States. These conversations are sourced from the hospital’s electronic health record (EHR) portal messaging platform. Patient portal messaging is an asynchronous healthcare communication service in which patients and their clinicians discuss a wide variety of patient health issues, including symptoms, medication efficacy, treatment planning, scheduling logistics, and more (North et al., 2019).

We begin with 610k total messages taken from the secure patient portal between 1/2020 - 9/2024. Our dataset includes messages from primary care, and thus includes a wide range of medical topics. We gather all patient-initiated messages which received a written clinician response to create 146k conversations, i.e. original patient message and response from a clinician. Our final data pool contains 10,105 unique patients, of which 64% are female and 36% are male, with ages ranging between 18-80. Each sample in our data pool consists of a patient message, a clinician response, and a summary of the patient’s chart or electronic health record (EHR) data<sup>2</sup>. We utilize 144k conversations from the data pool as training data, and gather evaluation datasets from the remaining 2k conversations.

<sup>1</sup><https://hf.co/collections/PortalPal-AI/evaluating-alignment-for-patient-message-response-drafting>

<sup>2</sup>See appendix A for full dataset details

Dataset	Source	Response	Clinician ct.	Size	Message	Response
IPPM	Patient Portal Message + EHR	Theme-Guided	4	300	83±54	53±32
SyPPM	Synthetic Message + EHR	Theme-Guided	3	100	110±51	70±26
SoCPPM	Patient Portal Message + EHR	Real-Time	196	300	69±45	55±78

Table 2: Summary of the three datasets. Patient messages in IPPM and SoCPPM, and EHR summaries for all datasets are sourced from a real EHR portal. SyPPM messages are semi-synthetic, generated using de-identified real patient messages for public release. Details on how clinician responses are collected and annotated are in Appendix G. We include mean  $\pm$  standard deviation of the word count of patient messages and clinician responses.

## 2.1 Thematic Analysis of Responses

We address RQ1 by carefully deriving elements of high-quality clinician responses to patient messages. Based on manual thematic analysis of our real patient-clinician conversations, and research workshops with a team of 13 expert primary care physicians, nurses, and triage nurses, we derive a set of clinically-relevant “themes” which can be used to characterize the quality of clinician responses to patient messages (Braun and Clarke, 2006; Sun et al., 2013). These themes can be found in Table 1. Appendix B gives full details of our mixed-methods approach to identify these themes.

## 2.2 Summary of Evaluation Datasets

Table 2 summarizes our three evaluation datasets. Here we briefly describe the three datasets derived from these 2k conversations and share additional dataset details in Appendix A.2. Each sample in each dataset is a tuple of strings  $\{m, c, r\}$  consisting of a patient message  $m$ , a summary  $c$  of the patient’s EHR chart and a single clinician response  $r$ . The Ideal Patient Portal Messaging (IPPM) dataset is created to evaluate LLMs in a setting where clinicians do not face the same resource constraints as in the real-world, thus responses are written by a team of paid expert clinicians who are guided by the themes derived in Section 2.1. The publicly-available Synthetic Patient Portal Message (SyPPM) contain semi-synthetic patient portal messages, paired with real de-identified patient EHR summaries, with responses collected via the same method as IPPM. The Standards of Care Patient Portal Messaging (SoCPPM) dataset is created to evaluate LLMs in a practical setting, where response drafts are compared with a clinician response which was sent via the portal in real time, thus responses are collected via the patient portal.

## 3 Scalable Evaluation of LLMs

We want to evaluate the reliability of LLM responses on the response drafting task (RQ2). Our evaluation seeks to identify: in order to achieve

the same quality of response, 1) how much content would the clinician need to *add* to the LLM draft? and 2) how much content would the clinician need to *remove* from the LLM draft? Hence, we use a reference-based approach which directly compares an LLM draft with a response written by an expert clinician (Li et al., 2024). Comparing what needs to be removed from and added to an LLM-drafted response to achieve an expert-written response, is analogous to measuring 1) *recall*, i.e. how much of the expert-written response is covered by the LLM-drafted response, and 2) *precision*, i.e. how much of the LLM-drafted response is matched in the expert’s response. As our goal is to identify the editing load of a clinician using a LLM-as-judge framework, we call this the **EditJudge Evaluation Framework** (Figure 2). This framework is a human-AI collaborative, task-specific, reference-based, LLM-as-judge evaluation framework (Li et al., 2025; Bavaresco et al., 2025).

We use two measures of editing load to capture complementary aspects of alignment between generated and reference responses. The *content-level edit-F1* score assesses whether a response drafting LLM reproduces specific clinical facts, instructions, or action items present in the reference, which is critical for safety and correctness. However, clinically appropriate drafts may vary substantially in wording or level of detail while addressing the same underlying intent. The *theme-level edit-F1* score captures higher-level alignment by measuring whether the response addresses thematically similar clinical goals, concerns, and communicative functions (e.g., reassurance, triage guidance, or follow-up planning), even when the granular content differs. Using both metrics distinguishes incomplete response drafts from those that are semantically (content-level) and thematically aligned but phrased differently, providing a more reliable evaluation of response draft quality.

### 3.1 Content-Level edit-F1 Score

Given an expert-written clinician response  $r_e$  and an LLM response draft  $r_d$ , the content-level edit-F1

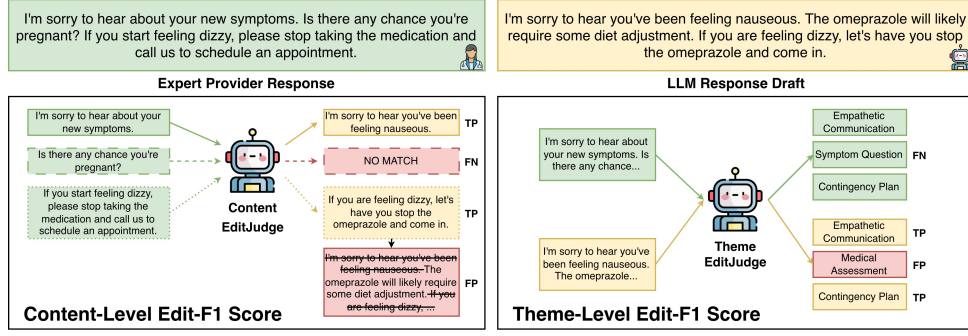


Figure 2: The EditJudge Evaluation Framework for evaluating LLM response drafts. The content-level edit-F1 score identifies matching content in the response draft ( $EM$ , i.e. true positives), along with expected deletions ( $ED$ , false positives) and expected additions ( $EA$ , false negatives) needed in order to align the LLM response draft with the clinician’s desired response. The theme-level edit score identifies matching themes, serving as a relaxed evaluation of the theme-level alignment.

score aims to identify how many *expected additions* ( $EA$ ) and *expected deletions* ( $ED$ ) are needed from the clinician, in order to unify  $r_d$  with  $r_e$ . Matching content in the response draft  $r_d$  is referred to as an *expected match* ( $EM$ ), meaning we would not expect the clinician to have to rewrite that content in order to achieve their desired response  $r_e$ , saving the clinician time and achieving reliability via LLM response drafting.

We give our algorithm for counting  $EA$ ,  $ED$ , and  $EM$  in Algorithm 1 in Appendix C. This algorithm splits an expert-written response  $r_e$  into atomic elements (sentences), then for each element uses a fine-tuned judge LLM (content-level edit-Judge) to either identify expected matches  $EM$  in the response draft  $r_d$ , or expected additions  $EA$  to the response draft to achieve that element. The content-level editJudge takes as input a sentence from the expert-written response  $s_e$  and the LLM-drafted response  $r_d$ , and outputs either the matching content from the LLM-drafted response  $s_d$ , or “NO MATCH” if there is no matching content. Finally, this algorithm identifies expected deletions  $ED$  in the response draft by quantifying the remaining amount of unmatched content. By treating expected matches, expected additions, and expected deletions as true positives, false negatives, and false positives respectively, we calculate recall, i.e. the percentage of the expert-written response  $r_e$  which does not need to be added to  $r_d$ , and precision, i.e. the percentage of the LLM response draft  $r_d$  which does not need to be removed. We calculate the harmonic mean of the content-level recall and precision scores (i.e.  $F_1$ ) and call this the **content-level edit-F1 score**. Assuming additions and deletions are evenly-weighted, content-level edit-F1 gives the expected reduction in editing load

for the clinician by using the LLM response draft.

We evaluate 10 variations of content-level editJudge models, and select a fine-tuned LLama-3-8B-Instruct model for use in our experiments in Section 5. This editJudge model achieves 96% agreement with expert human annotators, including 92% overlap with expert-annotated matching content decisions. We discuss data annotation, training, and evaluation of editJudge models in Appendix C.

### 3.2 Theme-Level edit-F1 Score

Given a clinician response  $r_e$  and an LLM response draft  $r_d$ , the theme-level edit-F1 score aims to identify the higher-level themes in the clinician response  $r_e$  which are correctly matched by the themes in the LLM response draft  $r_d$ . To identify themes in each response, we develop and evaluate a theme-level editJudge model. Given a sentence from either the clinician response  $s_e \in r_e$  or the LLM drafted response  $s_d \in r_d$ , the theme-level editJudge model assigns a theme label  $l_s$ . Predicting clinician response themes is a 9-class multi-label classification task, as there are 8 high-level themes (see Table 1) and an “Other” class to capture miscellaneous themes not captured in the main 8 classes. Using the theme labels  $l_{s_d}$  assigned to sentences  $s_d$  from the LLM drafted response  $r_d$  as predictions for the theme labels  $l_{s_e}$  assigned to sentences  $s_e$  from the clinician response  $r_e$ , the **theme-level edit-F1 score** is the micro average  $F_1$  of theme predictions. We develop and evaluate a fine-tuned theme-level editJudge theme classification model which achieves an  $F_1$  score of 0.82 on expert-annotated dataset (details in Appendix C).



## 4 Experimental Setup

We are interested in how LLMs might be more closely aligned with expert clinicians, to increase the reliability and responsibility of LLMs in response drafting (RQ3). We describe the models used in our evaluation, a measure of inter-annotator predictability (IAP) to contextualize our results, and a measurement of theme frequency in clinician and LLM response drafts.

### 4.1 Models and Adaptation Methods

#### 4.1.1 Local and Frontier LLMs

Locally-hosted LLMs are often preferable in clinical settings due to the sensitive nature of protected health information (PHI) and the frequency with which PHI occurs in patient portal messages (Sal-lam et al., 2023; Zhou et al., 2023). Token throughput and hosting memory constraints are also important considerations (Lorencin et al., 2025). As such, we are interested in evaluating 7-8b parameter LLMs on the response drafting task. We use three models: (i) the instruction-tuned Llama3-8B model (AI@Meta, 2024), (ii) a healthcare-specific version of the same model Aloe-8B (Gururajan et al., 2024), and (iii) Qwen3-8B (Team, 2025) from a different model family. We also test three commercial models on the SyPPM dataset, our public dataset: (i) Claude 4.5 Sonnet (Anthropic, 2025), (ii) Gemini 2.5 Pro (Comanici et al., 2025), and (iii) GPT-OSS (Agarwal et al., 2025).

#### 4.1.2 Adaptation Techniques

We are interested in exploring several avenues for aligning LLMs with expert clinicians to improve reliability and responsibility. We briefly describe each adaptation strategy here, providing full details in Appendix D, and prompts in Appendix H.

**0-Shot.** Minimally-guided responses from each model are evaluated to identify how closely-aligned the LLM is with expert clinicians.

**Thematic.** Some prior work has shown that prompting techniques can improve LLM performance on patient messaging tasks (Genovese et al., 2025). We are interested in whether the themes derived in Section 2.1 can align LLMs more closely with expert clinicians. The thematic prompt includes a brief explanation of each of the 8 themes, to guide the LLM with context.

**RAG.** Retrieval augmented generation has been used in other patient messaging tasks to improve style and content of LLM responses (Chen et al.,

2025). We perform 5-shot RAG prompting.

**SFT.** Supervised fine-tuning on prior patient-clinician conversations has proven to be an effective way to adapt LLM for patient message response drafting (Liu et al., 2024). We perform SFT using all 144k training messages.

**TADPOLE.** We develop a novel Thematic Agentic Direct Preference Optimization for Learning Enhancement strategy for creating theme-driven preference training data for DPO (Rafailov et al., 2023). TADPOLE uses response enhancement agents designed for each theme derived in Section 2.1. We test several preference pair creation strategies (details in Appendix D.3), and report the results of models trained using the best-performing strategy.

### 4.2 Inter-Annotator Predictability

A key consideration when evaluating the LLM-clinician alignment is how closely-aligned clinicians are with each other. Clinician alignment may vary based on experience factors (e.g. role, years of experience, specialty), personality factors (e.g. writing style), and interpersonal factors (e.g. relationship with the patient). We gather 3 expert responses to 40 samples from the SyPPM dataset to quantify inter-annotator predictability (IAP). We calculate IAP using the editJudge framework to compare inter-human alignment on patient message response drafting. IAP gives us a measure of how useful a different clinician’s response might be when used as a response draft. We also report inter-annotator agreement of ground truth in Tables 10-11 in Appendix E.1.

### 4.3 Estimated Theme Frequency

As manually annotating sentence themes in all responses would be inefficient, we use our empirically-validated sentence-level theme classifier (theme-level editJudge LLM, achieves 0.82 F1 on test set in Appendix C) to classify themes in all clinician responses (i.e., ground truth) and all LLM response drafts to estimate thematic tendencies (see Table 5).

## 5 Results

We evaluate six LLMs and five adaptation techniques on the IPPM and SyPPM response drafting evaluation datasets and discuss our findings. Due to space constraints, we discuss results on the SoCPPM dataset in Appendix F.

Table 3 contains both content-level and theme-level edit-F1 scores, **averaged** across the three lo-

Dataset	Model	Content-Level			Theme-Level		
		Precision	Recall	Edit-F1	Precision	Recall	Edit-F1
IPPM	0-Shot	0.07±0.02	0.26±0.04	0.10±0.02	0.49±0.03	0.74±0.03	0.58±0.02
	Theme	0.06±0.01	<b>0.30±0.05</b>	0.09±0.01	0.47±0.01	<b>0.80±0.02</b>	0.58±0.01
	RAG	0.11±0.03	0.30±0.17	0.13±0.01	0.48±0.20	0.66±0.09	0.56±0.02
	SFT	<b>0.15±0.01</b>	0.16±0.00	<b>0.14±0.01</b>	<b>0.64±0.01</b>	0.57±0.01	<b>0.60±0.01</b>
	TADPOLE	0.13±0.01	0.18±0.01	<b>0.14±0.01</b>	0.54±0.00	0.65±0.02	0.59±0.01
SyPPM	0-Shot	0.12±0.04	0.31±0.03	0.16±0.04	0.47±0.02	0.46±0.03	0.47±0.02
	Theme	0.11±0.00	<b>0.33±0.10</b>	0.15±0.02	0.50±0.01	<b>0.58±0.01</b>	0.54±0.00
	RAG	0.17±0.08	0.28±0.07	0.18±0.05	0.47±0.03	0.43±0.02	0.45±0.02
	SFT	<b>0.22±0.01</b>	0.17±0.01	0.18±0.0	<b>0.64±0.01</b>	0.41±0.01	0.50±0.01
	TADPOLE	0.21±0.01	0.20±0.02	<b>0.20±0.01</b>	0.62±0.01	0.54±0.02	<b>0.58±0.01</b>
	<i>Gemini</i>	<i>0.20</i>	<i>0.43</i>	<i>0.26</i>	<i>0.58</i>	<i>0.69</i>	<i>0.64</i>
IAP		0.26	0.25	0.24	0.61	0.63	0.62

Table 3: Edit-F1 scores for LLM adaptations on the IPPM and SyPPM patient message response drafting datasets. Each model adaptation is performed on three underlying LLMs, we report scores as average±standard deviation. We report content-level precision, recall, and edit-F1 (Section 3.1), as well as theme-level precision, recall, and edit-F1 (Section 3.2). We include the best commercial model (Gemini + theme prompting) scores on the publicly-available SyPPM dataset. Finally, we report content-level inter-annotator predictability (IAP), comparing LLM performance and expert human alignment.

cal LLMs described in Section 4.1.1, alongside standard deviation. Table 4 contains content- and theme-level edit-F1 scores for Claude 4.5 Sonnet, Gemini 2.5 Pro, and GPT-OSS reasoning models, using both 0-shot and thematic prompting adaptation. In Tables 3 and 4 we report micro average precision, recall, and edit-F1 at the content and theme levels. Table 5 contains theme frequencies for clinician responses and adapted LLM drafts, averaged across all evaluation datasets.

## 5.1 Content-Level Results

**Usefulness of Thematic Context:** We find that fine-tuned models achieve highest precision, theme-prompted models achieve highest recall, and the TADPOLE adaptation strategy offers the best blend of precision and recall with the highest average content-level edit-F1 scores. We find that added context improves LLM alignment with individual clinicians, and that edit-F1 performance generally scales with the amount of added context. Examining theme-specific content-level recall (Table 14 in Appendix F.2), TADPOLE-adapted models blend precision with empathetic communication content (0.30 average recall vs 0.28 average among other adaptations) and contingency planning content (0.27 vs 0.21)—two themes which tend to appear more in “ideal” response drafts. Among commercial models, thematic prompting adaptation improves performance of all three LLMs. We find that the best frontier-level model in our evaluation is Gemini 2.5 Pro adapted with thematic prompting, achieving 0.26 content-level and 0.64 theme-level edit-F1. Our single best-performing TADPOLE model (Qwen3-8B trained on the “cor-

rupted” preference pairs<sup>3</sup>) achieves comparable performance (0.25 content-level edit-F1 score) to the best-performing frontier model (Gemini 2.5 Pro + theme prompt, 0.26). Our evaluation suggests that using one of these models in patient message response drafting would lead to a 25-26% reduction in clinician edits.

**Epistemic Uncertainty:** Individual variation stemming from epistemic uncertainty is often observed in medicine (Han et al., 2021), including patient message response drafting (Chen et al., 2024b; Garcia et al., 2024; Laukka et al., 2020; Baxter et al., 2024; English et al., 2024a). Our results support this finding (see Tables 10-12 in Appendix E.1). When one clinician’s responses are used as drafts for another clinician, we find an average content-level edit-F1 score of 0.24—meaning that using another clinician’s response as a draft only reduces clinician edits by 24%. This indicates substantial epistemic uncertainty at the content level of clinician responses, i.e., LLMs specialized at the task level are subject to performance loss due to inter-clinician variation in judgment and preferences. This highlights the need for LLMs to be specialized at the expert level in order to further improve clinician efficiency with response drafts.

## 5.2 Theme-Level Results

**LLMs Generate Quality Empathetic Content:** Evaluating at the theme level shows that LLMs are capable of generating some themes accurately, while other themes are more challenging. For example, LLMs tend to generate the empathetic

<sup>3</sup>See TADPOLE results in Table 9 in Appendix D.3

Prompt	Model	Content-Level			Theme-Level		
		Pr	Re	Edit-F1	Pr	Re	Edit-F1
0-Shot	GPT	0.03	0.21	0.05	0.45	<b>0.64</b>	0.53
	Gemini	0.17	<b>0.40</b>	0.23	0.52	0.56	<b>0.54</b>
	Claude	<b>0.20</b>	0.38	<b>0.25</b>	<b>0.52</b>	0.54	0.53
	Avg	<i>0.13</i>	<i>0.33</i>	<i>0.18</i>	<i>0.50</i>	<i>0.58</i>	<i>0.53</i>
Theme	GPT	0.06	0.30	0.09	0.49	<b>0.77</b>	0.60
	Gemini	<b>0.20</b>	<b>0.43</b>	<b>0.26</b>	0.56	0.69	<b>0.64</b>
	Claude	0.16	0.37	0.22	<b>0.58</b>	0.69	0.63
	Avg	<i>0.14</i>	<i>0.37</i>	<i>0.19</i>	<i>0.54</i>	<i>0.72</i>	<i>0.62</i>
IAP		<i>0.26</i>	<i>0.25</i>	<i>0.24</i>	<i>0.61</i>	<i>0.63</i>	<i>0.62</i>

Table 4: Edit-F1 results for Claude 4.5 Sonnet, Gemini 2.5 Pro, and GPT-OSS reasoning models on the publicly-available SyPPM evaluation dataset. We evaluate each model using 0-shot and thematic prompts, and average scores for each prompt. We report precision, recall, and edit-F1 at both the content and theme levels. We report content-level IAP, comparing LLM performance and expert human alignment at the content level.

communication theme frequently (Table 5), and they perform well overall at generating this theme—e.g. TADPOLE-adapted models achieve an average theme-level edit-F1 score of 0.99 on the empathetic communication theme in SyPPM (see Table 15 in Appendix F.2). This finding supports English et al. (2024b), which finds that nurses report that LLM response drafts improve empathy and tone. On the contrast, Table 5 shows that unaligned LLMs will rarely ask follow-up questions. Unaligned LLMs tend to be misaligned with clinicians on question asking themes—e.g. 0-shot models achieve only 0.17 and 0.08 average theme-level edit-F1 scores on SyPPM symptom and medication question-asking themes (Table 15). Contextual adaptation greatly improves LLM performance at question asking, with TADPOLE-adapted LLMs improving to 0.79 and 0.49 average theme-level edit-F1 scores on SyPPM symptom and medication question-asking themes.

**Individuality of Expert Clinicians:** In general, IAP is much higher at the theme level than at the content level, indicating that theme-level alignment is a more achievable goal when drafting clinician responses. However, some individual themes have very low IAP, e.g. treatment planning (0.07 IAP theme-level edit-F1 score in Table 15) and contingency planning (0.06). Discussions with various clinicians, including our annotators, highlight that different clinicians tend to think differently about how content will be perceived by patients – e.g. some clinicians indicate that the benefits of providing contingency plans do not outweigh the burden it places on patients. This again underscores the need for LLMs to be able to be adapted at an individual level, in order to draft useful responses for individual clinicians with different roles (triage nurse, medical assistant, residents), specialties (in-

ternal medicine vs family medicine), years of experiences, and preferences. Individual alignment is vital for *reliable* and *responsible* use of LLM-mediated tools in high-stakes professional workflows like healthcare.

### 5.3 Implications of Results

**Reliable LLM Adaptation:** We find that unadapted LLMs tend to generate medical assessment themes more successfully than contextually-adapted LLMs. This is supported by our estimate of theme proportions (Table 5), which finds that unadapted LLMs generate far more medical assessment and treatment planning themes than clinicians and contextually-adapted LLMs. These themes cover utterances related to medical decision making and communication, i.e., explaining test results, symptoms, and potential diagnoses; and recommending various forms of treatment. Intuitively, unadapted LLMs generate these themes more frequently as they relate to general LLM alignment principles, e.g., safety and helpfulness (Ji et al., 2023). However, such behavior can lead to over-diagnosis and over-treatment (Kale and Korenstein, 2018), an emerging concern about using AI in medicine (Scott et al., 2024). Responses drafted by unadapted models also tend to be longer (Garcia et al., 2024; Hu et al., 2025; Tai-Seale et al., 2024), which may introduce more cognitive burden for clinicians, defeating the purpose of saving clinicians’ time spent in responding to messages.

**Importance of Evaluation:** Our evaluation measures how many edits a clinician would make to the LLM-generated draft before sending the response. This is different from the goal of measuring response quality along pre-defined axes, and influences our decision to define a ground truth as a single clinician response, rather than a strategy

Response	Emp	Sym Q	Med Q	Assess	Plan	Logis	Coord	Cont	Oth
Clinicians	0.85	0.36	0.30	0.34	0.19	0.56	0.45	0.22	0.02
0-Shot	0.94	0.02	0.05	0.89	0.82	0.59	0.78	0.18	0.14
Theme	0.95	<b>0.26</b>	0.13	0.94	0.79	0.64	0.82	0.18	0.20
RAG	<b>0.77</b>	0.01	0.05	0.79	0.65	<b>0.56</b>	<b>0.69</b>	0.11	0.19
SFT	0.97	0.02	0.02	0.23	0.26	0.38	<b>0.69</b>	0.02	<b>0.02</b>
TADPOLE	0.99	0.29	<b>0.20</b>	<b>0.28</b>	<b>0.31</b>	0.36	0.83	<b>0.25</b>	0.01

Table 5: Proportion of responses containing different thematic content, found in responses written by clinicians and various model adaptations. Clinician theme proportion is averaged across the IPPM, SyPPM, and SoCPPM datasets. LLM adaptation theme proportion is averaged over the three underlying LLMs as well as the three datasets. Bold proportions highlight the adaptation that was closest to clinician proportions.

such as rubric-based evaluation (Arora et al., 2025) or surveying expert feedback (Liu et al., 2024) on a generated response. Results from our targeted evaluation highlight the challenge of aligning models with individual clinicians’ judgment, tone, and preferences when responding to patients. It also yields insights for future work to explore alternatives to response drafting to improve clinician efficiency, e.g., suggesting clinicians theme-based “nudges” — rather than content— for themes with higher epistemic uncertainty.

## 6 Related Works

**Patient Message Response Drafting.** Several works have studied the usefulness of LLMs in drafting clinician responses to patient messages. Most evaluate drafts via only clinician feedback, limiting the scale of evaluation, and employ only 0-shot frontier-level LLMs (most commonly OpenAI GPT-4) (Biro et al., 2025; Sharma et al.; English et al., 2024b; Small et al., 2024; Tai-Seale et al., 2024; Hu et al., 2025; Bootsma-Robroeks et al., 2025). Our work extends prior work in two ways: (1) large scale evaluation of adapted LLMs and (2) inclusion of EHR data with message to situate generated responses. Results from prior studies are mixed, with some showing the potential of LLM drafts in promoting empathy and giving health advice (English et al., 2024b; Eschler et al., 2015), while others show that there is room for improvement in LLM draft completeness, tone, and simplicity (Garcia et al., 2024; Small et al., 2024; Chen et al., 2024b). Among studies that go beyond 0-shot evaluation, Hu et al. (2025) and Kim et al. (2024) explore prompting strategies to improve LLM response drafts. Our thematic prompting strategy builds on prior work by incorporating a more granular-level understanding of LLM behavior in response generation across the constituent

themes of a clinician’s response. Liu et al. (2024) is perhaps most similar to our work in that they perform SFT of a Llama model and evaluate on a small test set (n=10) using clinician feedback and BERTScore. Our work in developing a thorough automated evaluation framework aims to build on this by enabling larger-scale automated evaluation. Our focus on large-scale evaluation enables deeper insight into the risks and benefits of LLM use in patient message response drafting.

**Evaluation based on LLM-As-Judge.** The use of LLMs as judges of LLM-generated content has grown significantly in recent years (Li et al., 2024; Lin and Chen, 2023; Li et al., 2025; Bavaresco et al., 2025), including in healthcare text generation contexts (Croxford et al., 2025; Bedi et al., 2025; Zhao et al., 2025; Krolik et al., 2024). Perhaps most similar to our work, Croxford et al. (2025) introduce an LLM-as-Judge framework for evaluating generated EHR summaries and use a rubric-based evaluation. In contrast, our novel edit-F1 framework is designed to estimate edit load, i.e., expected deletions/additions to LLM-generated draft.

## 7 Conclusion

We have evaluated LLMs on the patient message response drafting task. We have developed a set of clinician response themes and used these to develop a novel evaluation framework for assessing clinician editing load given LLM response drafts. We have performed a large-scale evaluation of contextually-adapted LLMs and frontier LLMs, finding that contextual adaptation improves LLM performance. We highlight that individual clinician preferences vary significantly, and that adaptation of LLMs to individual clinicians is required to further increase the reliability and responsibility of LLM use for patient message response drafting.



## 8 Limitations

**Dataset** Our data is drawn from a single hospital system and patient portal platform, which may limit generalizability to other healthcare settings with different workflows, patient populations, and communication norms. This is a rural hospital system. Future work may explore safety, bias and robustness of adapted LLMs in such settings. The judge LLM and thematic classification models we developed in Section 3 are tuned specifically for our evaluation datasets and would require additional validation before application in other contexts (Wu and Aji, 2025; Chen et al., 2024a).

**Automated Evaluation** Some prior evaluations of minimally-adapted LLM use in the patient portal suggest that reduction in clinician time via LLM response drafting is minimal (Hu et al., 2025; Tai-Seale et al., 2024; Bootsma-Robroeks et al., 2025). Our evaluation seeks to fill a critical research gap by automating the evaluation of how much a clinician would edit these responses, which we hope will enable progress towards better LLM alignment with individual clinicians and meaningful reduction in clinician workload. Our evaluation suggests that best-performing response drafting LLMs would reduce clinician edits by 25-26%. This is a modest reduction, potentially due to the complexity of our data which covers real messages from general primary care and a wide range of medical topics and patient intents. Our focus on this automated evaluation limits us from performing in-depth qualitative analysis by clinicians and patients. While our hospital network is not an early adopter of LLM use in clinic which prohibits the use of our models for live patient messages, we hope to perform further studies with clinicians and patients in future work.

**Ethical Considerations** Real patient data used in our evaluations is highly sensitive, and extreme caution should be taken when using LLMs on real patient data to ensure patient privacy. We carefully design our evaluations to promote the responsible use of this data in our evaluation. Our data cleaning process ensures sensitive patients, e.g. patients under the age of 18, were not included in our final dataset. We host all real data on a secure server and perform all IPPM and SoCPPM experiments on this server. We only use proprietary LLMs on semi-synthetic data (SyPPM) which was created via completely de-identified patient charts and messages.

## References

- Sandhini Agarwal, Lama Ahmad, Jason Ai, Sam Altman, Andy Applebaum, Edwin Arbus, Rahul K Arora, Yu Bai, Bowen Baker, Haiming Bao, and 1 others. 2025. gpt-oss-120b & gpt-oss-20b model card. *arXiv preprint arXiv:2508.10925*.
- AI@Meta. 2024. [Llama 3 model card](#).
- Anthropic. 2025. [Claude 4.5 sonnet](#).
- Rahul K Arora, Jason Wei, Rebecca Soskin Hicks, Preston Bowman, Joaquin Quiñero-Candela, Foivos Tsimpourlas, Michael Sharman, Meghan Shah, Andrea Vallone, Alex Beutel, and 1 others. 2025. Healthbench: Evaluating large language models towards improved human health. *arXiv preprint arXiv:2505.08775*.
- Elizabeth Baltaro, Wendy Henderson, and Karen M Goldstein. 2022. Patient electronic messaging: 12 tips to save time. *Family Practice Management*, 29(6):5–9.
- Anna Bavaresco, Raffaella Bernardi, Leonardo Bertolazzi, Desmond Elliott, Raquel Fernández, Albert Gatt, Esam Ghaleb, Mario Giulianelli, Michael Hanna, Alexander Koller, Andre Martins, Philipp Mondorf, Vera Neplenbroek, Sandro Pezzelle, Barbara Plank, David Schlangen, Alessandro Suglia, Aditya K Surikuchi, Ece Takmaz, and Alberto Testoni. 2025. [LLMs instead of human judges? a large scale empirical study across 20 NLP evaluation tasks](#). In *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 238–255, Vienna, Austria. Association for Computational Linguistics.
- Sally L Baxter, Christopher A Longhurst, Marlene Millen, Amy M Sitapati, and Ming Tai-Seale. 2024. Generative artificial intelligence responses to patient messages in the electronic health record: early lessons learned. *JAMIA open*, 7(2):ooae028.
- Suhana Bedi, Hejie Cui, Miguel Fuentes, Alyssa Unell, Michael Wornow, Juan M. Banda, Nikesh Kotecha, Timothy Keyes, Yifan Mai, Mert Oez, Hao Qiu, Shrey Jain, Leonardo Schettini, Mehr Kashyap, Jason Alan Fries, Akshay Swaminathan, Philip Chung, Fateme Nateghi, Asad Aali, and 62 others. 2025. Medhelm: Holistic evaluation of large language models for medical tasks. *ArXiv*, abs/2505.23802.
- Joshua M. Biro, Jessica L. Handley, J. Malcolm McCurry, Adam Visconti, Jeffrey M Weinfeld, J. Gregory Trafton, and Raj M. Ratwani. 2025. Opportunities and risks of artificial intelligence in patient portal messaging in primary care. *NPJ Digital Medicine*, 8.
- Charlotte MHT Bootsma-Robroeks, Jessica D Workum, Stephanie CE Schuit, Anne Hoekman, Tarannom Mehri, Job N Doornberg, Tom P van der Laan, and Rosanne C Schoonbeek. 2025. AI-generated draft replies to patient messages: exploring effects of implementation. *Frontiers in Digital Health*, 7:1588143.

- Virginia Braun and Victoria Clarke. 2006. Using thematic analysis in psychology. *Qualitative research in psychology*, 3(2):77–101.
- Jeffrey Budd. 2023. Burnout related to electronic health record use in primary care. *Journal of Primary Care & Community Health*, 14.
- Guiming Hardy Chen, Shunian Chen, Ziche Liu, Feng Jiang, and Benyou Wang. 2024a. Humans or llms as the judge? a study on judgement biases. *arXiv preprint arXiv:2402.10669*.
- Shan Chen, Marco Guevara, Shalini Moningi, Frank Hoebers, Hesham Elhalawani, Benjamin H Kann, Fallon E Chipidza, Jonathan Leeman, Hugo JWL Aerts, Timothy Miller, and 1 others. 2024b. The effect of using a large language model to respond to patient messages. *The Lancet Digital Health*, 6(6):e379–e381.
- Wenyuan Chen, Fateme Nateghi Haredasht, Kameron C Black, Francois Grolleau, Emily Alsentzer, Jonathan H Chen, and Stephen P Ma. 2025. Retrieval-augmented guardrails for ai-drafted patient-portal messages: Error taxonomy construction and large-scale evaluation. *arXiv preprint arXiv:2509.22565*.
- Gheorghe Comanici, Eric Bieber, Mike Schaeckermann, Ice Pasupat, Naveen Sachdeva, Inderjit Dhillon, Marcel Blistein, Ori Ram, Dan Zhang, Evan Rosen, Luke Marris, Sam Petulla, Colin Gaffney, Asaf Aharoni, Nathan Lintz, Tiago Cardal Pais, Henrik Jacobsson, Idan Szpektor, Nan-Jiang Jiang, and 3416 others. 2025. [Gemini 2.5: Pushing the frontier with advanced reasoning, multimodality, long context, and next generation agentic capabilities](#). *Preprint*, arXiv:2507.06261.
- Emma Croxford, Yanjun Gao, Elliot First, Nicholas Pellegrino, Miranda Schnier, John Caskey, Madeline Oguss, Graham Wills, Guanhua Chen, Dmitriy Dligach, and 1 others. 2025. Automating evaluation of ai text generation in healthcare with a large language model (llm)-as-a-judge. *medRxiv*, pages 2025–04.
- Cathal Doyle, Laura Lennox, and Derek Bell. 2013. A systematic review of evidence on the links between patient experience and clinical safety and effectiveness. *BMJ open*, 3(1):e001570.
- Eden English, Janelle Laughlin, Jeffrey Sippel, Matthew DeCamp, and Chen-Tan Lin. 2024a. Utility of artificial intelligence–generative draft replies to patient messages. *JAMA Network Open*, 7(10):e2438573–e2438573.
- Eden F. English, Janelle Laughlin, Jeffrey Sippel, Matthew DeCamp, and Chen-Tan Lin. 2024b. [Utility of artificial intelligence–generative draft replies to patient messages](#). *JAMA Network Open*, 7.
- Jordan Eschler, Leslie S Liu, Lisa M Vizer, Jennifer B McClure, Paula Lozano, Wanda Pratt, and James D Ralston. 2015. Designing asynchronous communication tools for optimization of patient-clinician coordination. In *AMIA Annual Symposium Proceedings*, volume 2015, page 543.
- Patricia Garcia, Stephen P. Ma, Shreya J. Shah, Margaret Smith, Yejin Jeong, Anna Devon-Sand, Ming Tai-Seale, Kevin Takazawa, Danyelle Clutter, Kyle Vogt, Carlene Lugtu, Matthew Rojo, Steven Lin, Tait Shanafelt, Michael A. Pfeffer, and Christopher Sharp. 2024. [Artificial intelligence–generated draft replies to patient inbox messages](#). *JAMA Network Open*, 7.
- Joseph Gatto, Parker Seegmiller, Timothy E. Burdick, Inas S. Khayal, Sarah DeLozier, and Sarah M. Preum. 2025. [Follow-up question generation for enhanced patient-provider conversations](#). In *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 25222–25240, Vienna, Austria. Association for Computational Linguistics.
- Joseph Gatto, Parker Seegmiller, Timothy E. Burdick, and Sarah M. Preum. 2024. In-context learning for preserving patient privacy: A framework for synthesizing realistic patient portal messages. In *Machine Learning for Health (ML4H) Findings*.
- Ariana Genovese, Sahar Borna, Cesar Abraham Gomez-Cabello, Syed Ali Haider, Srinivasagam Prabha, Maissa Trabilisy, Cui Tao, Keith T Aziz, Peter M. Murray, and AJ Forte. 2025. [Artificial intelligence for patient support: Assessing retrieval-augmented generation for answering postoperative rhinoplasty questions](#). *Aesthetic surgery journal*.
- Ashwin Kumar Gururajan, Enrique Lopez-Cuena, Jordi Bayarri-Planas, Adrian Tormos, Daniel Hinojos, Pablo Bernabeu-Perez, Anna Arias-Duart, Pablo Agustin Martin-Torres, Lucia Urcelay-Ganzabal, Marta Gonzalez-Mallo, Sergio Alvarez-Napagao, Eduard Ayguadé-Parra, and Ulises Cortés Dario Garcia-Gasulla. 2024. [Aloe: A family of fine-tuned open healthcare llms](#). *Preprint*, arXiv:2405.01886.
- Paul KJ Han, Tania D Strout, Caitlin Gutheil, Carl Germann, Brian King, Eirik Ofstad, Pål Gulbrandsen, and Robert Trowbridge. 2021. How physicians manage medical uncertainty: a qualitative study and conceptual taxonomy. *Medical decision making*, 41(3):275–291.
- Paul A Harris, Robert Taylor, Robert Thielke, Jonathon Payne, Nathaniel Gonzalez, and Jose G Conde. 2009. [Research electronic data capture \(redcap\)—a metadata-driven methodology and workflow process for providing translational research informatics support](#). *Journal of Biomedical Informatics*, 42(2):377–381.
- Di Hu, Yawen Guo, Yiliang Zhou, Lidia Flores, and Kai Zheng. 2025. A systematic review of early evidence on generative ai for drafting responses to patient messages. *npj Health Systems*, 2(1):27.

- Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, Weizhu Chen, and 1 others. 2022. Lora: Low-rank adaptation of large language models. *ICLR*, 1(2):3.
- Jiaming Ji, Mickel Liu, Josef Dai, Xuehai Pan, Chi Zhang, Ce Bian, Boyuan Chen, Ruiyang Sun, Yizhou Wang, and Yaodong Yang. 2023. Beavertails: Towards improved safety alignment of llm via a human-preference dataset. *Advances in Neural Information Processing Systems*, 36:24678–24704.
- Thinking Machines John Schulman. 2025. [Lora without regret](#).
- Minal S Kale and Deborah Korenstein. 2018. Overdiagnosis in primary care: framing the problem and finding solutions. *Bmj*, 362.
- Jiyeong Kim, Michael L Chen, Shawheen J Rezaei, April S Liang, Susan M Seav, Sonia Onyeka, Julie J Lee, Shivam C Vedak, David Mui, Rayhan A Lal, and 1 others. 2024. Perspectives on artificial intelligence-generated responses to patient messages. *JAMA Network Open*, 7(10):e2438535–e2438535.
- Jack Krolik, Herprit Mahal, Feroz Ahmad, Gaurav Trivedi, and Bahador Saket. 2024. Towards leveraging large language models for automated medical q&a evaluation. *arXiv preprint arXiv:2409.01941*.
- Elina Laukka, Moona Huhtakangas, Tarja Heponiemi, Sari Kujala, Anu-Marja Kaihlanen, Kia Gluschkoff, and Outi Kanste. 2020. Health care professionals’ experiences of patient-professional communication over patient portals: systematic review of qualitative studies. *Journal of Medical Internet Research*, 22(12):e21623.
- Dawei Li, Bohan Jiang, Liangjie Huang, Alimohammad Beigi, Chengshuai Zhao, Zhen Tan, Amrita Bhat-tacharjee, Yuxuan Jiang, Canyu Chen, Tianhao Wu, and 1 others. 2025. From generation to judgment: Opportunities and challenges of llm-as-a-judge. In *Proceedings of the 2025 Conference on Empirical Methods in Natural Language Processing*, pages 2757–2791.
- Haitao Li, Qian Dong, Junjie Chen, Huixue Su, Yujia Zhou, Qingyao Ai, Ziyi Ye, and Yiqun Liu. 2024. Llms-as-judges: a comprehensive survey on llm-based evaluation methods. *arXiv preprint arXiv:2412.05579*.
- Shuyue Stella Li, Jimin Mun, Faeze Brahman, Pedram Hosseini, Bryceton G Thomas, Jessica M Sin, Bing Ren, Jonathan S Ilgen, Yulia Tsvetkov, and Maarten Sap. Alfa: Aligning llms to ask good questions a case study in clinical reasoning. In *Second Conference on Language Modeling*.
- Yen-Ting Lin and Yun-Nung Chen. 2023. Llm-eval: Unified multi-dimensional automatic evaluation for open-domain conversations with large language models. In *Proceedings of the 5th Workshop on NLP for Conversational AI (NLP4ConvAI 2023)*, pages 47–58.
- Siru Liu, Allison B McCoy, Aileen P Wright, Babatunde Carew, Julian Z Jenkins, Sean S Huang, Josh F Peterson, Bryan Steitz, and Adam Wright. 2024. Leveraging large language models for generating responses to patient messages—a subjective analysis. *Journal of the American Medical Informatics Association*, 31(6):1367–1379.
- Ivan Lorencin, Nikola Tankovic, and Darko Etinger. 2025. Optimizing healthcare efficiency with local large language models. *Intelligent Human Systems Integration (IHSI 2025): Integrating People and Intelligent Systems*, 160(160).
- Kathryn A. Martinez, Rebecca Schulte, Michael B. Rothberg, Maria C Tang, and Elizabeth R. Pfoh. 2023. [Patient portal message volume and time spent on the ehr: an observational study of primary care clinicians](#). *Journal of General Internal Medicine*, 39:566 – 572.
- Frederick North, Kristine E Luhman, Eric A Mallmann, Toby J Mallmann, Sidna M. Tulledge-Scheitel, Emily J North, and Jennifer L. Pecina. 2019. [A retrospective analysis of provider-to-patient secure messages: How much are they increasing, who is doing the work, and is the work happening after hours?](#) *JMIR Medical Informatics*, 8.
- Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D Manning, Stefano Ermon, and Chelsea Finn. 2023. Direct preference optimization: Your language model is secretly a reward model. *Advances in neural information processing systems*, 36:53728–53741.
- Nils Reimers and Iryna Gurevych. 2019. Sentence-bert: Sentence embeddings using siamese bert-networks. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 3982–3992.
- Matthew Sakumoto and Aditi U. Joshi. 2023. [Digital empathy 2.0: Connecting with patients using the written word](#). *Telehealth and Medicine Today*, 8(5).
- Malik Sallam, Nesreen A Salim, Muna Barakat, and Ala’a B Al-Tammemi. 2023. Chatgpt applications in medical, dental, pharmacy, and public health education: A descriptive study highlighting the advantages and limitations. *Narra j*, 3(1):e103.
- Ian A Scott, Anton Van Der Vegt, Paul Lane, Steven McPhail, and Farah Magrabi. 2024. Achieving large-scale clinician adoption of ai-enabled decision support. *BMJ Health & Care Informatics*, 31(1):e100971.
- Rahul Sharma, Pragnya Ramjee, Kaushik Murali, and Mohit Jain. Editing with ai: How doctors refine llm-generated answers to patient queries. In *The Second Workshop on GenAI for Health: Potential, Trust, and Policy Compliance*.
- William R Small, Batia Wiesenfeld, Beatrix Brandfield-Harvey, Zoe Jonassen, Soumik Mandal, Elizabeth R

- Stevens, Vincent J Major, Erin Lostraglio, Adam Szerencsy, Simon Jones, and 1 others. 2024. Large language model-based responses to patients' in-basket messages. *JAMA network open*, 7(7):e2422399–e2422399.
- Moir Stewart. 1995. [Effective physician-patient communication and health outcomes: a review](#). *CMAJ : Canadian Medical Association journal = journal de l'Association medicale canadienne*, 152 9:1423–33.
- Si Sun, Xiaomu Zhou, Joshua C Denny, Trent S Rosenbloom, and Hua Xu. 2013. Messaging to your doctors: understanding patient-provider communications via a portal system. In *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems*, pages 1739–1748.
- Ming Tai-Seale, Sally L Baxter, Florin Vaida, Amanda Walker, Amy Sitapati, Chad Osborne, Joseph Diaz, Nimit Desai, Sophie Webb, Gregory Polston, Teresa Helsten, Erin Gross, Jessica Thackaberry, Ammar Mandvi, Dustin Lillie, Steve Li, Geneen T Gin, Suraj A Achar, Heather Hofflich, and 3 others. 2024. [Ai-generated draft replies integrated into health records and physicians' electronic communication](#). *JAMA Network Open*, 7.
- Qwen Team. 2024. [Qwen2.5: A party of foundation models](#).
- Qwen Team. 2025. [Qwen3 technical report](#). *Preprint*, arXiv:2505.09388.
- Louise Underdahl, Mary Ditri, and Lunthita M Duthely. 2024. [Physician burnout: Evidence-based roadmaps to prioritizing and supporting personal wellbeing](#). *Journal of Healthcare Leadership*, 16:15 – 27.
- Minghao Wu and Alham Fikri Aji. 2025. Style over substance: Evaluation biases for large language models. In *Proceedings of the 31st International Conference on Computational Linguistics*.
- Qi Yan, Zheng Jiang, Zachary Harbin, Preston H Tolbert, and Mark G Davies. 2021. Exploring the relationship between electronic health records and provider burnout: a systematic review. *Journal of the American Medical Informatics Association*, 28(5):1009–1021.
- Justin Zhao, Timothy Wang, Wael Abid, Geoffrey Angus, Arnav Garg, Jeffery Kinnison, Alex Sherstinsky, Piero Molino, Travis Addair, and Devvret Rishi. 2024. Lora land: 310 fine-tuned llms that rival gpt-4, a technical report. *arXiv preprint arXiv:2405.00732*.
- M. Zhao, I. Y. Oh, A. Gupta, S. Cohen-Cutler, K. M. Harmon, A. M. Lai, and B. A. Sisk. 2025. Automating evaluation of llm-generated responses to patient questions about rare diseases. In *medRxiv*.
- Juexiao Zhou, Xiaonan He, Liyuan Sun, Jiannan Xu, Xiuying Chen, Yuetan Chu, Longxi Zhou, Xingyu Liao, Bin Zhang, and Xin Gao. 2023. Skingpt-4: an interactive dermatology diagnostic system with visual large language model. *arXiv preprint arXiv:2304.10691*.



## A Dataset Details

### A.1 Data Collection and Formatting

As described in Section 2, the patient-clinician conversations used in our experiments are collected from a large academic hospital in the Eastern United States. These conversations are sourced from the hospital’s electronic health record (EHR) portal messaging platform. 610k total messages are taken from the secure patient portal between 1/2020 - 9/2024. Our dataset includes messages from primary care, and thus includes a wide range of medical topics. We gather all patient-initiated messages which received a written clinician response to create 146k conversations, i.e. original patient message and response from a clinician. Our final data pool contains 10,105 unique patients, of which 64% are female and 36% are male, with ages ranging between 18-80. Each sample in our data pool consists of a patient message, a clinician response, and a summary of the patient’s chart before the sending of the patient message. We designate 144k conversations from the data pool as training data, and we gather evaluation datasets from the remaining 2k conversations.

Details from throughout the EHR are summarized into four categories. First, the patient’s age range and gender are given as **Demographics**. Next, the patient’s active problems are listed under **Full Active Problem List**. The patient’s recent encounters (with a maximum of 10 entries), including diagnoses, surgeries, visits, etc. are listed under **Recent Encounters**. Finally, a patient’s outpatient medications are summarized in **Medications**. An example de-identified chart from SyPPM is provided in Figure 3.

### A.2 Evaluation Dataset Details

Designating 144k training conversations, we gather evaluation datasets from the remaining 2k conversations. We create three evaluation sets, designed to evaluate LLM alignment with experts according to different standards of care. Each sample in each dataset is a tuple of strings  $\{m, c, r\}$  consisting of a patient message  $m$ , a summary  $c$  of the patient’s EHR chart and a single clinician response  $r$ .

**IPPM** The Ideal Patient Portal Messaging (IPPM) dataset is created to evaluate LLMs in a setting where clinicians do not face the same resource constraints as in the real-world. In this evaluation dataset, ground-truth responses are written by a paid team of 4 expert primary care nurses who

work daily in the patient portal, collected via REDCap surveys (Harris et al., 2009). In addition to giving ample time to write a full response to each message/EHR summary, experts were asked “if you had unlimited time, what would be included in your response to this patient?” To provoke quality responses, clinicians were given a separate text entry box for each of the themes derived in Section 2.1. For example, the *Treatment Contingency Planning* text box included the prompt “please outline a backup/red flag plan for the patient, if applicable.” An example REDCap survey is given in Appendix G for reproducibility. The IPPM dataset is comprised of 300 patient messages and corresponding EHR charts, with one expert clinician response per sample.

**SyPPM** As the other datasets use real patient data containing protected health information (PHI), they are not suitable for public release. We create the Synthetic Patient Portal Messaging (SyPPM) as a public benchmark to promote open-source research in clinician response drafting. We begin by taking 100 semi-synthetic patient portal messages which are created using a small number of de-identified patient portal messages in an in-context synthesis prompt (Gatto et al., 2025, 2024) and pair them with real de-identified patient EHR summaries. Ground-truth responses to each patient message are then provided by a primary care clinician, using the same theme-guided REDCap survey used for IPPM.

**SoCPPM** The Standards of Care Patient Portal Messaging (SoCPPM) dataset is created to evaluate LLMs in a practical setting, in which response drafts are compared with the clinician response which was sent via the secure portal in real time. This dataset is comprised of 300 patient messages and corresponding EHR summaries, where ground-truth responses are sourced from the patient portal. We evaluate LLM response drafts with respect to these real responses from the patient portal to study how LLM responses might perform in real-world settings, against the current standards of care in the patient portal.

## B Thematic Analysis Details

We carefully derive elements of high-quality clinician responses to patient messages. Based on prior work, manual thematic analysis of real patient-clinician conversations, and consultation with expert primary care physicians, nurses, and triage

```

####Demographics####
Age: Between 35 - 40
Gender: Female
####Full Active Problem List####:
DVT prophylaxis - Pelvic floor dysfunction - Chronic left-sided low back pain with left-sided
sciatica - Hypothyroidism due to Hashimoto's thyroiditis - IUD (intrauterine device) in place - H/O
abnormal cervical Papanicolaou smear - Paresthesia of left leg - Hypothyroid - Vitamin D deficiency -
Uncontrolled pain - Healthcare maintenance - CREATED BY INTERFACE - Disorder of muscle,
ligament, and fascia - Postoperative pain - Altered bowel habits
####Recent Encounters (Max 10)####
Diagnoses (Past Year): Cough, persistent - Acne vulgaris - Abnormal uterine bleeding - Encounter for
cosmetic procedure - Changing skin lesion - Acne scarring - Fatigue, unspecified type - COVID-19
ruled out
Diagnoses (Older):
####Medications (Outpatient)####
Active (Start Date Before Message, Not Yet Ended):
-TRETINOIN 0.025 % TOPICAL CREAM
-LEVOTHYROXINE 75 MCG TABLET
-CLINDAMYCIN 1 % LOTION
-SPIRONOLACTONE 100 MG TABLET
-SPIRONOLACTONE 50 MG TABLET
-LORAZEPAM 1 MG TABLET

```

Figure 3: Example de-identified EHR chart summary from our SyPPM patient message response drafting evaluation dataset

nurses, we derive a set of “themes” which can be used to characterize the quality of clinician responses to patient messages (Braun and Clarke, 2006; Sun et al., 2013). Below, we present our hybrid (top-down and bottom-up) approach to identify these themes.

As the quality of patient-clinician communication has a significant impact on patient health outcomes, characterizing quality response elements is important preliminary work for evaluating LLMs on the patient message response drafting task (Stewart, 1995; Doyle et al., 2013). Our goal is to derive themes that should occur in clinician responses to patient messages. We are interested in both empirically-derived themes, sourced from real patient-clinician conversations, as well as theoretically-derived themes, sourced from expert consultation and clinician communication theory (Stewart, 1995; Sakumoto and Joshi, 2023). Empirical themes are indicative of the current standards of care in patient portal communication, whereas theoretical themes may not be found in real-world clinician communication due to time, system, and resource constraints often experienced in asyn-

chronous patient-clinician communication in the patient portal (North et al., 2019; Martinez et al., 2023). We therefore employ a hybrid top-down (theoretical), bottom-up (empirical) approach to identifying themes of quality clinician communication within the patient portal.

### B.1 Theoretical Response Themes

We collaborate with a team of 11 clinicians to identify “ideal” clinician response themes to various patient messages. This iterative process involved 1-1 interviews with 2 primary care physicians and 9 primary care nurses, all of whom regularly interact with patients on the EHR portal from which our data pool (Appendix A) is sourced. These interviews consisted of discussions based on open-ended questions, e.g. “what are your primary goals when writing responses to patient messages in the patient portal?” as well as discussions guided by examples of patient messages, e.g. “what would you want to say to this patient?” or “how would your response vary based on a <specific change> in the patient-initiated message?” Through these interviews, we derive an initial set of theoretical

clinician response themes based on suggested best practices.

## B.2 Empirical Response Themes

Using notes from these conversations as a back-drop, a team of three authors<sup>4</sup>, including a primary care physician, performed a comprehensive, iterative thematic analysis (Braun and Clarke, 2006) using a random sample of 100 patient messages, alongside a summary of the patient’s electronic health record and the clinician’s response. This process involved hand-labeling each sentence-length element of 25 clinician responses with a “frame,” then grouping those frames into “themes,” and repeating this process with new samples. In total we repeated this process four times.

After performing the bottom-up thematic analysis, additional input from two primary care physicians guided the final, comprehensive list of eight clinician response themes comprised of 67 frames. Descriptions and examples of each response theme can be found in Table 1.

## C EditJudge Framework Details

In Figure 2 we see an example of how the content-level and theme-level edit-F1 scores are calculated given a clinician response and an LLM response draft. In Algorithm 1 we give the algorithm for counting expected matches  $EM$ , expected additions  $EA$ , and expected deletions  $ED$  in an LLM-drafted response, in order to calculate content-level edit-F1 scores.

### C.1 Content-Matching Judge Model

Here we describe the process used to fine-tune the content-level editJudge model used in Algorithm 1 to calculate content-level edit-F1. First, three authors hand-label 450 training samples and 50 evaluation samples. Each sample input is a response draft written by the Aloe-8B (Gururajan et al., 2024) 0-shot model, along with a sentence drawn from an expert-written response to a sample from the publicly-available SyPPM evaluation dataset. The annotators either wrote “NO MATCH” if there was no matching content from the response draft, or copy/pasted the matching content from the response draft if applicable. The prompt to identify matches was “if the expert clinician would not have to rewrite this content in order to achieve the same

**Algorithm 1** Counting expected matches  $EM$ , expected additions  $EA$ , and expected deletions  $ED$  in an LLM-drafted response

**Require:**  $r_e$  (expert-written response),  $r_d$  (LLM-drafted response)

**Ensure:**  $EM, ED, EA$

```

1: Split  $r_e$  into atomic elements (sentences)
2: Initialize  $EM \leftarrow 0, EA \leftarrow 0$ 
3: for all sentence  $s_e$  in  $r_e$  do
4:   if MATCH  $s_e$  with content in  $r_d$  then
5:      $EM \leftarrow EM + 1$ 
6:   else
7:      $EA \leftarrow EA + 1$ 
8:   end if
9: end for
10:  $r_d^- \leftarrow$  Remove matching content from  $r_d$ 
11: Split  $r_d^-$  into sentences
12:  $ED \leftarrow$  number of sentences in  $r_d^-$ 
13: return  $EM, ED, EA$ 

```

meaning as their given sentence, this is matching content.” Author annotators were asked to flag all samples about which they were unsure or which required clinical expertise, and two expert clinicians were consulted on these samples to make a final decision.

This matching decision is not always straightforward. For example, in Figure 2 we see that the clinician-written sentence “I’m sorry to hear about your new symptoms” matches with the LLM-drafted sentence “I’m sorry you’ve been feeling nauseous.” While expert clinicians in our evaluation agreed that they would not need to rewrite this LLM-drafted sentence, in order to achieve the same meaning as the clinician-written sentence, this is not always trivial and may vary from clinician to clinician. Examples of clinician-verified matches and non-matches from our training samples can be found in Table 6.

Given a sentence  $s_e$  from an expert-written response  $r_e$  and an LLM-drafted response  $r_d$  the content-level editJudge model was tasked with outputting either the matching content from the LLM draft  $s_d$ , or the string “NO MATCH”. Since the matching content  $s_d$  is later removed from  $r_d$  to identify expected deletions  $ED$ , the output of the judge model  $\hat{s}_d$  must match verbatim to the matching content in the draft  $s_d$  in order to remove  $s_d$  in Algorithm 1. We therefore evaluate the editJudge model by identifying whether it outputs exactly-matching content  $s_d$  identified by the annotators.

<sup>4</sup>Each team member is well-versed in health informatics and qualitative thematic analysis

Clinician Sentence	LLM Draft (Excerpt)	Match Decision	Clinician Reasoning
Thank you for touching base - I'm sorry you have been experiencing these troubling symptoms.	Hi Sarah, Thank you for reaching out and bringing this to my attention. Postmenopausal bleeding can indeed be a sign of something more serious, so it's important that we investigate further.	Thank you for reaching out and bringing this to my attention.	Both sentences thank the patient for being proactive and express sympathy about their symptoms. The clinician would not have to rewrite this component in order to achieve the same meaning.
Please let us know if you have any issues with getting the medication.	Thanks for reaching out to me about your decision to switch. I would recommend a different medication. Please let me know if you have questions or concerns. We can discuss this further in your upcoming appointment.	NO MATCH	While issues with getting the medication may be classified as a concern, the draft is not specific enough and the clinician would have to rewrite.
Have you eaten anything out of the ordinary for you?	Have you experienced any other symptoms, such as stomach pain, bloating, or changes in bowel movements? Are you taking any new medications or supplements that could be causing the nausea? Have you recently changed your diet or experienced any significant stress?	Have you recently changed your diet or experienced any significant stress?	Both symptom-related follow-up questions ask the patient about recent diet changes, and the clinician would not have to rewrite the drafted sentence in order to achieve the same meaning.

Table 6: Selected examples from the content-level editJudge evaluation dataset. The editJudge model is given the LLM draft (an excerpt from each is shown in this table to preserve space) and a sentence from a clinician-written response, and is tasked with outputting either the matching content from the LLM draft, or the string “NO MATCH”. We show two matching decisions, one from the empathetic communication theme and another from the symptom-related follow-up question theme, as well as a close non-match from the contingency planning theme.

Model	Type	Avg Agr	Avg Non-Match	Avg Match	% Match
Qwen2.5-7B-Instruct	0-Shot	0.74	1.00	0.07	0.07
Llama-3-8B-Instruct	0-Shot	0.17	0.11	0.32	0.50
Qwen2.5-7B-Instruct	5-Shot	0.71	0.93	0.14	0.14
Llama-3-8B-Instruct	5-Shot	0.63	0.88	0.00	0.00
Qwen2.5-3B	SFT	0.76	0.97	0.21	0.21
Qwen2.5-3B-Instruct	SFT	0.80	0.94	0.43	0.50
Llama-3.2-3B-Instruct	SFT	0.85	1.00	0.46	0.57
Qwen2.5-7B	SFT	0.87	0.97	0.61	0.71
Qwen2.5-7B-Instruct	SFT	0.89	0.97	0.68	0.71
Llama-3-8B-Instruct	SFT	<b>0.96</b>	<b>1.00</b>	<b>0.84</b>	<b>0.92</b>

Table 7: EditJudge model performance across different configurations. We find that SFT is superior to either 0-shot or 5-shot editJudge models. We find that the best model, the fine-tuned instruction-tuned Llama3-8B model, achieves 96% agreement with clinician-guided author annotations. 84% of the *matching* author annotations were exactly matched by this judge model, and 92% of match decisions contained at least some overlap.



We first identify whether the editJudge model correctly makes the matching decision (either by outputting “NO MATCH” or some substring  $\hat{s}_d$  from the LLM draft  $r_d$ ), and call this **agreement**, i.e. the proportion of evaluation samples on which the judge model makes the correct matching decision. We further score the editJudge model by identifying **non-match agreement**, i.e. the proportion of non-matches correct identified by the judge model, and **match agreement**, the proportion of annotated which are exactly matched by the editJudge model outputs. To get a granular estimate of judge model outputs, we also score **match overlap**, i.e. the proportion of evaluation responses in which editJudge model output  $\hat{s}_d$  and annotated matching content  $s_d$  overlap. We evaluated 6 judge models, testing 0-shot, 5-shot, and supervised fine-tuning adaptation strategies for this content-level matching task.

We see content-level judge results in Table 7. In general, SFT is far superior to either 0-shot or 5-shot judge models. We find that the best model, the instruction-tuned Llama3-8B model (AI@Meta, 2024) fine-tuned on the 450 training samples, achieves 96% agreement with clinician-guided author annotations. 84% of the *matching* author annotations were exactly matched by this judge model, meaning the exact correct content would be removed from the LLM draft  $r_d$  to identify exact expected deletions  $ED$ , and 92% of match decisions contained at least some overlap.

## C.2 Sentence Theme Classification Model

We now similarly describe the fine-tuning the sentence-level theme classification model, used to calculate the theme-level edit-F1 score described in Section 2.1. First, one author hand-labeled 175 training samples and 50 evaluation samples. Each sample was a sentence-length string taken from responses to SyPPM samples generated by the Aloe-8B (Gururajan et al., 2024) 0-shot. Consulting with two expert clinicians, each sample was assigned a single theme label, including the 8 themes and an “Other” label, to set up a 9-class classification task. Example sentences from each theme can be found in Table 1.

Following the results of the content-level edit-Judge training, we choose to fine-tune a Llama3-8B model (AI@Meta, 2024) to perform the sentence classification, where the task is to output the class label (e.g. “Symptom-Related Follow-Up Question”) given the response sentence. Class-wise performance and micro average F1 of this sentence

Theme	F1
Empathetic Communication	0.94
Symptom-Related Follow-Up Questions	1.00
Medication-Related Follow-Up Questions	0.67
Medical Assessment Explanation	0.67
Medical Planning Instruction	0.71
Logistics: Scheduling, Billing, Operations	0.82
Care Coordination	0.80
Contingency Planning	0.67
Other	1.00
<i>Micro Avg</i>	<i>0.82</i>

Table 8: Sentence classification model results. Using a fine-tuned Llama3-8B model (AI@Meta, 2024), we report class-wise performance and micro average F1. We see that the sentence classification model performs well overall, with a micro average F1 of 0.82, and that it predicts all individual classes competently ( $> 0.67$  F1).

classification model are reported in Table 8. We see that the sentence classification model performs well overall, with a micro average F1 of 0.82, and that it predicts all individual classes competently ( $> 0.67$  F1). We note that this task is subjective on some level, given that theme classes are not necessarily disjoint. For example, there are valid reasons to argue that a question such as “have you noticed any diarrhea while on your amoxicillin?” could be both a symptom- and medication-related follow-up question. However, we enforce a single-class label for simplicity in our evaluations.

## D LLM Adaptation Details

As described in Section 4.1.2, here we provide details for the supervised fine-tuning (SFT) and thematic agentic direct preference optimization for learning enhancement (TADPOLE) LLM adaptation strategies which we use in our evaluation in Section 5. Prompts for the 0-shot and thematic adaptations can be found in Appendix H. Further details for the RAG, SFT, and TADPOLE adaptations can be found below.

### D.1 RAG Details

Using the training dataset (144k) as a RAG database, we encode patient messages and EHR summaries using S-BERT<sup>5</sup> (Reimers and

<sup>5</sup>all-MiniLM-L6-v2

Gurevych, 2019), and include the 5 most similar message + EHR strings, along with their real clinician responses in the prompt to guide the LLM, alongside the instruction from the 0-shot prompt.

## D.2 SFT Details

We perform supervised fine-tuning using all 144k training messages. The LLM is trained to output the clinician response  $r$ , given the patient message  $m$  and a summary of the patient’s EHR  $c$  contextualized with the 0-shot prompt (see Appendix H for this prompt).

Each time a model is fine-tuned, both for the SFT models in Section 4.1.2 and for the fine-tuned judge models in Section 3, we train for 1 epoch using a batch size of 4 on a single Nvidia A40 GPU (48GB RAM). We train using low-rank adaptation (LoRA) (Hu et al., 2022) for efficiency, which has shown to be a performant fine-tuning strategy (John Schulman, 2025; Zhao et al., 2024). We use LoRA with rank 8 and an alpha scaling factor of 16. We use the AdamW optimizer with weight decay of 0.01, linear learning rate scheduler with warmup over 10% of the training steps, and gradient clipping at a norm of 1.0. We apply mixed precision training using float16 to optimize memory usage and training speed.

## D.3 TADPOLE Details

For each theme, TADPOLE takes a base response  $r$  and creates both an “enhanced” response  $r^+$  and “corrupted” response  $r^-$  by either adding or removing thematic content from the response. First, we take 8k training samples and generate base responses using the fine-tuned (SFT) model. For *enhancing* a response  $r$  with content from a given theme  $t$ , we use a response enhancing agent to get an enhanced response  $r_t^+$ . Each thematic enhancement agent is a simple 3-shot prompt. For *corrupting* a response  $r$  with content from a given theme  $t$ , we use a standard corruption agent contextualized with the theme  $t$  to obtain a corrupted response  $r_t^-$ . Enhancement prompts and the corruption prompts are developed for and passed to the Qwen2.5-32B-Instruct<sup>6</sup> (Team, 2024) model. We obtain 1k enhanced responses for each theme and 1k corrupted responses for each theme for a total of 8k enhanced, base, and corrupted response  $\{r^+, r, r^-\}$  tuples.

Following Li et al., we test several preference

pair creation strategies using these tuples. **Enhanced** pairs  $\{r^+, r\}$  use enhanced responses and base responses as chosen and rejected responses, respectively. **Corrupted** pairs  $\{r, r^-\}$  choose base responses over corrupted responses. **Hard-Corrupted** pairs  $\{r^+, r^-\}$  choose enhanced responses over corrupted responses. We also investigate a **Blend** which contains an even amount of all three pairs. We perform DPO (Rafailov et al., 2023) on the fine-tuned model using 8k TADPOLE preference pairs. We perform DPO on the SFT model using a beta of 0.01. Similarly with SFT, we perform DPO by training for 1 epoch using a batch size of 1 on a single Nvidia A40 GPU (48GB RAM). We apply mixed precision training using float16 to optimize memory usage and training speed.

We report average content-level and theme-level edit-F1 scores on IPPM for each TADPOLE strategy in Table 9. The hard-corrupted strategy achieves best performance at the content-level, as well as overall when weighting evenly between content- and theme-level edit-F1 scores. Hence we report the results of the models trained on hard-corrupted pairs in Section 5.

## E Measures of Inter-Clinician Variation

### E.1 Inter-Annotator Agreement

Clinician responses to patient messages may vary based on experience factors (e.g. role, years of experience, specialty), personality factors (e.g. writing style), and interpersonal factors (e.g. relationship with the patient). Table 12 gives examples of different clinician responses to the same patient message within the SyPPM dataset.

As noted in Section 4.2, we gather 3 expert responses to 40 samples from the SyPPM dataset. Of the 3 experts, 1 is a primary care physician with 15+ years of experience and 2 are primary care nurses, each with 5+ years of experience. In Section 4.2 we describe how we might use multiple responses to understand inter-annotator predictability (IAP). Here we describe three measures of inter-annotator agreement (IAA), using these same samples.

We are interested in measuring how similarly clinicians would respond to the same patient message in the same conditions. We start by identifying, for each theme, the proportion of patient messages to which all three annotator responses either included that theme (**strict inclusion**), or did not include that theme (**strict exclusion**). Taken together (**strict agreement**), we can estimate the

<sup>6</sup>Qwen/Qwen2.5-32B-Instruct

	Content-Level			Theme-Level		
Pairs	Pr	Re	Edit-F1	Pr	Re	Edit-F1
Blend	0.13	<b>0.19</b>	<b>0.14</b>	0.53	<b>0.65</b>	0.58
Enhanced	0.09	0.14	0.10	0.45	0.62	0.52
Corrupted	0.13	0.16	0.12	<b>0.60</b>	0.62	<b>0.61</b>
Hard-Corrupted	<b>0.13</b>	0.18	<b>0.14</b>	0.54	<b>0.65</b>	0.59
<i>IAP</i>	<i>0.26</i>	<i>0.25</i>	<i>0.24</i>	<i>0.61</i>	<i>0.63</i>	<i>0.62</i>

Table 9: Content-level and theme-level edit-F1 scores for varying TADPOLE preference pair creation strategies on the IPPM dataset. The hard-corrupted strategy achieves best performance at the content-level, as well as overall when weighting evenly between content- and theme-level edit-F1 scores.

IAA Measure	Emp	Sym Q	Med Q	Asse	Plan	Log	Coord	Cont
Strict Inclusion	0.53	0.53	0.20	0.03	0.00	0.57	0.00	0.00
Strict Exclusion	0.00	0.00	0.00	0.33	0.93	0.00	0.33	0.47
Strict Agreement	0.53	0.53	0.20	0.36	0.93	0.57	0.33	0.47

Table 10: Inter-annotator agreement (IAA) measured at the theme-level by identifying cases when all three annotators either included (strict inclusion) or excluded (strict exclusion) each theme in their response. We find that some themes are unanimously found in all clinician responses to most ( $> 50\%$ ) patient messages. Interestingly, we also find that the medical treatment theme is almost never found in any clinician response to most patient messages ( $< 7\%$ ). This speaks to the reluctance of these clinicians to treat patients via the portal, instead favoring information seeking (e.g. follow-up questions) responses.

Clinician	A	B	C
<b>A</b>	1.00	0.51	0.59
<b>B</b>	0.51	1.00	0.45
<b>C</b>	0.59	0.45	1.00

Table 11: Inter-annotator agreement measured at the content-level between clinician pairs using cosine similarity. We find that agreement between clinician pairs varies substantially, with some (clinicians A and C) more aligned than others (clinicians B and C).

extent to which each response theme is clinician-independent.

These theme-level IAA measurements can be found in Table 10. We find that themes such as empathetic communication, symptom-related follow-up questions, and logistical information are unanimously found in all clinician responses to most ( $> 50\%$ ) patient messages in SyPPM. Interestingly, we also find that the medical treatment theme is almost never found in any clinician response to most patient messages ( $< 7\%$ ). This speaks to the reluctance of these clinicians to treat patients via the portal, instead favoring information seeking (e.g. follow-up questions) responses.

For a simpler measure of IAA, we also measure the average pairwise cosine similarity of each clinician’s responses, comparing each pair of clinicians in Table 11. We find that agreement between

clinician pairs varies substantially, with some (clinicians A and C, 0.59) more aligned than others (clinicians B and C, 0.45).

## E.2 Inter-Annotator Predictability

We calculate IAP using both content-level and theme-level edit-F1 scores to enable direct comparison to our model results in Section 5. To estimate the amount of agreement between two expert clinicians in our evaluation framework, we assign the first clinician the role of expert and the second the role of drafting responses. Treating the first clinician’s response as the expert response  $r_e$  and the second’s response as the response draft  $r_d$ , we calculate content-level and theme-level edit-F1 scores using the editJudge described in Section 3. Assigning each ordered pair ( $N = 6$ ) of expert responses as ground-truth responses and response drafts, we compare  $6 \times 40 = 240$  total responses, and take the average results. Tables 3, 14, and 15 give IAP estimates for content-level and theme-level edit-F1 scores, class average content-level recall scores, and class average theme-level edit-F1 scores, respectively.

Patient Message	Clinician A Response	Clinician B Response
I'm not feeling quite myself lately. I've been experiencing some weakness that's making everyday activities a bit more challenging. I was wondering if we could touch base about what might be causing this?	I'm sorry you have been experiencing these troubling symptoms. Could you describe where you experience the weakness? How does it impede your daily activities? Does it come and go? How long has this been going on? Do you have any other symptoms such as dizziness or lightheadedness? Have you checked your blood pressure at home? Have you had any changes to your medications recently? Please call the office to schedule an appointment for urgent evaluation. If your symptoms worsen acutely, including any dizziness or lightheadedness, or syncopal episodes (fainting), you should call 911 and be seen emergently in the ER.	Sorry to hear you aren't feeling well. Are you having any other symptoms? How long have these symptoms been going on? Have you ever had symptoms like this before? Are you having any nausea, vomiting, diarrhea, or constipation? Are you having any fevers? Are you losing weight without trying? Are you having any blood in bowel movements? Are you having abdominal pain? Have you noticed any particular foods that trigger the symptoms? Have you started any new medications or supplements? Have you recently changed dosing or timing of medications you take? Have you tried any medications that have helped? Please give us a call to schedule an appointment. You should be seen in the ED if you have worsening or sudden abdominal pain, severe vomiting, dizziness, chest pain, or shortness of breath.
I'm having a pretty rough time with my seasonal allergies right now. My eyes are itchy, I'm congested, and I just can't seem to stop sneezing. I've been using some over-the-counter meds, but they're not really giving me the relief I need. I was wondering if you could recommend something a bit stronger or if I should come in for an appointment.	I'm sorry you have been experiencing these troubling symptoms. Which medications have you tried, and what has helped you in the past?	Are you having any other symptoms? Are you having any fevers? Are you having any shortness of breath? Have you started any new medications or supplements? Have you recently changed dosing or timing of medications you take? Have you tried any medications that have helped? Please give us a call to schedule an appointment. Give our triage nurses a call if your symptoms are worsening.
I've been dealing with itchy eyes for weeks now, and I'm guessing it's just my allergies acting up again. I was wondering if I could get your thoughts on it - should I just stick with my usual meds or is there something else I can try?	I'm sorry that you have been experiencing these troubling symptoms. Have you been having any other symptoms? Have you had any recent changes in your medications? Have you tried anything that may have helped alleviate your symptoms? If your symptoms are persisting on your usual allergy medications, or symptoms are worsening, please call the office to schedule an appointment.	Thanks for checking in. Are you having any other symptoms? How long have these symptoms been going on? Have you ever had symptoms like this before? Have you started any new medications or supplements? Have you recently changed dosing or timing of medications you take? Have you tried any medications that have helped? Please call to schedule an appointment. You should be seen in the ED if you have worsening or sudden shortness of breath, vision changes, or chest pain.

Table 12: Examples of different clinician responses to the same patient message within the SyPPM dataset. We collect responses from three separate annotators to 40 messages within the SyPPM dataset, and show selected examples from two annotators here.



		Content-Level			Theme-Level		
Dataset	Model	Precision	Recall	Edit-F1	Precision	Recall	Edit-F1
SoCPPM	0-Shot	0.06±0.01	0.29±0.08	0.10±0.01	0.48±0.01	0.83±0.08	0.61±0.01
	Theme	0.06±0.00	0.32±0.11	0.09±0.01	0.44±0.00	<b>0.85±0.11</b>	0.58±0.01
	RAG	0.11±0.03	<b>0.33±0.18</b>	0.14±0.01	0.49±0.03	0.75±0.18	0.59±0.01
	SFT	<b>0.15±0.01</b>	0.18±0.00	<b>0.15±0.01</b>	<b>0.63±0.01</b>	0.62±0.00	<b>0.62±0.01</b>
	TADPOLE	0.12±0.01	0.19±0.01	0.14±0.01	0.51±0.01	0.69±0.01	0.59±0.01
IAP		0.26	0.25	0.24	0.61	0.63	0.62

Table 13: Edit-F1 scores for LLM adaptations on the SoCPPM patient message response drafting dataset. Each model adaptation is performed on three underlying LLMs, we report scores as average±standard deviation. We report content-level precision, recall, and edit-F1 (Section 3.1), as well as theme-level precision, recall, and edit-F1 (Section 3.2). We report content-level inter-annotator predictability (IAP), comparing LLM performance and expert human alignment.

## F Additional Results

### F.1 SocPPM Results

The SoCPPM dataset is created to evaluate LLMs in a practical setting, in which response drafts are compared with the clinician response which was sent via the secure portal in real time. In some ways this is a less-ideal form of the patient message response drafting task, because real-time clinician responses tend to contain a high degree of variation which is challenging to filter automatically. For example, real-time clinician responses frequently contain standardized responses (“dot phrases”) which offer commonly-repeated instructions, e.g. “please call the COVID-19 hotline if you are experiencing any of the following symptoms...” (Baltaro et al., 2022). Additionally, real-time responses are written under more duress due to workforce constraints and growing use of the patient portal (Budd, 2023; Underdahl et al., 2024; Martinez et al., 2023; Yan et al., 2021).

In Table 13 we report the content-level and theme-level precision, recall and edit-F1 scores for adapted LLMs on the SoCPPM dataset. We find that LLMs in general perform more poorly on this dataset than the ideal IPPM and SyPPM datasets. The best-performing model adaptation on SyPPM (TADPOLE) achieves 0.20 content-level edit-F1 scores on SyPPM (see Table 3). The best-performing model adaptation on SoCPPM (SFT) achieves only 0.15 content-level edit-F1 on SoCPPM (Table 13). We hypothesize that this is because the SoCPPM dataset represents a version of the patient message response drafting task that is both more challenging, due to the existence of situated knowledge scattered throughout the EHR system that is unknowable for the response drafting

LLM, and less ideal, given that frequently clinician responses in practical settings can be messy and often sent under time pressure (Budd, 2023; Underdahl et al., 2024; Martinez et al., 2023; Yan et al., 2021).

We also note that SFT outperforms TADPOLE on the SoCPPM dataset, with SFT achieving 0.15 and 0.62 content- and theme-level edit-F1 scores, respectively, and TADPOLE achieving only 0.14 and 0.59 (Table 13). As TADPOLE adaptation uses thematic preference pairs to further fine-tune SFT models, we hypothesize that the themes used to generate these preference pairs are less suitable for the lower-quality, higher-variation responses found in real-time clinician responses.

### F.2 Class-Average Edit-F1 Scores

In Tables 3 and 13 we report content-level edit-F1 scores across the IPPM-SyPPM, and SoCPPM datasets, respectively. To investigate theme-specific performance of LLM response drafts, we also report theme class-specific scores at the content and theme levels. At the content level, in Table 14 we report the average recall of theme-labeled content within the expert responses of a given evaluation dataset. At the theme level, in Table 15 we report the class-average edit-F1 scores when predicting expert response themes with LLM response draft themes. We discuss these results in Section 5.

In Table 4 in Section 5 we give content-level and theme-level edit-F1 scores for the Claude 4.5 Sonnet (Anthropic, 2025), Gemini 2.5 Pro (Comanici et al., 2025) and GPT-OSS (Agarwal et al., 2025) reasoning models. In Tables 16 and 17 we similarly report the content-level average recall of theme-labeled content and the theme-level class-average edit-F1 scores.

Dataset	Model	Emp	SymQ	MedQ	Assess	Plan	Logis	Coord	Cont
SoCPPM	<i>Proportion</i>	<i>0.81</i>	<i>0.05</i>	<i>0.02</i>	<i>0.38</i>	<i>0.27</i>	<i>0.42</i>	<i>0.58</i>	<i>0.03</i>
	0-Shot	0.29±0.02	0.07±0.00	0.07±0.00	0.16±0.01	0.23±0.01	0.21±0.02	0.24±0.01	0.21±0.02
	Theme	0.30±0.02	0.07±0.01	0.07±0.00	<b>0.16±0.00</b>	0.23±0.00	0.23±0.03	<b>0.25±0.00</b>	0.24±0.04
	RAG	0.30±0.02	0.07±0.00	0.07±0.00	0.16±0.01	<b>0.25±0.03</b>	0.22±0.03	0.25±0.01	0.23±0.04
	SFT	<b>0.30±0.01</b>	<b>0.08±0.00</b>	0.07±0.00	0.15±0.00	0.21±0.01	<b>0.24±0.00</b>	0.24±0.00	<b>0.27±0.00</b>
	TADPOLE	0.30±0.02	<b>0.08±0.00</b>	0.07±0.00	0.15±0.01	0.23±0.03	0.23±0.02	0.24±0.00	0.25±0.04
IPPM	<i>Proportion</i>	<i>0.76</i>	<i>0.23</i>	<i>0.09</i>	<i>0.31</i>	<i>0.24</i>	<i>0.51</i>	<i>0.67</i>	<i>0.07</i>
	0-Shot	0.28±0.02	0.07±0.00	<b>0.07±0.00</b>	0.15±0.01	0.23±0.02	0.21±0.03	0.24±0.00	0.22±0.04
	Theme	0.29±0.02	0.07±0.02	0.06±0.01	<b>0.15±0.00</b>	0.28±0.09	0.23±0.02	0.25±0.01	0.22±0.07
	RAG	0.24±0.06	0.04±0.03	0.04±0.03	<b>0.15±0.00</b>	<b>0.33±0.09</b>	0.19±0.06	<b>0.26±0.01</b>	0.20±0.06
	SFT	<b>0.30±0.00</b>	<b>0.08±0.00</b>	<b>0.07±0.00</b>	<b>0.15±0.00</b>	0.21±0.00	<b>0.24±0.00</b>	0.24±0.00	<b>0.27±0.00</b>
	TADPOLE	<b>0.30±0.00</b>	<b>0.08±0.00</b>	<b>0.07±0.00</b>	0.15±0.01	0.21±0.02	0.24±0.01	0.24±0.01	0.26±0.02
SyPPM	<i>Proportion</i>	<i>0.99</i>	<i>0.79</i>	<i>0.79</i>	<i>0.33</i>	<i>0.05</i>	<i>0.75</i>	<i>0.11</i>	<i>0.56</i>
	0-Shot	0.28±0.03	0.06±0.02	0.07±0.01	0.15±0.00	<b>0.27±0.05</b>	0.22±0.01	0.24±0.00	0.22±0.02
	Theme	0.30±0.02	<b>0.08±0.00</b>	<b>0.08±0.00</b>	<b>0.16±0.01</b>	0.24±0.01	0.24±0.03	<b>0.25±0.01</b>	0.25±0.05
	RAG	<b>0.31±0.01</b>	<b>0.08±0.00</b>	0.07±0.00	<b>0.16±0.01</b>	0.23±0.02	<b>0.25±0.01</b>	0.24±0.01	0.26±0.01
	SFT	0.29±0.03	0.06±0.02	0.06±0.02	0.15±0.01	0.27±0.06	0.24±0.01	<b>0.25±0.01</b>	0.22±0.05
	TADPOLE	0.30±0.00	<b>0.08±0.00</b>	0.07±0.00	0.15±0.01	0.21±0.01	0.24±0.01	0.24±0.01	<b>0.27±0.00</b>
<i>Gemini</i>		<i>0.30</i>	<i>0.07</i>	<i>0.07</i>	<i>0.16</i>	<i>0.23</i>	<i>0.24</i>	<i>0.24</i>	<i>0.24</i>
<i>IAP</i>		<i>0.30</i>	<i>0.07</i>	<i>0.07</i>	<i>0.16</i>	<i>0.23</i>	<i>0.24</i>	<i>0.24</i>	<i>0.24</i>

Table 14: Class average content-level recall scores for adapted LLMs. Each model adaptation is performed on three underlying LLMs, we report average results  $\pm$  standard deviation. We report micro average recall scores for each theme class. We also report the proportion of responses which contain each theme in each dataset. We include SyPPM results of the best commercial model (Gemini with theme prompting) for comparison. Finally, we report theme-level IAP, comparing LLM performance and expert human alignment at the theme level.

## G Example REDCap Survey

In the IPPM evaluation dataset, ground-truth responses are written by a paid team of 4 expert primary care nurses who work daily in the patient portal, collected via REDCap surveys (Harris et al., 2009). In the SyPPM evaluation dataset, ground-truth responses are written by a paid primary care doctor with 15+ years of experience. Each clinician was paid \$50 for every 10 responses (estimated to take 1 hour), in order to give ample time to write a full response to each message/EHR summary. While writing responses, experts were prompted “if you had unlimited time, what would be included in your response to this patient?” To provoke quality responses, clinicians were given a separate text entry box for each of the themes derived in Section 2.1. For example, the *Treatment Contingency Planning* text box included the prompt “please outline a backup/red flag plan for the patient, if applicable.” Screenshots of an example REDCap survey question can be found in Figure 4 and Figure 5.

## H Prompts

In Section 3 we describe several methods for adapting LLMs for the patient message response drafting task. We give the 0-shot and thematic prompts in Figure 6 and Figure 7, respectively. The thematic prompt guides the model to use our derived themes when drafting responses to patient messages. In

Section 5 we see that thematic prompting, and other forms of contextual adaptation such as RAG, SFT, and our novel TADPOLE DPO-based strategy, improve LLM performance on the response drafting task.

## Message Labeling

AAA  
+ -

### Patient Message 1

Please refer to the following patient's message and chart summary while drafting a response to the patient.

#### Patient Chart Summary:

###Demographics###

...

###Full Active Problem List###:

...

###Recent Encounters (Max 10)###

...

###Medications (Outpatient)###

...

#### Patient Message:

Hey Doc,

I've been having some stomach issues lately. Sometimes I get these sharp pains and I'm really gassy. It's pretty annoying and I don't know what's causing it. It doesn't happen all the time, just every now and then.

Should I come in or something?

#### Empathetic Communication

*Empathetic response elements can include acknowledging the patient's problem, encouraging their positive effort, requesting patiences, or thanking the patient.*

Please write any empathetic statements you wish to include in your response to this patient, if applicable.

Expand

Figure 4: Screenshot of the beginning of a REDCap survey question used to collect clinician responses to patient messages in the SyPPM dataset. The patient's EHR chart and message are first given, then the clinician is prompted with a series of text entry boxes for each response theme described in Section 2.1.

<p><b>Logistics: Scheduling, Billing, Operations</b></p> <p>Please provide necessary scheduling, billing, or other operations information to the patient, if applicable. Give any logistical instructions to the patient. Ask the patient any logistics-related follow-up questions.</p>	<div></div> <p>Expand</p>
<p><b>Planning: Care Coordination and Communication Planning</b></p> <p>Please outline communication plans for/with the patient, and/or inter- or intra-clinic communication plans, if applicable.</p>	<div></div> <p>Expand</p>
<p><b>Planning: Treatment Contingency Planning</b></p> <p>Please outline a backup/red flag plan for the patient, if applicable.</p>	<div></div> <p>Expand</p>
<p><b>Other</b></p> <p>Please add anything else that should be said to this patient, but doesn't fit in the previous categories, if applicable.</p>	<div></div> <p>Expand</p>
<p><b>Assumptions/Notes/Additional Thoughts</b></p> <p><i>Please feel free to make any necessary assumptions about the patient's case based on their available chart or your best guess/experience with similar patients.</i></p> <p>Here you may detail any assumptions you needed to make when drafting your response, and any notes you wish to leave as part of your annotation. You may also detail parts of your thought process as you drafted your response.</p>	<div></div> <p>Expand</p>

Figure 5: Screenshot of the end of a REDCap survey response used to collect clinician responses to patient messages in the SyPPM dataset. After seeing the patient's EHR chart and message, the clinician is prompted with a series of text entry boxes for each response theme described in Section 2.1. The clinician is also prompted to give any additional thoughts or assumptions they made while drafting their response.



Dataset	Model	Emp	SymQ	MedQ	Assess	Plan	Logis	Coord	Cont
SoCPPM	<i>Proportion</i>	<i>0.81</i>	<i>0.05</i>	<i>0.02</i>	<i>0.38</i>	<i>0.27</i>	<i>0.42</i>	<i>0.58</i>	<i>0.03</i>
	0-Shot	0.88±0.02	0.12±0.06	0.08±0.07	<b>0.57±0.02</b>	0.46±0.02	<b>0.58±0.03</b>	0.68±0.02	<b>0.12±0.07</b>
	Theme	0.88±0.02	0.13±0.05	<b>0.12±0.05</b>	0.55±0.00	0.44±0.02	<b>0.58±0.03</b>	0.69±0.02	0.09±0.01
	RAG	0.82±0.05	0.12±0.05	0.07±0.06	0.55±0.01	<b>0.47±0.01</b>	0.57±0.03	<b>0.70±0.03</b>	0.11±0.09
	SFT	0.88±0.01	0.10±0.16	0.00±0.00	0.42±0.02	0.36±0.07	0.51±0.04	0.64±0.02	0.09±0.08
	TADPOLE	<b>0.89±0.00</b>	<b>0.18±0.04</b>	0.10±0.04	0.45±0.02	0.35±0.07	0.53±0.05	0.68±0.01	0.09±0.03
IPPM	<i>Proportion</i>	<i>0.76</i>	<i>0.23</i>	<i>0.09</i>	<i>0.31</i>	<i>0.24</i>	<i>0.51</i>	<i>0.67</i>	<i>0.07</i>
	0-Shot	0.85±0.02	0.05±0.05	0.09±0.06	<b>0.51±0.02</b>	0.42±0.02	<b>0.59±0.02</b>	0.76±0.03	0.12±0.05
	Theme	0.85±0.02	0.45±0.05	0.20±0.06	0.49±0.01	0.41±0.00	<b>0.59±0.02</b>	<b>0.77±0.00</b>	0.15±0.04
	RAG	0.77±0.05	0.05±0.02	0.08±0.10	0.47±0.00	<b>0.46±0.05</b>	0.58±0.03	0.75±0.02	0.11±0.03
	SFT	0.86±0.00	0.08±0.02	0.08±0.08	0.32±0.03	0.38±0.05	0.49±0.03	0.73±0.01	0.07±0.07
	TADPOLE	<b>0.87±0.00</b>	<b>0.45±0.02</b>	<b>0.21±0.07</b>	0.30±0.05	0.36±0.05	0.46±0.05	0.77±0.01	<b>0.15±0.02</b>
SyPPM	<i>Proportion</i>	<i>0.99</i>	<i>0.79</i>	<i>0.79</i>	<i>0.33</i>	<i>0.05</i>	<i>0.75</i>	<i>0.11</i>	<i>0.56</i>
	0-Shot	0.98±0.02	0.17±0.10	0.08±0.07	<b>0.50±0.00</b>	0.10±0.01	0.54±0.23	0.19±0.08	0.42±0.06
	Theme	<b>0.99±0.00</b>	0.72±0.01	0.32±0.01	0.50±0.01	0.06±0.04	0.52±0.05	0.17±0.01	0.33±0.15
	RAG	0.93±0.03	0.19±0.10	0.05±0.00	0.49±0.01	0.12±0.04	0.57±0.17	0.20±0.02	0.32±0.21
	SFT	0.98±0.01	0.38±0.04	0.09±0.04	0.33±0.08	<b>0.24±0.13</b>	<b>0.58±0.09</b>	<b>0.22±0.02</b>	0.13±0.03
	TADPOLE	<b>0.99±0.00</b>	<b>0.79±0.03</b>	<b>0.49±0.01</b>	0.16±0.04	0.17±0.09	0.19±0.02	0.22±0.03	<b>0.46±0.16</b>
	<i>Gemini</i>	<i>0.99</i>	<i>0.71</i>	<i>0.28</i>	<i>0.50</i>	<i>0.14</i>	<i>0.76</i>	<i>0.33</i>	<i>0.71</i>
	<i>IAP</i>	<i>0.80</i>	<i>0.80</i>	<i>0.53</i>	<i>0.38</i>	<i>0.07</i>	<i>0.73</i>	<i>0.15</i>	<i>0.06</i>

Table 15: Class average theme-level edit-F1 scores for LLM adaptations. Each model adaptation is performed on three underlying LLMs, we report average results  $\pm$  standard deviation. We report micro average edit-F1 scores for each theme class. We also report the proportion of responses which contain each theme in each dataset. We include SyPPM results of the best commercial model (Gemini with theme prompting) for comparison. Finally, we report theme-level IAP, comparing LLM performance and expert human alignment at the theme level.

Prompt	Model	Emp	SymQ	MedQ	Assess	Plan	Logis	Coord	Cont
0-Shot	GPT	<b>0.31</b>	0.07	<b>0.07</b>	0.15	0.21	0.24	0.24	0.27
	Gemini	0.30	<b>0.08</b>	<b>0.07</b>	0.15	0.21	<b>0.25</b>	0.24	<b>0.28</b>
	Claude	<b>0.31</b>	0.07	<b>0.07</b>	0.15	<b>0.23</b>	0.24	0.24	0.27
	Avg	<i>0.31</i>	<i>0.07</i>	<i>0.07</i>	<i>0.15</i>	<i>0.22</i>	<i>0.24</i>	<i>0.24</i>	<i>0.27</i>
Theme	GPT	0.28	0.07	<b>0.08</b>	0.14	<b>0.34</b>	0.21	0.24	0.19
	Gemini	0.30	0.07	0.07	<b>0.16</b>	0.23	0.24	0.24	0.24
	Claude	<b>0.31</b>	<b>0.08</b>	0.07	0.14	0.20	0.24	0.24	<b>0.27</b>
	Avg	<i>0.30</i>	<i>0.07</i>	<i>0.07</i>	<i>0.15</i>	<i>0.26</i>	<i>0.23</i>	<i>0.24</i>	<i>0.23</i>
<i>IAP</i>		<i>0.30</i>	<i>0.21</i>	<i>0.21</i>	<i>0.13</i>	<i>0.27</i>	<i>0.37</i>	<i>0.15</i>	<i>0.64</i>

Table 16: Class average content-level recall scores for Claude 4.5 Sonnet, Gemini 2.5 Pro and GPT-OSS reasoning models on the publicly-available SyPPM evaluation dataset. We evaluate each model using 0-shot and thematic prompts. Classifying elements in clinician responses into themes, we report response draft recall scores averaged across each theme. We also report content-level IAP, comparing LLM performance and expert human alignment at the content level.

Prompt	Model	Emp	SymQ	MedQ	Assess	Plan	Logis	Coord	Cont
0-Shot	GPT	<b>0.99</b>	0.42	<b>0.22</b>	<b>0.50</b>	0.10	<b>0.81</b>	0.27	0.64
	Gemini	<b>0.99</b>	0.16	0.03	<b>0.50</b>	<b>0.12</b>	0.79	0.29	<b>0.66</b>
	Claude	<b>0.99</b>	<b>0.49</b>	<b>0.22</b>	<b>0.50</b>	0.09	0.53	<b>0.34</b>	0.52
	Avg	<i>0.99</i>	<i>0.36</i>	<i>0.16</i>	<i>0.50</i>	<i>0.10</i>	<i>0.71</i>	<i>0.30</i>	<i>0.61</i>
Theme Prompting	GPT	<b>0.99</b>	0.80	<b>0.43</b>	<b>0.50</b>	0.10	<b>0.82</b>	0.27	0.60
	Gemini	<b>0.99</b>	0.71	0.28	<b>0.50</b>	<b>0.14</b>	0.76	<b>0.33</b>	<b>0.71</b>
	Claude	<b>0.99</b>	<b>0.87</b>	0.39	<b>0.50</b>	0.12	0.63	0.15	0.56
	Avg	<i>0.99</i>	<i>0.79</i>	<i>0.37</i>	<i>0.50</i>	<i>0.12</i>	<i>0.74</i>	<i>0.25</i>	<i>0.62</i>
<i>Theme Proportion</i>		<i>0.99</i>	<i>0.79</i>	<i>0.79</i>	<i>0.33</i>	<i>0.05</i>	<i>0.75</i>	<i>0.11</i>	<i>0.56</i>
<i>IAP</i>		<i>0.80</i>	<i>0.80</i>	<i>0.53</i>	<i>0.38</i>	<i>0.07</i>	<i>0.73</i>	<i>0.15</i>	<i>0.06</i>

Table 17: Class average theme-level edit-F1 scores for Claude 4.5 Sonnet, Gemini 2.5 Pro and GPT-OSS reasoning models on the publicly-available SyPPM evaluation dataset. We evaluate each model using 0-shot and thematic prompts. We report edit-F1 scores for each theme class. Additionally, we report the proportion of responses which contain each theme (theme proportion) in the SyPPM dataset. Finally, we also report theme-level IAP, comparing LLM performance and expert human alignment at the theme level.

You will be given a patient portal message and the patient's EHR chart. Respond to the message as though you were their doctor.  
Patient Chart: {chart}  
Patient Message: {message}  
Doctor Response:

Figure 6: 0-shot prompt for patient message response drafting

You will be given a patient portal message and the patient's EHR chart. Respond to the message as though you were their doctor.  
Good doctor responses are typically concise, and may contain some of the following themes:

- Empathetic communication, e.g. showing understanding and care for the patient's situation
- Symptom-related follow-up questions, e.g. asking about the patient's symptoms or history
- Medication-related follow-up questions, e.g. asking about the patient's medication history or plans
- Medical assessment, e.g. explaining symptom causes or why you are prescribing a medication
- Medical instruction, e.g. outlining medical plans for the patient or providing instructions to the patient
- Logistical information, e.g. giving details about scheduling, billing, operations
- Care coordination, e.g. outlining communication plans for/with the patient, and/or inter- or intra-clinic communication plans
- Contingency planning, e.g. outlining a backup/red flag plan for the patient

Patient Chart: {chart}  
Patient Message: {message}  
Doctor Response:

Figure 7: Thematic prompt for patient message response drafting