

# T-FIX: Text-Based Explanations with Features Interpretable to eXperts

Shreya Havaladar<sup>§\*</sup> Helen Jin<sup>§\*</sup> Chaehyeon Kim<sup>§\*</sup> Anton Xue<sup>§\*</sup> Weiqiu You<sup>§\*</sup>

Gary Weissman<sup>†</sup> Rajat Deo<sup>†</sup> Sameed Khatana<sup>†</sup>

Helen Qu<sup>★</sup> Marco Gatti<sup>★</sup> Daniel A. Hashimoto<sup>†</sup> Amin Madani<sup>‡</sup>

Masao Sako<sup>★</sup> Bhuvnesh Jain<sup>★</sup> Lyle Ungar<sup>§</sup> Eric Wong<sup>§</sup>

{shreyah, helenjin, chaenyk, antonxue, weiqiuy, ungar, exwong}@seas.upenn.edu

## Abstract

As LLMs are deployed in knowledge-intensive settings (e.g., surgery, astronomy, therapy), users expect not just answers, but also meaningful explanations for those answers. In these settings, users are often domain experts (e.g., doctors, astrophysicists, psychologists) who require confidence that a model’s explanation reflects expert-level reasoning. However, current evaluation schemes primarily emphasize plausibility or internal faithfulness of the explanation, often neglecting whether the content of the explanation truly aligns with expert intuition. We formalize *expert alignment* as a criterion for evaluating explanations with T-FIX, a benchmark spanning seven knowledge-intensive domains. T-FIX includes datasets and novel alignment metrics developed in collaboration with domain experts, so an LLM’s explanations can be scored directly against expert judgment.<sup>1</sup>

## 1 Introduction

LLMs are increasingly used for domain-specific tasks, which require substantial background knowledge from specialized fields. It is foreseeable that LLM-powered systems will soon assist in high-stakes environments such as operating rooms, astronomical observatories, and therapeutic settings.

For LLMs to be trustworthy and reliable in these critical applications, users require not only correct answers but also **good explanations** (Rudin, 2019; Pedreschi et al., 2019).

What constitutes a “good explanation”? This largely depends on *the explanation’s target audience* (Ribera and Lapedriza, 2019; Sokol and Flach, 2020). As LLMs are increasingly adopted for specialized tasks like surgical assistance or supernova analysis, the primary users are often domain experts, such as doctors and astrophysicists. Consequently, a “good explanation” in these specialized

**QUESTION:** Analyze this patient’s electronic health record and determine if they are at risk for developing sepsis in the next 12 hours. <PATIENT DATA>

**PREDICTION:** High risk

**EXPLANATION:** The patient is at high risk because they are 37 years old, have elevated glucose levels, are hypotensive, and seem visibly distressed.

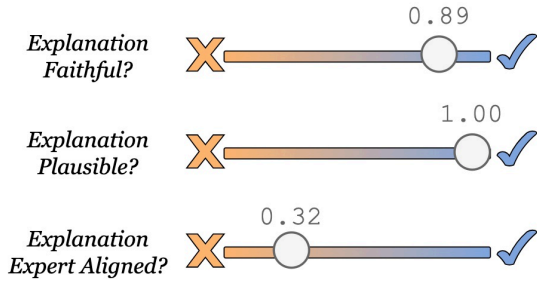


Figure 1: Most current evaluations for LLM explanations consider two dimensions: the overall plausibility and the faithfulness to the reasoning process. However, a crucial third dimension, **expert alignment**, asks: Does the LLM reason like a domain expert would? For example, an LLM correctly predicts sepsis risk with a plausible, faithful explanation, but because the explanation emphasizes features that clinicians rarely use for sepsis diagnosis, the expert alignment score is low.

contexts must *offer insights that are valuable and interpretable to these domain experts*.

Existing evaluations of LLM explanations predominantly focus on two dimensions: (1) plausibility, ensuring that the answer logically follows from the provided explanation; and (2) faithfulness, verifying that the answer accurately reflects the LLM’s actual reasoning process. (Zhou et al., 2021; Agarwal et al., 2024; Parcalabescu and Frank, 2023).

While these dimensions are necessary, they are not sufficient for knowledge-intensive applications. Domain experts often need highly specific informa-

<sup>1</sup><https://anonymous.4open.science/r/FIX-2-BE33/>

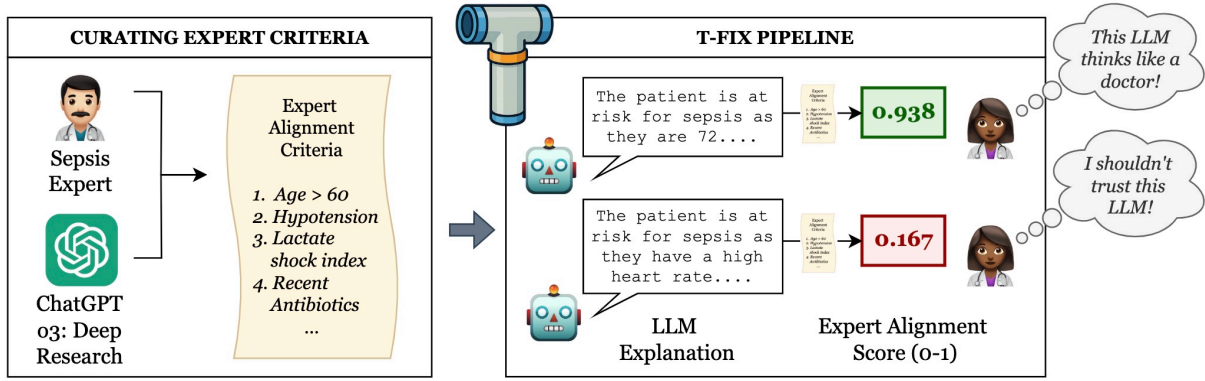


Figure 2: An overview of the T-FIX construction process. For each dataset, we first establish expert alignment criteria – features deemed important by domain experts for a specific task – through collaboration with these experts and LLM-based deep research tools. These criteria form the basis of the T-FIX evaluation pipeline, which processes an LLM-generated explanation to output an expert alignment score. A high score suggests the explanation reflects reasoning aligned with domain experts (i.e., the LLM “thinks like an expert”), while a low score indicates the explanation may rely on aspects that experts would deem irrelevant.

tion regarding how a prediction was derived (Wang and Yin, 2021), particularly whether **the LLM considered aspects of the input that they themselves deem critical**.

To address this, we propose a third dimension for evaluating LLM-generated explanations: **Expert Alignment**. This dimension measures the extent to which an LLM-generated explanation for a given input and prediction focuses on criteria that a domain expert would deem important when making the same prediction.

An LLM can generate a correct answer with a plausible and faithful explanation, yet still rely on features that domain experts consider irrelevant or low-priority, as shown in Figure 1. Such misaligned reasoning can undermine trust in the model, even when the output is technically correct.

While alignment with domain expert reasoning has been explored in machine learning, for example, by identifying meaningful feature groups (Jin et al., 2024), such approaches are primarily suited for interpreting traditional, non-generative neural networks. Modern LLMs typically generate free-form text explanations that are not directly based on these explicit feature groups. To our knowledge, no benchmark currently exists to evaluate the expert alignment of such free-form textual explanations.

To fill this gap, we introduce the T-FIX benchmark: a collection of datasets spanning seven distinct domains, accompanied by an evaluation framework. Designed in collaboration with domain experts, T-FIX assesses the expert alignment of LLM-generated explanations within each domain.

Our contributions are as follows:

- We introduce *expert alignment as a desired attribute of LLM-generated explanations* and create T-FIX, the first benchmark designed to evaluate this.
- We release a pipeline to *evaluate how well any LLM “thinks like an expert,”* designed to be easily extendable to new domains.
- We demonstrate that current LLMs often *struggle to generate explanations that align with expert intuition*, highlighting this as a significant area for their future improvement.
- We find that LLMs generally perform better when they reason over multiple expert criteria, yet modern high-performing LLMs *do not appear to rely on expert reasoning*.

## 2 Expert Alignment Criteria

The development of the T-FIX benchmark was a highly collaborative and interdisciplinary process. For each of our seven domains (see Figure 4), our first step was to identify the **expert criteria most relevant to making a prediction**, detailed in the left of Figure 2.

When answering knowledge-intensive questions like “Will this patient develop sepsis in the next 12 hours?” or “What kind of supernova produced these wavelengths?”, doctors and astrophysicists rely on domain-specific heuristics, prioritizing certain features over others based on training and experience. For instance, in sepsis classification, an experienced clinician would typically emphasize

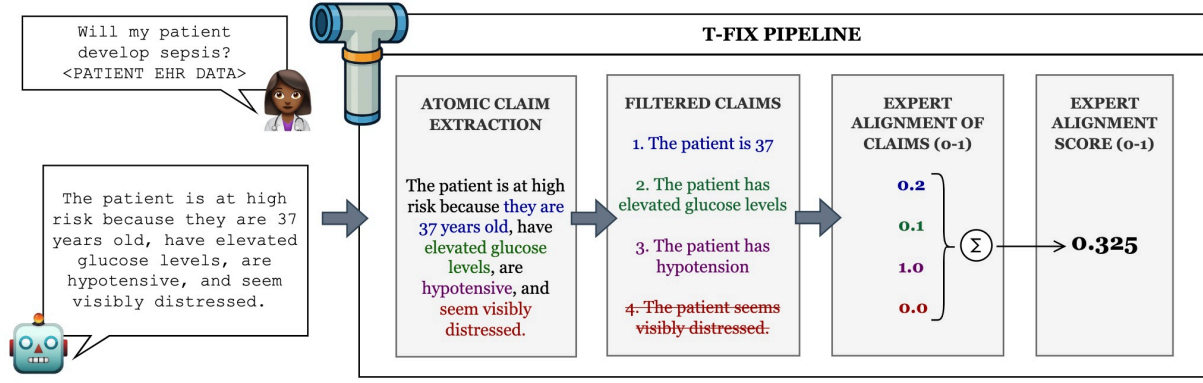


Figure 3: Our T-FIX pipeline. To evaluate an LLM-generated explanation, we first decompose it into atomic claims. Next, we filter out irrelevant claims, such as unsupported or speculative statements. Each remaining claim is then scored against the domain-specific expert alignment criteria on a 0–1 scale: a score of 1 indicates perfect overlap with at least one criterion, while 0 indicates no overlap. Filtered-out claims are automatically assigned a score of 0. We compute the final expert-alignment score for the explanation by averaging across all claim scores.

features like advanced age and hypotension, while assigning lower importance to signals like glucose levels or patient demeanor, which are less directly indicative of sepsis risk.

Thus, an LLM that makes the correct prediction by attending to age and hypotension is *more expert-aligned* than one that arrives at the same answer by focusing on glucose and demeanor.

We define the subset of features that experts prioritize most highly when performing a task as the task’s **expert alignment criteria**.

**Step 1: Surveying the Field.** To seed our initial list of expert criteria, we prompt OpenAI’s o3 model to perform a comprehensive literature review of the relevant field. Each prompt includes a task description, example input-output pairs from the dataset, and instructions to generate a list of criteria considered important for performing the task – accompanied by reputable citations.

We begin with this deep research approach to *avoid over-reliance on any single expert’s perspective*. Our goal is to synthesize insights from a broad array of books, journals, and academic publications to produce as comprehensive a list as possible.

**Step 2: Iteration with Domain Experts.** To validate and improve the output from Step 1, we present the preliminary criteria list to a domain expert (see Figure 4 for details on each expert per domain). We ask the expert to (1) remove any incorrect or irrelevant criteria, (2) add any important ones that were missed, and (3) ensure that the list reflects a consensus that their peers would agree with. The expert then refines the list until it accu-

ately captures the field’s knowledge.

An example criterion for sepsis classification is as follows: Advanced age (over 65 years) markedly increases susceptibility to rapid sepsis progression and higher mortality after infection.

All Deep Research prompt templates and final expert alignment criteria lists for all domains are available in our GitHub repository.

### 3 T-FIX Pipeline

LLM-generated explanations contain a mix of reasoning steps – some aligned with expert judgment, and others based on irrelevant information.

To systematically evaluate such complex explanations, we first break them down into atomic claims, or standalone “features” that can be individually assessed for expert alignment. By scoring each feature separately and then aggregating these scores, we can compute an overall expert alignment score for the full explanation. See Figure 3 for an example of this multi-step process.

Our T-FIX pipeline for evaluating expert alignment consists of three main components:

- Claim Extraction:** Decomposing a free-form explanation into standalone, atomic claims.
- Relevancy Filtering:** Removing claims that are unsupported, speculative, or otherwise irrelevant to the model’s prediction.
- Alignment Scoring:** Measuring the degree of overlap between each remaining claim and domain expert criteria on a 0–1 scale.

We build our pipeline using GPT-4o, as it is both fast and cost-effective.



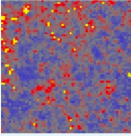
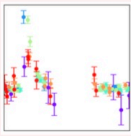
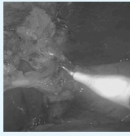


DOMAIN	Cosmology		Psychology		Medical		
DATASET	Mass Maps	Supernova	Politeness	Emotion	Cholecystectomy	Cardiac	Sepsis
ADAPTED FROM	[Kacprzak et al., 2023]	[Team et al., 2018]	[Havaladar et al., 2023a]	[Demszyk et al., 2020]	[Madani et al., 2022]	[Kansal et al., 2025]	[Kansal et al., 2025]
MOTIVATION	Discovering relationships between cosmological structures and the initial state of the universe	Identifying time periods with high astronomical signal to optimize telescope observations	Understanding differences in politeness expression to improve cross-cultural communication.	Understanding the nuances of emotion expression in online settings.	Helping surgeons identify which incisions optimize patient safety while operating	Helping clinicians identify patents who are risk of cardiac arrest during ER admission	Helping clinicians identify which variables contribute to sepsis development
TASK	Predicting cosmological parameters $\Omega_m$ and $\sigma_8$ given an image representing weak lensing maps data.	Classifying the type of astronomical object (SNIa, TDE, etc.) given time-series flux measurements across multiple wavelengths	Classifying the politeness of a text conversation snippet in English, Japanese, Chinese, or Spanish.	Detecting which of 28 core emotions is most reflected by the speaker of a text Reddit comment.	Determining safe/unsafe organ regions to cut into during cholecystomy surgery given a laparoscopic image of a patient's abdomen.	Determining whether a patient is at high risk of soon experiencing cardiac arrest given time-series Electrocardiogram (ECG) data.	Determining whether a patient is at high risk of developing sepsis in the near future given time-series Electronic Health Record (EHR) data.
INPUT → OUTPUT	Weak lensing map image → $\Omega_m$ , $\sigma_8$ values	Multiband time series data → astronomical object class	Conversation snippet → politeness level on a 1-5 scale	Reddit comment → emotion label	Image from laparoscopic camera → description of safe and unsafe regions	ECG time series data → Yes/No cardiac arrest classification	EHR time series data → Yes/No sepsis risk prediction
INPUT EXAMPLE			"I totally didn't realize this was a vandalized page. Please accept my apology"	"Thanks for your reply:) until then hubby and I will anxiously wait "			
DOMAIN EXPERT	Astronomy professor at an American university	Astrophysics professor at an American university	Psychology professor at an American university	Psychology professor at an American university	Gastrointestinal surgeon in an American hospital	Professor of cardiovascular medicine at an American university	Pulmonary care physician at an American hospital
EXPERT ALIGNMENT CRITERIA	A set of cosmological lensing features such as cluster peaks, voids, filaments, clumpiness, connectivity, and contrast — used to infer parameters through matter distribution patterns.	A classification framework for astrophysical transients based on flux continuity, light curve shape, amplitude, duration, periodicity, spectral features, and photometric evolution trends.	A taxonomy of politeness strategies including honorifics, apologies, indirectness, and discourse cues across social, emotional, and linguistic contexts.	A taxonomy of emotional cues from valence, arousal, and direct emotion markers to signals of confusion, blame, praise, and relief — used to infer nuanced affective states.	A checklist of expert surgical safety criteria for cholecystectomy, emphasizing precise anatomical identification, dissection landmarks, and caution in high-risk variations.	A set of ECG indicators including HR deceleration, ST changes, QRS abnormalities, atrial arrhythmias, and conduction delays — signaling imminent arrest risk.	A sepsis risk framework combining age, vital sign criteria (SIRS, qSOFA, NEWS), lactate, shock index, hypotension, SOFA changes, and early clinical actions to flag severity.

Figure 4: Overview of datasets and domains in T-FIX. We evaluate LLM expert alignment across seven diverse domains, spanning cosmology, psychology, and medicine. For each dataset, we highlight the motivating task, input–output format, representative example, and the expert responsible for validating alignment criteria. The final row summarizes the expert alignment criteria used for scoring explanations in each domain. The column colors reflect dataset modality: blue indicates vision, yellow indicates language, and pink indicates time-series.

### 3.1 Stage 1: Atomic Claim Extraction

Given a free-form text explanation accompanying an LLM’s prediction, our first goal is to identify and extract the distinct reasoning steps, i.e. “features”, used by the LLM. We achieve this by decomposing the explanation into *atomic claims*.

An atomic claim is defined as a self-contained, indivisible statement that conveys a single veri-

fiable fact, and can be fully understood without reference to the surrounding context.

To extract atomic claims, we adapt prompting techniques from the claim decomposition literature (Wanner et al., 2024; Gunjal and Durrett, 2024) and prompt GPT-4o to transform a free-form explanation into a list of fully decontextualized atomic claims. We treat each claim as representing a single “feature” in the LLM’s explanation.

### 3.2 Stage 2: Relevancy Filtering

Not all extracted claims contribute meaningfully to expert reasoning. Some may be unsupported (i.e., references to content not present in the input), speculative (i.e., unfounded hypotheses), or otherwise irrelevant (e.g., repeating the model’s final prediction or citing unrelated information).

Given that domain experts heavily prefer succinct, informative explanations, we prompt GPT-4o to remove such noisy claims by evaluating each atomic claim based on the original input. A claim is retained if it satisfies the following two criteria:

- (1) Clearly grounded in and supported by the input (i.e., not unfounded or speculative)
- (2) Directly contributes to explaining *why* the model made its prediction.

On average, 72% of the claims generated in Stage 1 pass this relevancy filter and are carried forward for alignment scoring.

### 3.3 Stage 3: Alignment Scoring

In the final stage of our pipeline, we evaluate each retained atomic claim by comparing it to the domain-specific expert alignment criteria (see Section 2). This step quantifies how closely the reasoning in the LLM’s explanation reflects expert judgment.

Given an atomic claim and a list of expert criteria, we prompt GPT-4o to measure the claim’s expert alignment in two steps:

1. **Identify the most aligned expert criterion.** The model selects the criterion whose focus and intent best match the core idea of the atomic claim. The model may also indicate that no criteria align with the claim.
2. **Assign an alignment score (0-1).** The model scores how well the claim aligns with the chosen criterion: 1 for complete overlap, and 0 for no alignment. Intermediate scores reflect partial alignment, such as when the claim touches on a relevant concept but lacks specificity. See Table 1 for details on intermediate scores.

For example, consider the expert criterion for sepsis classification: Advanced age (over 65 years). The claim “The patient is at risk as they are 72 years old” would receive an alignment score of 1.0, as it directly and fully supports the criterion. In contrast, the claim “The patient is at risk as they are 37” may receive a score of 0.2: while it discusses patient age, the specific value does not

Score Range	Meaning
(0, 0.25]	The claim references an unrelated or misleading feature, or misinterprets the criterion’s meaning
(0.25, 0.5]	The claim loosely refers to the correct concept but lacks key details, thresholds, or uses vague language
(0.5, 0.75]	The claim references a relevant feature but only partially reflects the criterion (e.g., omits thresholds, is overly general, contains noise)
(0.75, 1]	The claim is specific, directly relevant, and fully captures the meaning and intent of the expert criterion

Table 1: Interpretation of alignment score ranges used in scoring atomic claims against expert criteria.

align with the expert threshold for elevated risk. In contrast, the claim “The patient is NOT at risk as they are 37” would also receive a score of 1.0.

Examples of claims with high and low alignment for each domain, along with rationale for why those scores were assigned, are provided in Table A3.

### 3.4 Final Aggregation

We assign an alignment score of 0 to the claims that were filtered out or did not align with any criteria. This ensures *LLM-generated explanations are penalized for unsupported or speculative statements, irrelevant information, and misaligned reasoning*. We then average the alignment scores across all claims to produce a final expert alignment score for the explanation.

The prompts for all three stages can be found in Appendix B and in our Github repository.

## 4 Pipeline Validation

Given our pipeline relies on multiple curated GPT-4o prompts, we want to ensure that the extracted and filtered claims are accurate, and that the final alignment scores match domain expert intuition.

To validate the outputs at each stage, we conduct an annotation study for 35 examples (5 per domain). This includes 295 extracted claims and 211 aligned claims. We recruit a total of six annotators, with two annotators per example<sup>2</sup>.

**Validating atomic claim extraction.** Annotators receive the original explanation and its extracted atomic claims from Stage 1. They classify each

<sup>2</sup>Annotators are PhD students who study machine learning at an American university and are previously familiar with evaluating LLM outputs for given criteria.

Pipeline Stage	$\mathcal{N}$	Accuracy	Cohen’s $\kappa$
Claim Extraction	35	0.943	0.717
Relevancy Filtering	295	0.871	0.402
Expert Alignment	211	0.923	0.405

Table 2: Pipeline validation: Accuracy averaged across all T-FIX domains and annotator agreement – Cohen’s  $\kappa$  for each stage in our pipeline. Domain-specific statistics are provided in Table A2.

extraction as: (A) Perfect – all claims correctly extracted, (B) Partially accurate – 1–3 claims missing or incorrect, or (C) Incorrect – 3+ claims missing or incorrect. We convert these labels to accuracy scores:  $A = 1.0$ ,  $B = 0.5$ ,  $C = 0.0$ .

**Validating relevancy filtering.** Annotators review the explanation, extracted claims, and filtered claims from Stage 2. For each claim, they assess whether: (A) It was correctly kept or filtered, (B) It was incorrectly kept or filtered, or (C) It is ambiguous or borderline. These are scored as:  $A = 1.0$ ,  $B = 0.0$ ,  $C = 0.5$ .

**Validating expert alignment scoring.** Annotators are shown the alignment criteria and the filtered, scored claims from Stage 2. We define *direction* as the alignment score category (high, neutral, low), and *magnitude* as the exact score (e.g., 0.1 vs. 0.3 for low alignment).

Annotators evaluate each score as: (A) Fully accurate – an expert would agree with the score; correct direction and magnitude, (B) Partially accurate – correct direction, but magnitude off by  $\leq 0.2$ , or (C) Incorrect – wrong direction and magnitude off by  $> 0.2$ . These are scored as:  $A = 1.0$ ,  $B = 0.5$ ,  $C = 0.0$ .

**Results & agreement.** Table 2 reports average accuracy at each stage across all seven T-FIX domains, along with Cohen’s  $\kappa$  for inter-annotator agreement. The  $\kappa$  scores fall in the moderate-to-substantial agreement range, suggesting consistent annotator judgments and supporting the validity of our T-FIX pipeline. Domain-specific metrics are shown in Table A2.

## 5 Included Datasets

T-FIX contains seven open-source datasets, spanning the fields of cosmology, psychology, and medicine. To assess LLM explanations across multiple modalities, we include text, vision, and time-series datasets. We select these seven datasets due

to the availability of domain experts willing to work with us for these tasks.

As running T-FIX requires querying LLMs, many of which follow a pay-as-you-go API structure, we keep the total size of our benchmark to 700 (100 per dataset) in order for T-FIX to be accessible to as many researchers as possible.

We select a subset of 100 examples from the test set of each open-source dataset in T-FIX, and balance this sampling across classes when possible. We provide an overview of the included open-source datasets in Figure 4.

See Appendix C for additional details about the motivation, task, and prompting procedure for each dataset.

## 6 Experiments

After building a pipeline to evaluate the expert alignment of an LLM explanation, we evaluate a suite of today’s top LLMs on T-FIX to determine how expert-aligned these models are on domain-specific tasks.

We use the following prompting techniques as baselines to generate explanations for each dataset in T-FIX.

1. **Vanilla:** The LLM is prompted to generate an explanation along with its answer, without any additional guidance or reasoning structure.
2. **Chain-of-Thought (CoT):** The LLM is prompted to reason step-by-step through intermediate steps before answering, supporting more accurate responses on complex, multi-step tasks.
3. **Socratic Prompting:** The LLM is instructed to question its own reasoning, encouraging reflection and the surfacing of uncertainties or assumptions.
4. **Subquestion Decomposition:** The LLM is guided to break down a complex task into simpler subquestions, answer them individually, and then synthesize a final response.

Domain-specific prompts are detailed in Appendix C, with templates for the above prompting strategies in Figure A5.

Results for GPT-4o, GPT-o1, Gemini-2.0-Flash, and Claude-3.5-Sonnet<sup>3</sup> are shown in Table 3.

<sup>3</sup>We only select LLMs with vision support and context windows long enough to accommodate our time-series datasets. All models are accessed in May 2025.



Baseline	Cosmology		Psychology		Medicine		
	Mass Maps	Supernova	Politeness	Emotion	Cholecystectomy	Cardiac	Sepsis
<i>GPT-4o</i>							
Vanilla	0.421	0.877	0.629	0.597	0.295	0.533	0.545
Chain-of-Thought	0.390	0.859	0.625	0.639	0.338	0.564	0.532
Socratic Prompting	0.412	0.859	0.596	0.612	0.369	0.569	0.539
SubQ Decomposition	0.354	0.881	0.596	0.531	0.358	0.519	0.563
<i>o1</i>							
Vanilla	0.616	0.778	0.615	0.609	0.443	0.501	0.515
Chain-of-Thought	0.595	0.766	0.620	0.658	0.473	0.481	0.552
Socratic Prompting	0.503	0.782	0.555	0.467	0.456	0.449	0.578
SubQ Decomposition	0.491	0.805	0.536	0.545	0.409	0.473	0.576
<i>Gemini-2.0-Flash</i>							
Vanilla	0.515	0.811	0.618	0.600	0.407	0.529	0.566
Chain-of-Thought	0.507	0.815	0.569	0.566	0.376	0.553	0.578
Socratic Prompting	0.281	0.815	0.559	0.554	0.394	0.475	0.581
SubQ Decomposition	0.405	0.789	0.566	0.520	0.393	0.494	0.584
<i>Claude-3.5-Sonnet</i>							
Vanilla	0.710	0.761	0.634	0.642	0.264	0.565	0.611
Chain-of-Thought	0.688	0.776	0.639	0.622	0.286	0.538	0.584
Socratic Prompting	0.698	0.764	0.590	0.580	0.292	0.549	0.592
SubQ Decomposition	0.628	0.754	0.631	0.617	0.271	0.555	0.584

Table 3: Evaluating top LLMs on T-FIX. We report the average expert alignment score across all examples in the dataset. Corresponding accuracies are in Table A1 and baseline prompting strategies are described in Section 6.

## 7 Analysis

In this section, we analyze how LLMs distribute reasoning across expert criteria and whether higher task accuracy indicates better expert alignment.

### 7.1 Coverage of Expert Alignment Criteria

Section 3 describes our pipeline for measuring the proportion of expert-aligned claims in LLM explanations. We now examine a complementary question: *How many expert alignment criteria does an LLM consider across its explanations?*

A single gold-standard explanation rarely requires reasoning over *all* expert criteria; most high-quality explanations reference only 3–5. Thus, assessing coverage at the question level is not meaningful. Instead, we analyze coverage at the dataset level – whether different prompting strategies lead to a broader utilization of expert criteria across all questions within a domain.

Figure 5 presents the Shannon entropy of GPT-4o’s covered expert criteria in each domain. We observe a correlation between entropy and performance: domains where GPT-4o underperforms (e.g., Cholecystectomy, Supernova) show lower entropy, indicating limited criteria coverage. In contrast, well-performing domains (e.g., Politeness, Sepsis) exhibit more uniform coverage, equally tak-

ing into account all expert criteria.

This suggests that **LLMs that reason uniformly over expert alignment criteria perform better** – a promising insight for future work in prompting or training models to incorporate a broader range of expert reasoning.

### 7.2 Expert-Alignment vs. Accuracy

T-FIX focuses on evaluating explanation quality, but we are also interested in understanding the relationship between expert alignment and prediction accuracy. Specifically, we ask: *Does higher answer accuracy correspond to stronger expert alignment?*

Figure 6 shows the Pearson correlation of expert alignment (see Table A3) with accuracy (see Table A1) for each domain, averaged across models. In some domains with higher performance, like Cholecystectomy and Emotion, we do observe higher expert alignment as well. However, the overall correlation is weak across domains.

The heatmap suggests **today’s high-performing LLMs do not appear to rely on expert reasoning**. Future research is needed to explore whether aligning model reasoning with expert criteria – via training objectives or prompting – can improve downstream performance.

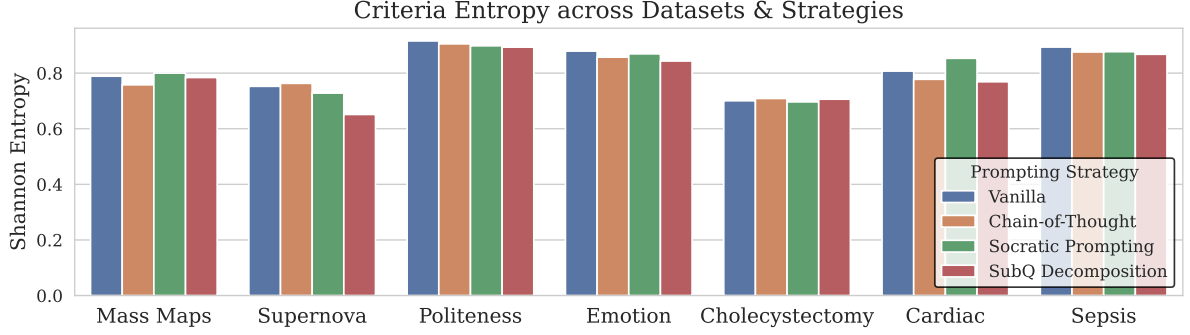


Figure 5: Shannon Entropy of expert alignment criteria for GPT-4o. For each prompting baseline, we show coverage of each domain’s explanations across all expert criteria – a high value indicates the LLM considers *many criteria across examples*, while a low value indicates the LLM *focuses on the same criteria repeatedly*.

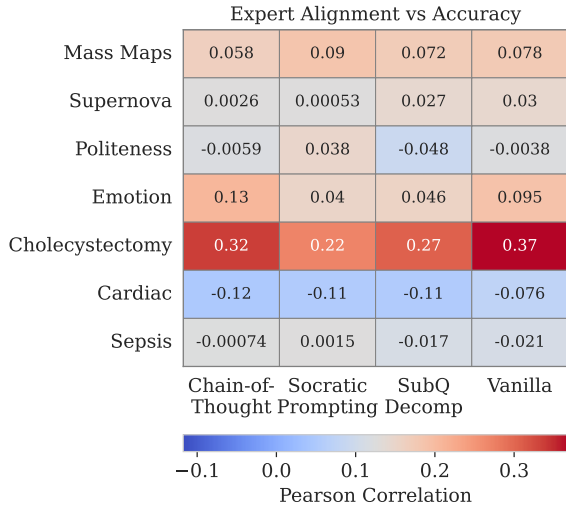


Figure 6: Expert-Alignment vs. Accuracy Correlation Heatmap, averaged across GPT-4o, o1, Gemini-2.0-Flash, and Claude-3.5-Sonnet. Red indicates positive correlation, blue is negative, gray is no correlation.

## 8 Related Work

**Evaluating LLM Explanations.** Common explanation methods for LLMs include feature attribution (e.g., LIME, SHAP (Ribeiro et al., 2016; Lundberg and Lee, 2017)), counterfactuals, and self-generated explanations (Im et al., 2023; Zhao et al., 2023). Some models are also trained to produce human-readable justifications (Camburu et al., 2018). To assess explanation quality and utility, recent work highlights criteria such as faithfulness (alignment with the model’s reasoning) and plausibility (how convincing it is to humans) (Jacovi and Goldberg, 2020; Zhou et al., 2021; Agarwal et al., 2024). Human studies show mixed outcomes: explanations sometimes aid understanding (Hase and Bansal, 2020; Bansal et al., 2021), but can also

offer little value or cause over-trust (Wang et al., 2023). A promising alternative is to use LLMs as automatic judges of explanation quality (Zheng et al., 2023; Chen et al., 2024), providing a scalable substitute for expensive human evaluation; we adopt this approach in T-FIX.

**Domain & Expert Alignment** Concept-based models constrain parts of the network to predict high-level, human-defined concepts, enabling incorporation of domain knowledge into final predictions (Koh et al., 2020). Extensions of concept bottlenecks and related methods aim to align latent representations with semantically meaningful features (Kim et al., 2018; Chen et al., 2020; Ghorbani et al., 2019), potentially grouped for expert interpretability (Jin et al., 2024). In NLP, integrating human knowledge has included collecting human-written explanation datasets to train models (Camburu et al., 2018) and using learned explanations to guide predictions (Bhatt et al., 2020). To our knowledge, no prior work explicitly evaluates text explanations for expert alignment like T-FIX.

## 9 Conclusion

We introduce T-FIX, the first benchmark designed to evaluate LLM explanations for expert alignment across seven knowledge-intensive domains. Our analysis reveals that today’s models struggle to generate explanations that experts would rely on, highlighting a critical area for improvement.

Future work may include exploring instruction-tuning LLMs to generate explanations with strong expert alignment, extending T-FIX to additional domains, and Human-Computer Interaction studies exploring how expert-aligned explanations affect real-world decision-making by practitioners.



## Limitations

As with any LLM-based system, the quality of the outputs is dependent on the input prompt. T-FIX is no exception – though we spend a significant amount of time analyzing outputs and prompt iterating, we do a finite amount of prompt iteration. There is a chance our benchmark could be marginally improved with additional prompt iteration. We hope the issue of prompt dependency diminishes with future models that are more robust and less susceptible to tiny prompt ablations.

While our evaluation pipeline currently uses GPT-4o for scoring, it is model-agnostic by design, and we encourage future work to apply or adapt the pipeline with other LLMs to improve robustness and reduce evaluator-model entanglement.

For pipeline validation, we conduct a user study where we annotate 35 examples. Though the annotation results on this subset suggest our pipeline is accurate, this work could have benefited from a larger and more robust annotation study. Future work should also involve domain experts vetting the pipeline in addition to recruited annotators.

In addition, we only have one expert to validate the expert alignment criteria for each domain. Though our usage of a deep research LLM minimizes over-reliance on a single domain expert, multiple experts would have been better to create the expert criteria. We were constrained by domain experts eager and available to collaborate with us.

Our experiments focus on a set of four models and four prompting strategies, and including additional models and strategies could provide a more comprehensive set of baseline results. Though many other high-performing LLMs and prompting techniques exist as of May 2025, we are conscious of budget and the environmental impact of running multiple experiments using T-FIX.

## Ethical Considerations

Using LLMs in the domains we describe in T-FIX, especially those relating to medicine, poses a unique set of risks and challenges. We do not advocate that LLMs should replace domain experts in these tasks; rather, T-FIX should serve as a step towards experts being able to use LLMs in a reliable and trustworthy way.

Additionally, LLMs are constantly changing, especially those that are company-owned and not open-source. This poses potential issues relating to the reproducibility of our baseline results as time

progresses and advances are made.

Lastly, nearly all LLMs contain biases – some harmful – that may propagate up in a system built off of these models. All users of T-FIX must be conscious of this risk.

## References

- T. M. C. Abbott, M. Agüena, A. Alarcon, S. Allam, O. Alves, A. Amon, F. Andrade-Oliveira, J. Annis, S. Avila, D. Bacon, E. Baxter, K. Bechtol, M. R. Becker, G. M. Bernstein, S. Bhargava, S. Birrer, J. Blazek, A. Brandao-Souza, S. L. Bridle, and 151 others. 2022. [Dark energy survey year 3 results: Cosmological constraints from galaxy clustering and weak lensing](#). *Physical Review D*, 105(2).
- Chirag Agarwal, Sree Harsha Tanneru, and Himabindu Lakkaraju. 2024. Faithfulness vs. plausibility: On the (un) reliability of explanations from large language models. *arXiv preprint arXiv:2402.04614*.
- Gagan Bansal, Tongshuang Wu, Joyce Zhou, Raymond Fok, Besmira Nushi, Ece Kamar, Marco Tulio Ribeiro, and Daniel S. Weld. 2021. [Does the whole exceed its parts? the effect of AI explanations on complementary team performance](#). In *Proceedings of the 2021 CHI Conference on Human Factors in Computing Systems (CHI)*, pages 1–16.
- Umang Bhatt, Alice Xiang, Shubham Sharma, Adrian Weller, Ankur Taly, Ying Jia, Joydeep Ghosh, Rajiv Puri, José M. F. Moura, and Peter Eckersley. 2020. Explainable machine learning in deployment. In *Proceedings of the 2020 Conference on Fairness, Accountability, and Transparency (FAT\*)*, pages 648–657.
- Oana-Maria Camburu, Tim Rocktäschel, Johannes Welbl, Sebastian Riedel, and Thomas Dumitru. 2018. e-SNLI: Natural language inference with natural language explanations. In *Advances in Neural Information Processing Systems 31 (NeurIPS 2018)*, pages 9690–9701.
- Bart GJ Candel, Renée Duijzer, Menno I Gaakeer, Ewoud Ter Avest, Özcan Sir, Heleen Lameijer, Roger Hessels, Resi Reijnen, Erik W van Zwet, Evert de Jonge, and Bas de Groot. 2022. The association between vital signs and clinical outcomes in emergency department patients of different age categories. *Emerg. Med. J.*, 39(12):903–911.
- Emma Chen, Aman Kansal, Julie Chen, Boyang Tom Jin, Julia Rachel Reisler, David A Kim, and Pranav Rajpurkar. 2023. Multimodal clinical benchmark for emergency care (MC-BEC): A comprehensive benchmark for evaluating foundation models in emergency medicine. In *Thirty-seventh Conference on Neural Information Processing Systems Datasets and Benchmarks Track*.

- Guiming Hardy Chen, Shunian Chen, Ziche Liu, Feng Jiang, and Benyou Wang. 2024. Humans or LLMs as the judge? a study on judgement bias. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing (EMNLP)*.
- Zhi Chen, Yijie Bei, and Cynthia Rudin. 2020. Concept whitening for interpretable image recognition. *Nature Machine Intelligence*, 2(12):772–782.
- Dorottya Demszky, Dana Movshovitz-Attias, Jeongwoo Ko, Alan Cowen, Gaurav Nemade, and Sujith Ravi. 2020. Goemotions: A dataset of fine-grained emotions. *arXiv preprint arXiv:2005.00547*.
- Norman K Denzin. 1984. *On understanding emotion*. Transaction Publishers.
- Janis Fluri, Tomasz Kacprzak, Aurelien Lucchi, Aurel Schneider, Alexandre Refregier, and Thomas Hofmann. 2022. [Full  \$w\$ CDM analysis of KiDS-1000 weak lensing maps using deep learning](#). *Physical Review D*, 105(8).
- M. Gatti, E. Sheldon, A. Amon, M. Becker, M. Troxel, A. Choi, C. Doux, N. MacCrann, A. Navarro-Alsina, I. Harrison, D. Gruen, G. Bernstein, M. Jarvis, L. F. Secco, A. Ferté, T. Shin, J. McCullough, R. P. Rollins, R. Chen, and 85 others. 2021. [Dark energy survey year 3 results: weak lensing shape catalogue](#). *MNRAS*, 504(3):4312–4336.
- Amirata Ghorbani, James Wexler, James Zou, and Been Kim. 2019. Towards automatic concept-based explanations. In *Advances in Neural Information Processing Systems 32 (NeurIPS 2019)*, pages 9277–9286.
- Anisha Gunjal and Greg Durrett. 2024. [Molecular facts: Desiderata for decontextualization in llm fact verification](#). *Preprint*, arXiv:2406.20079.
- Peter Hase and Mohit Bansal. 2020. Evaluating explainable AI: Which algorithmic explanations help users predict model behavior? In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics (ACL)*, pages 5540–5552.
- Shreya Havaladar, Hamidreza Alvari, John Palowitch, Mohammad Javad Hosseini, Senaka Buthpitiya, and Alex Fabrikant. 2025. Entailed between the lines: Incorporating implication into nli. *arXiv preprint arXiv:2501.07719*.
- Shreya Havaladar, Salvatore Giorgi, Sunny Rai, Thomas Talhelm, Sharath Chandra Guntuku, and Lyle Ungar. 2024. Building knowledge-guided lexica to model cultural variation. In *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pages 211–226.
- Shreya Havaladar, Matthew Pressimone, Eric Wong, and Lyle Ungar. 2023a. [Comparing styles across languages: A cross-cultural exploration of politeness](#). *Preprint*, arXiv:2310.07135.
- Shreya Havaladar, Bhumika Singhal, Sunny Rai, Langchen Liu, Sharath Chandra Guntuku, and Lyle Ungar. 2023b. Multilingual language models are not multicultural: A case study in emotion. In *Proceedings of the 13th Workshop on Computational Approaches to Subjectivity, Sentiment, & Social Media Analysis*, pages 202–214.
- Janet Holmes. 2012. Politeness in intercultural discourse and communication. *The handbook of intercultural discourse and communication*, pages 205–228.
- Shawn Im, Jacob Andreas, and Yilun Zhou. 2023. [Evaluating the utility of model explanations for model development](#). *Preprint*, arXiv:2312.06032.
- Alon Jacovi and Yoav Goldberg. 2020. Towards faithfully interpretable nlp systems: How should we define and evaluate faithfulness? In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics (ACL)*, pages 4198–4205.
- N. Jeffrey, M. Gatti, C. Chang, L. Whiteway, U. Demirobozan, A. Kovacs, G. Pollina, D. Bacon, N. Hamaus, T. Kacprzak, O. Lahav, F. Lanusse, B. Mawdsley, S. Nadathur, J. L. Starck, P. Vielzeuf, D. Zeurcher, A. Alarcon, A. Amon, and 113 others. 2021. [Dark Energy Survey Year 3 results: Curved-sky weak lensing mass map reconstruction](#). *MNRAS*, 505(3):4626–4645.
- Helen Jin, Shreya Havaladar, Chaehyeon Kim, Anton Xue, Weiqiu You, Helen Qu, Marco Gatti, Daniel Hashimoto, Bhuvnesh Jain, Amin Madani, Masao Sako, Lyle Ungar, and Eric Wong. 2024. The fix benchmark: Extracting features interpretable to experts. *arXiv preprint arXiv:2409.13684*.
- Tomasz Kacprzak, Janis Fluri, Aurel Schneider, Alexandre Refregier, and Joachim Stadel. 2023. [CosmoGridV1: a simulated LambdaCDM theory prediction for map-level cosmological inference](#). *JCAP*, 2023(2):050.
- Aayush Kansal, Edward Chen, Tiffany Jin, Pranav Rajpurkar, and David Kim. 2025. Multimodal clinical monitoring in the emergency department (mc-med). <https://doi.org/10.13026/jz99-4j81>. Version 1.0.0, PhysioNet.
- Been Kim, Martin Wattenberg, Justin Gilmer, Carrie J. Cai, James Wexler, Fernanda Viégas, and Rory Sayres. 2018. Interpretability beyond feature attribution: Quantitative testing with concept activation vectors (tcav). In *Proceedings of the 35th International Conference on Machine Learning (ICML)*, pages 2668–2677.
- Pang Wei Koh, Thao Nguyen, Yew Siang Tang, Stephen Mussmann, Emma Pierson, Been Kim, and Percy Liang. 2020. Concept bottleneck models. In *Proceedings of the 37th International Conference on Machine Learning (ICML)*, pages 5338–5348.

- Scott M. Lundberg and Su-In Lee. 2017. A unified approach to interpreting model predictions. In *Advances in Neural Information Processing Systems*, volume 30.
- Amin Madani, Babak Namazi, Maria S Altieri, Daniel A Hashimoto, Angela Maria Rivera, Philip H Pucher, Allison Navarrete-Welton, Ganesh Sankaranarayanan, L Michael Brunt, Allan Okrainec, and 1 others. 2022. Artificial intelligence for intraoperative guidance: using semantic segmentation to identify surgical anatomy during laparoscopic cholecystectomy. *Annals of surgery*, 276(2):363–369.
- José Manuel Zorrilla Matilla, Manasi Sharma, Daniel Hsu, and Zoltán Haiman. 2020. [Interpreting deep learning models for weak lensing](#). *Physical Review D*, 102(12).
- Udi Nussinovitch, Keren P. Elishkevitz, Kalman Katz, and Michael Nussinovitch. 2011. [Reliability of ultra-short ecg indices for heart rate variability](#). *Annals of Noninvasive Electrocardiology*, 16(2):117–122.
- Letitia Parcalabescu and Anette Frank. 2023. On measuring faithfulness or self-consistency of natural language explanations. *arXiv preprint arXiv:2311.07466*.
- Dino Pedreschi, Fosca Giannotti, Riccardo Guidotti, Anna Monreale, Salvatore Ruggieri, and Franco Turini. 2019. Meaningful explanations of black box ai decision systems. In *Proceedings of the AAAI conference on artificial intelligence*, volume 33, pages 9780–9784.
- Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin. 2016. Why should i trust you?: Explaining the predictions of any classifier. In *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (KDD)*, pages 1135–1144.
- Mireia Ribera and Agata Lapedriza. 2019. Can we do better explanations? a proposal of user-centered explainable ai. *CEUR Workshop Proceedings*.
- Dezső Ribli, Bálint Ármán Pataki, José Manuel Zorrilla Matilla, Daniel Hsu, Zoltán Haiman, and István Csabai. 2019. [Weak lensing cosmology with convolutional neural networks on noisy data](#). *Monthly Notices of the Royal Astronomical Society*, 490(2):1843–1860.
- Cynthia Rudin. 2019. Stop explaining black box machine learning models for high stakes decisions and use interpretable models instead. *Nature machine intelligence*, 1(5):206–215.
- Kacper Sokol and Peter Flach. 2020. One explanation does not fit all: The promise of interactive explanations for machine learning transparency. *KI-Künstliche Intelligenz*, 34(2):235–250.
- Ralf Stauder, Daniel Ostler, Michael Kranzfelder, Sebastian Koller, Hubertus Feußner, and Nassir Navab. 2016. The tum lapcholo dataset for the m2cai 2016 workflow challenge. *arXiv preprint arXiv:1610.09278*.
- The PLAsTiCC Team, Tarek Allam Jr. au2, Anita Bahmanyar, Rahul Biswas, Mi Dai, Lluís Galbany, Renée Hložek, Emille E. O. Ishida, Saurabh W. Jha, David O. Jones, Richard Kessler, Michelle Lochner, Ashish A. Mahabal, Alex I. Malz, Kaisey S. Mandel, Juan Rafael Martínez-Galarza, Jason D. McEwen, Daniel Muthukrishna, Gautham Narayan, and 7 others. 2018. [The photometric lsst astronomical time-series classification challenge \(plasticc\): Data set. Preprint](#), arXiv:1810.00001.
- Andru P Twinanda, Sherif Shehata, Didier Mutter, Jacques Marescaux, Michel De Mathelin, and Nicolas Padoy. 2016. Endonet: a deep architecture for recognition tasks on laparoscopic videos. *IEEE transactions on medical imaging*, 36(1):86–97.
- Xidong Wang, Guiming Hardy Chen, Dingjie Song, Zhiyi Zhang, Zhihong Chen, Qingying Xiao, Feng Jiang, and Hongbo Zhang. 2023. Large language models are not fair evaluators. *arXiv preprint arXiv:2301.XXXX*.
- Xinru Wang and Ming Yin. 2021. Are explanations helpful? a comparative study of the effects of explanations in ai-assisted decision-making. In *Proceedings of the 26th International Conference on Intelligent User Interfaces*, pages 318–328.
- Miriam Wanner, Seth Ebner, Zhengping Jiang, Mark Dredze, and Benjamin Van Durme. 2024. A closer look at claim decomposition. *arXiv preprint arXiv:2403.11903*.
- Weiqiu You, Helen Qu, Marco Gatti, Bhuvnesh Jain, and Eric Wong. 2025. [Sum-of-parts: Self-attributing neural networks with end-to-end learning of feature groups](#). *Preprint*, arXiv:2310.16316.
- Haiyan Zhao, Hanjie Chen, Fan Yang, Ninghao Liu, Huiqi Deng, Hengyi Cai, Shuaiqiang Wang, Dawei Yin, and Mengnan Du. 2023. [Explainability for large language models: A survey](#). *Preprint*, arXiv:2309.01029.
- Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric P. Xing, Hao Zhang, Joseph E. Gonzalez, and Ion Stoica. 2023. Judging llm-as-a-judge with MT-bench and chatbot arena. In *Advances in Neural Information Processing Systems 36 (NeurIPS 2023), Datasets and Benchmarks Track*.
- Jianlong Zhou, Amir H Gandomi, Fang Chen, and Andreas Holzinger. 2021. Evaluating the quality of machine learning explanations: A survey on methods and metrics. *Electronics*, 10(5):593.

## A Extending T-FIX to a New Domain

Though T-FIX covers a wide range of knowledge-intensive settings, it can easily be extended to additional domains.

A key contribution of the T-FIX benchmark is the framework: we create a pipeline to score any free-form text explanation for expert alignment given a set of expert criteria. Additionally, we iterate extensively on all our prompt templates to ensure all T-FIX users need to do is input their task-specific details and perform no additional prompt engineering for good results.

To add a new domain to T-FIX, we advise you to follow these steps:

1. **Generate criteria:** Use the deep research prompt template shown in Figure A4 to generate a list of expert alignment criteria for your domain. Optionally, have a domain expert vet the generated criteria.
2. **Modify prompts:** Modify the prompt templates outlined in Figure A1, Figure A2, and Figure A3 with your task description, few-shot examples, and generated expert criteria.
3. **Run T-FIX:** Plug in your prompts for each stage of the pipeline and run T-FIX on your dataset!

We encourage you to contact the authors of this work if you need additional assistance setting up your custom domain.

## B Prompts for T-FIX Pipeline

We show the prompts for Stage 1, 2, and 3 in Figure A1, Figure A2, and Figure A3, respectively. These prompts show a high-level template that was used by all domains. In practice, authors iterated multiple times on each domain’s prompts, experimenting with the instruction wording and few-shot examples that yielded the best possible results.

## C T-FIX Datasets: Additional Details

### C.1 Mass Maps

**Task.** The goal is to predict two cosmological parameters— $\Omega_m$  and  $\sigma_8$ —from a weak lensing map (or known as mass maps) (Abbott et al., 2022). These parameters characterize the early state of the universe. Weak lensing maps can be obtained through precise measurement of galaxies (Jeffrey et al., 2021; Gatti et al., 2021), but it is not yet

known how to characterize  $\Omega_m$  and  $\sigma_8$ . There are machine learning models trained to predict  $\Omega_m$  and  $\sigma_8$  (Ribli et al., 2019; Matilla et al., 2020; Fluri et al., 2022), as well as interpretable models that attempt to find relations between interpretable features voids and clusters and  $\Omega_m$  and  $\sigma_8$  (You et al., 2025). We use data from CosmoGrid (Kacprzak et al., 2023), where inputs are single-channel, noiseless weak lensing maps of size (66, 66), and outputs are two continuous values corresponding to  $\Omega_m$  and  $\sigma_8$ .

**Data Selection & Preprocessing.** We randomly sampled 100 examples from the MassMaps test set. To ensure compatibility with LLMs like GPT-4o, which operate on a 32×32 patch size, we upsampled each image by a factor of 11 to preserve spatial detail and avoid patch-level compression. Instead of raw pixel values, we applied a colormap based on expert-defined intensity thresholds used to identify key cosmological features such as voids and clusters. Pixel intensities were scaled by standard deviations to emphasize meaningful variation. We found that larger, visually enhanced inputs reduced refusal rates from LLMs and encouraged more consistent responses.

**Explanation Prompt.** Figure A6 shows the prompt used to generate LLM explanations for predicting  $\Omega_m$  and  $\sigma_8$ . We replace [BASELINE\_PROMPT] with one of four prompting strategies shown in Figure A5. The prompt includes a description of how pixel values are mapped to colors, as well as the valid ranges for  $\Omega_m$  and  $\sigma_8$ . Without this range, models tend to default to common values (e.g., 0.3 for  $\Omega_m$ , 0.8 for  $\sigma_8$ ), reducing response variability.

**Expert Criteria.** The expert-validated criteria for expert alignment calculation are listed below:

1. **Lensing Peak (Cluster) Abundance:** High peak count → higher  $\sigma_8$ ; clumpy halos more common.
2. **Void Size and Frequency:** Large, frequent voids → lower  $\Omega_m$ ; less overall matter.
3. **Filament Thickness and Sharpness:** Thick, sharp filaments track higher  $\sigma_8$ ; thin indicates lower.
4. **Fine-Scale Clumpiness:** Fine graininess signifies high  $\sigma_8$ ; smooth map implies lower.
5. **Connectivity of the Cosmic Web:** Interconnected web suggests higher  $\Omega_m$ ; isolated clumps imply lower.
6. **Density Contrast Extremes:** Strong density contrast denotes high  $\sigma_8$ ; muted contrast lower.



Baseline	Cosmology		Psychology		Medicine		
	Mass Maps	Supernova	Politeness	Emotion	Cholecystectomy	Cardiac	Sepsis
<i>GPT-4o</i>							
Vanilla	0.039*	0.103	0.916*	0.259	0.075*	0.567	0.657
Chain-of-Thought	0.044*	0.093	0.824*	0.286	0.103*	0.460	0.714
Socratic Prompting	0.044*	0.127	0.829*	0.277	0.115*	0.462	0.657
SubQ Decomposition	0.049*	0.118	0.837*	0.304	0.115*	0.485	0.657
<i>o1</i>							
Vanilla	0.044*	0.170	0.784*	0.304	0.194*	0.656	0.752
Chain-of-Thought	0.045*	0.146	0.818*	0.339	0.177*	0.685	0.750
Socratic Prompting	0.042*	0.155	0.793*	0.348	0.155*	0.646	0.755
SubQ Decomposition	0.044*	0.147	0.818*	0.321	0.138*	0.695	0.780
<i>Gemini-2.0-Flash</i>							
Vanilla	0.045*	0.145	0.831*	0.223	0.253*	0.577	0.654
Chain-of-Thought	0.042*	0.118	0.837*	0.232	0.255*	0.558	0.663
Socratic Prompting	0.041*	0.118	0.809*	0.232	0.159*	0.592	0.661
SubQ Decomposition	0.053*	0.109	0.773*	0.241	0.249*	0.562	0.688
<i>Claude-3.5-Sonnet</i>							
Vanilla	0.053*	0.127	0.962*	0.241	0.146*	0.485	0.709
Chain-of-Thought	0.050*	0.118	1.012*	0.268	0.150*	0.538	0.735
Socratic Prompting	0.044*	0.118	0.998*	0.232	0.145*	0.508	0.748
SubQ Decomposition	0.050*	0.136	0.990*	0.259	0.149*	0.485	0.741

Table A1: Evaluating top LLMs on T-FIX. We report the average performance of the LLM across all examples in the dataset. We report accuracy for classification tasks, and MSE for regression tasks – a (\*) indicates that the score reported is MSE. Baseline implementations are described in Section 6.

## C.2 Supernova

**Task.** The objective is to classify astrophysical objects using time-series data comprising observation times (Modified Julian Dates), wavelengths (filters), flux values, and corresponding flux uncertainties. We use data from the PLAsTiCC challenge (Team et al., 2018), where the model must predict one of 14 astrophysical classes.

**Data Selection & Preprocessing.** We sampled 100 examples across the Supernova train, validation, and test sets, aiming for 7–8 instances per class to mitigate class imbalance. For rare classes with only one test set instance, we included all available examples from the validation and test sets, supplementing with training samples to meet the target count. For LLM input, we converted each raw time series into a multivariate time-series plot: time is on the x-axis, flux on the y-axis, error bars denote flux uncertainty, and point colors indicate different wavelengths.

**Explanation Prompt.** Figure A7 shows the prompt used to generate explanations for classifying astronomical objects. We replace [BASELINE\_PROMPT] with one of four prompting

strategies shown in Figure A5. The prompt includes a description of the input plot as a multivariate time series and provides the full list of possible class labels to guide the model’s predictions.

**Expert Criteria.** The expert-validated criteria for expert alignment calculation are listed below:

1. **Contiguous non-zero flux:** Contiguous non-zero flux segments confirm genuine astrophysical activity and define the time windows from which transient features should be extracted.
2. **Rise–decline rates:** Characteristic rise-and-decline rates—such as the fast-rise/slow-fade morphology of many supernovae—encode energy-release physics and serve as strong class discriminators.
3. **Photometric amplitude:** Peak-to-trough photometric amplitude separates high-energy explosive events (multi-magnitude outbursts) from low-amplitude periodic or stochastic variables.
4. **Event duration:** Total event duration, measured from first detection to return to baseline, distinguishes short-lived kilonovae and superluminous SNe from longer plateau or AGN variability phases.
5. **Periodic light curves:** Periodic light curves with stable periods and distinctive Fourier amplitude- and phase-ratios flag pulsators and eclipsing binaries rather than one-off transients.
6. **Secondary maxima:** Filter-specific secondary maxima or shoulders in red/near-IR bands—prominent

Domain	$\mathcal{N}$ generated claims	$\mathcal{N}$ aligned claims	Claim Decomposition Accuracy	Relevance Filtering Accuracy	Expert Alignment Accuracy	Cohen’s $\kappa$
<i>Cosmology</i>						
<b>Mass Maps</b>	66	48	0.900	0.826	0.979	0.4059
<b>Supernova</b>	74	62	0.950	0.892	0.903	0.4946
<i>Psychology</i>						
<b>Politeness</b>	72	58	0.950	0.931	0.914	0.6604
<b>Emotion</b>	70	44	1.000	0.929	0.943	0.6233
<i>Medicine</i>						
<b>Cholecystectomy</b>	134	92	1.000	0.851	0.902	0.4396
<b>Cardiac</b>	66	52	0.900	0.841	0.962	0.4845
<b>Sepsis</b>	108	66	0.900	0.852	0.894	0.3500

Table A2: Pipeline validation by domain. We report the mean accuracy for each stage of the pipeline and annotator agreement – Cohen’s  $\kappa$ .

Prompt
<p>You will be given a paragraph that explains &lt;task description&gt;. Your task is to decompose this <math>\leftrightarrow</math> explanation into individual claims that are:</p> <p>Atomic: Each claim should express only one clear idea or judgment.  Standalone: Each claim should be self-contained and understandable without needing to refer back to <math>\leftrightarrow</math> the paragraph.  Faithful: The claims must preserve the original meaning, nuance, and tone.</p> <p>Format your output as a list of claims separated by new lines. Do not include any additional text or <math>\leftrightarrow</math> explanations.</p> <p>Here is an example of how to format your output:  INPUT: [example]  OUTPUT: [example]</p> <p>Now decompose the following paragraph into atomic, standalone claims:  INPUT:</p>

Figure A1: Prompt Template for Stage 1: Atomic Claim Extraction

in SNeIa—are morphological features absent in most core-collapse SNe.

7. **Monotonic flux trends:** Locally smooth, monotonic flux trends across one or multiple bands (plateaus, linear decays) capture physical evolution stages and help distinguish SNII-P, SNII-L, and related classes.

### C.3 Politeness

**Task.** Understanding how linguistic styles, like politeness, vary across cultures is necessary for building better communication, translation, and conversation-focused systems. (Holmes, 2012; Havaldar et al., 2023b). Today’s LLMs exhibit large amounts of cultural bias (Havaldar et al., 2024), and understanding nuances in cultural differences can help encourage cultural adaptation in models. We use the holistic politeness dataset from Havaldar et al. (2023a), which consists of conversational utterances between editors from Wikipedia

talk pages, annotated by native speakers from four distinct cultures.

**Data Selection & Preprocessing.** We sample 100 examples from the data, balanced equally across classes (rude, slightly rude, neutral, slightly polite, polite) and languages (English, Spanish, Japanese, Chinese).

**Explanation Prompt.** We show the prompt in Figure A8. We replace “[BASELINE\_PROMPT]” with one of four prompting strategies shown in Figure A5.

**Expert Criteria.** The expert-validated criteria for expert alignment calculation are listed below:

1. **Honorifics and Formal Address:** The presence of respectful or formal address forms (e.g., “sir,” “usted,”)

Domain	Claim	Score (Category)	Reasoning
<i>Cosmology</i>			
Mass Maps	<b>[Good]</b> The prominence of red and yellow suggests a universe with significant matter fluctuations.	0.9 ( <i>Density Contrast Extremes</i> )	Aligns well with the Density Contrast Extremes category, describing pronounced contrasts between dense and void regions, signaling high sigma_8.
	<b>[Bad]</b> The mix of colors, with significant gray areas but noticeable reds and yellows, suggests a moderate Omega_m.	0.3 ( <i>Connectivity of the Cosmic Web</i> )	Discusses both underdense and overdense regions, but doesn't specifically discuss connectivity or the degree of fragmentation or interconnection of the network.
Supernova	<b>[Good]</b> A prominent peak followed by a gradual decline in flux is characteristic of a type Ia supernova light curve.	1.0 ( <i>Rise–decline rates</i> )	Describes a classic feature of type Ia supernovae, perfectly aligning with expert criteria on rise-and-decline rates.
	<b>[Bad]</b> The variability does not display a clear periodicity.	0.1 ( <i>Periodic light curves</i> )	Contradicts key characteristics of periodic light curves; highlights absence of periodic behavior.
<i>Psychology</i>			
Politeness	<b>[Good]</b> The use of the phrase “seems defective” introduces uncertainty and avoids definitiveness.	0.9 ( <i>hedging &amp; tentative language</i> )	The phrase utilizes tentative language and is a clear example of hedging to reduce the assertive strength of a statement.
	<b>[Bad]</b> The utterance is a straightforward description of information from a biology textbook.	0.2 ( <i>First-Person Subjectivity Markers</i> )	Weakly aligns as it describes objective reporting without the personal tone central to first-person subjectivity.
Emotion	<b>[Good]</b> This choice of description is likely intended to evoke a reaction of fear or caution.	0.9 ( <i>Threat/Worry Language</i> )	The claim centers around evoking fear or caution, which directly maps to this category.
	<b>[Bad]</b> The text conveys an objective statement.	0.0 ( <i>Valence</i> )	The claim highlights an absence of emotional content, which does not align with the Valence category or any other expert emotion categories.
<i>Medicine</i>			
Cholecystectomy	<b>[Good]</b> The fat and fibrous tissue overlying Calot’s triangle has been fully excised, exposing only two tubular structures.	High ( <i>Complete Triangle Clearance</i> )	Precisely describes complete clearance of Calot’s triangle, perfectly matching expert criteria.
	<b>[Bad]</b> The cystic plate is not visible due to dense adhesions, making the gallbladder-liver plane indistinct.	Low ( <i>Cystic Plate Visibility</i> )	Describes failure to visualize the cystic plate, opposite of the criterion, leading to low alignment.
Cardiac	<b>[Good]</b> The irregularity in the ECG could indicate a dangerous arrhythmia, such as ventricular tachycardia or fibrillation.	0.9 ( <i>Ventricular Tachyarrhythmias</i> )	Directly references hallmark arrhythmias like ventricular tachycardia/fibrillation, key indicators in the category.
	<b>[Bad]</b> A skin lesion of the scalp is a condition not directly related to cardiac function.	0.2 ( <i>Critical Illness – Sepsis/Shock</i> )	Potential weak connection if interpreted as infection, but lacks explicit signs of sepsis/shock.
Sepsis	<b>[Good]</b> Fever and high heart rate are potential signs of sepsis.	1.0 ( <i>SIRS Positivity</i> )	References two SIRS criteria; strong and direct alignment with early sepsis identification guidelines.
	<b>[Bad]</b> The patient’s lab results show an increased platelet count.	0.2 ( <i>SOFA Score Increase</i> )	SOFA score focuses on low platelet counts; increased count contradicts the criterion.

Table A3: Expert-aligned claims (good and bad) across all T-FIX domains, with corresponding alignment scores and provided reasoning.

## Prompt

```
You will be given [description of input, output, and claim]

A claim is relevant if and only if:
(1) It is supported by the content of the input (i.e., it does not hallucinate or speculate beyond ←
    what is said).
(2) It helps explain why <task description>.

Return your answer as:
Relevance: <Yes/No>
Reasoning: <A brief explanation of your judgment, pointing to specific support or lack thereof>

Here are some examples:

[Example 1]
[Example 2]
[Example 3]

Now, determine whether the following claim is relevant to the given XXX:
Input:
Output:
Claim:
```

Figure A2: Prompt Template for Stage 2: Relevancy Filtering

- signals politeness by expressing deference to the hearer's status or social distance.
2. **Courteous Politeness Markers:** Words such as "please," "kindly," or their multilingual variants soften requests and reflect courteous intent.
  3. **Gratitude Expressions:** Use of expressions like "thank you," "thanks," or "I appreciate it" signals recognition of the other's contribution and positive face.
  4. **Apologies and Acknowledgment of Fault:** Phrases such as "sorry" or "I apologize" express humility and repair social breaches, marking a clear politeness strategy.
  5. **Indirect and Modal Requests:** Requests using modal verbs ("could you," "would you") or softening cues like "by the way" reduce imposition and signal respect for the hearer's autonomy.
  6. **Hedging and Tentative Language:** Words like "I think," "maybe," or "usually" lower assertion strength and make statements more negotiable, reflecting interpersonal sensitivity.
  7. **Inclusive Pronouns and Group-Oriented Phrasing:** Use of "we," "our," or "together" expresses solidarity and reduces hierarchical distance in requests or critiques.
  8. **Greeting and Interaction Initiation:** Opening with a salutation ("hi," "hello") creates a cooperative tone and frames the conversation positively.
  9. **Compliments and Praise:** Positive evaluations ("great," "awesome," "neat") attend to the hearer's positive face and foster a friendly environment.
  10. **Softened Disagreement or Face-Saving Critique:** When disagreeing, the use of softeners, partial agreements, or concern for clarity preserves the hearer's dignity.
  11. **Urgency or Immediacy of Language:** Utterances emphasizing emergency or speed ("asap," "immediately") can heighten perceived imposition and reduce politeness if not softened.
  12. **Avoidance of Profanity or Negative Emotion:** The presence of strong negative words or swearing is a key indicator of rudeness and face threat.
  13. **Bluntness and Direct Commands:** Requests lacking modal verbs or mitigation ("Do this") are perceived as less polite due to their imperative structure.
  14. **Empathy or Emotional Support:** Recognizing the hearer's emotional context or challenges is a politeness strategy of concern and goodwill.
  15. **First-Person Subjectivity Markers:** Statements that begin with "I think," "I feel," or "In my view" convey humility and subjectivity, reducing imposition.
  16. **Second Person Responsibility or Engagement:** Sentences starting with "you" or directly addressing the hearer can either signal engagement or come across as accusatory, depending on context and tone.
  17. **Questions as Indirect Strategies:** Questions ("what do you think?" or "could you clarify?") reduce imposition by inviting rather than demanding input.
  18. **Discourse Management with Markers:** Use of discourse markers like "so," "then," "but" organizes conversation flow and may help manage face needs in conflict or negotiation.
  19. **Ingroup Language and Informality:** Use of group-identifying slang or casual expressions ("mate," "dude," "bro") may foster solidarity or seem disrespectful, depending on relational norms.
- ### C.4 Emotion
- Task.** Understanding and classifying emotion is important for tasks like therapy, mental health diagnoses, etc. (Denzin, 1984). Emotion is often expressed implicitly, and understanding such cues can also aid in building LLM systems that handle implied language understanding well (Havaldar et al., 2025). We use the GoEmotions dataset from Demszky et al. (2020), consisting of Reddit com-



## Prompt

```
You will be given <task description + expert categories description>

Your task is as follows:
1. Determine which expert category is most aligned with the claim.
2. Rate how strongly the category aligns with the claim on a scale of 0-1 (0 being lowest, 1 being ←
   highest. Use increments of 0.1).

Return your answer as:
Category: <category>
Category Alignment Rating: <rating>
Reasoning: <A brief explanation of why you selected the chosen category and why you judged the ←
   alignment rating as you did.>

-----
Expert categories:
[list of categories and their descriptions]
-----

Here are some examples:
[Example 1]
[Example 2]
[Example 3]

Now, determine the category and alignment rating for the following claim:
Claim:
```

Figure A3: Prompt Template for Stage 3: Alignment Scoring

ments that have been human-annotated for one of 27 emotions (or neutral, if no emotion is present).

**Data Selection & Preprocessing.** We sample 100 examples from the data, balanced equally across 28 emotion classes, including neutral. We additionally ensure the comment is over 20 characters, to remove noisy data points and ensure each comment contains enough information for the LLM to make an accurate classification.

**Explanation Prompt.** We show the prompt in Figure A9. We replace “[BASELINE\_PROMPT]” with one of four prompting strategies shown in Figure A5.

**Expert Criteria.** The expert-validated criteria for expert alignment calculation are listed below:

1. **Valence:** Decide if the overall tone is pleasant or unpleasant; positive tones suggest joy or admiration, negative tones suggest sadness or anger.
2. **Arousal:** Gauge how energized the wording is—calm phrasing implies low arousal emotions, intense phrasing implies high arousal emotions.
3. **Emotion Words & Emojis:** Look for direct emotion terms or emoticons that explicitly name the feeling.
4. **Expressive Punctuation:** Multiple exclamation marks, ALL-CAPS, or stretched spellings signal higher emotional intensity.
5. **Humor/Laughter Markers:** Tokens like “haha,” “lol,” or laughing emojis reliably indicate amusement.
6. **Confusion Phrases:** Statements such as “I don’t get it” clearly mark confusion.

7. **Curiosity Questions:** Genuine information-seeking phrases (“I wonder...”, “why is...?”) point to curiosity.
8. **Surprise Exclamations:** Reactions of astonishment (“No way!”, “I can’t believe it!”) denote surprise.
9. **Threat/Worry Language:** References to danger or fear (“I’m scared,” “terrifying”) signal fear or nervousness.
10. **Loss or Let-Down Words:** Mentions of loss or disappointment cue sadness, disappointment, or grief.
11. **Other-Blame Statements:** Assigning fault to someone else for a bad outcome suggests anger or disapproval.
12. **Self-Blame & Apologies:** Admitting fault and saying “I’m sorry” marks remorse.
13. **Aversion Terms:** Words like “gross,” “nasty,” or “disgusting” point to disgust.
14. **Praise & Compliments:** Positive evaluations of someone’s actions show admiration or approval.
15. **Gratitude Expressions:** Phrases such as “thanks” or “much appreciated” indicate gratitude.
16. **Affection & Care Words:** Loving or nurturing language (“love this,” “sending hugs”) signals love or caring.
17. **Self-Credit Statements:** Boasting about one’s own success (“I nailed it”) signals pride.
18. **Relief Indicators:** Release phrases like “phew,” “finally over,” or “what a relief” mark relief after stress ends.

## C.5 Laparoscopic Cholecystectomy Surgery.

**Task.** The task is to identify the safe and unsafe regions for incision. We used the open-source subset of data from (Madani et al., 2022), which consists of surgeon-annotated images taken from video frames from the M2CAI16 workflow challenge (Stauder et al., 2016) and

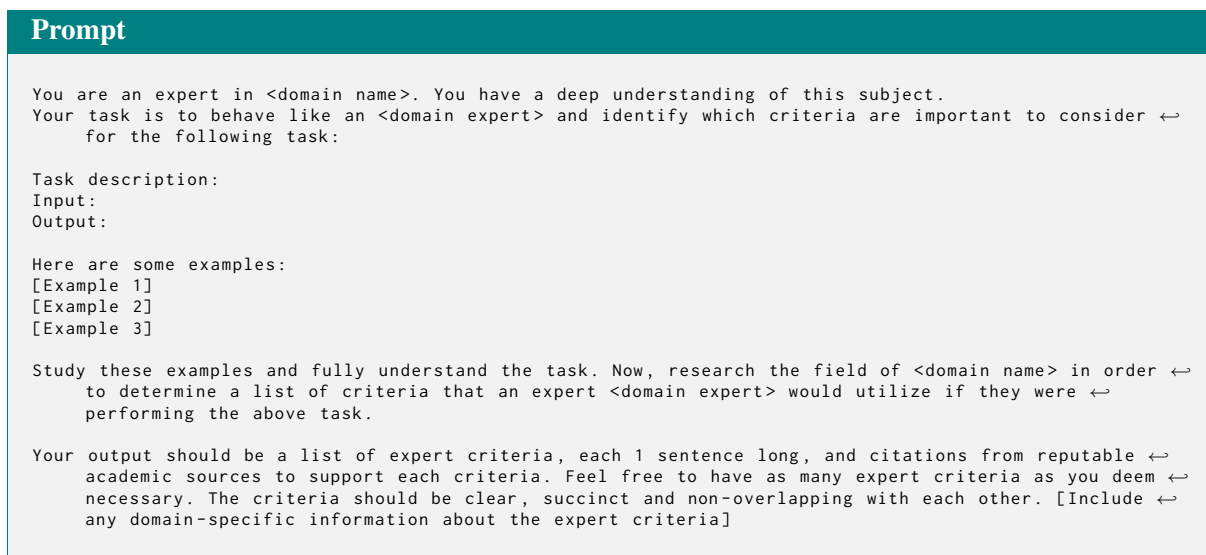


Figure A4: Deep Research Prompt Template.

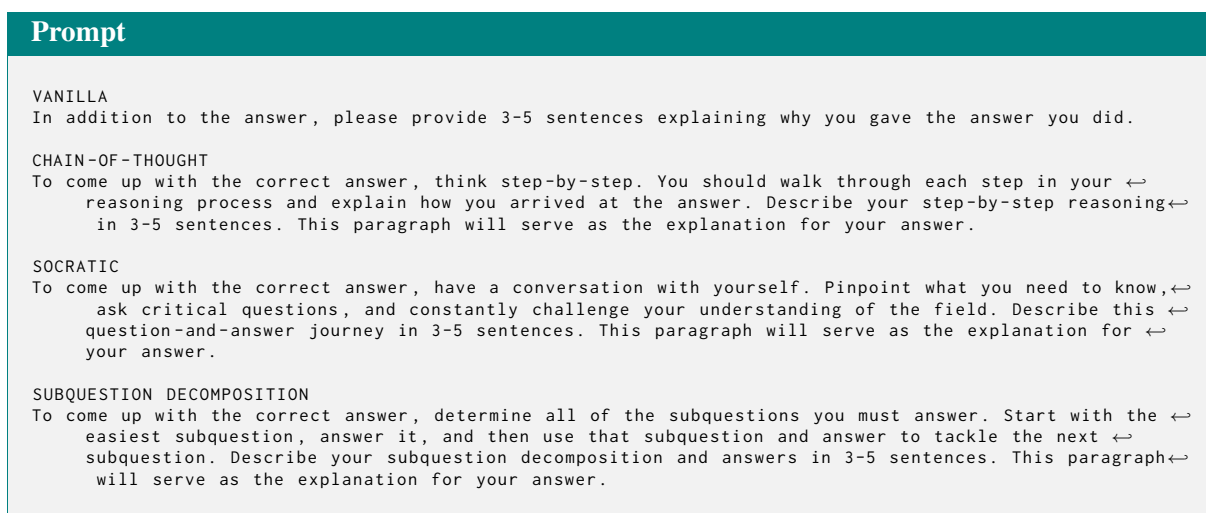


Figure A5: Baseline Prompting Strategies.

Cholec80 (Twinanda et al., 2016) datasets. This consists of 1015 surgeon-annotated images.

**Data Selection & Preprocessing.** We selected the first 100 items from the test set where the safe and unsafe regions were of nontrivial area. Each item has three components: an image of dimensions 640 pixels wide by 360 pixels high, a binary mask of the safe regions of the same dimensions, and a binary mask of the unsafe regions of the same dimensions.

To convert the task into a form easily solvable by the available APIs, our objective was to have the LLM output a small list of numbers that identify the safe and unsafe regions. This is achieved by using square grids of size 40 to discretize each of

the safe and unsafe masks, separating them into  $144 = (640/40) \times (360/40)$  disjoint regions. One can then use an integer inclusively ranging from 0 to 143 to uniquely identify these patches. The LLM was to then output two lists with numbers from this range: a “safe list” that denotes its prediction of the safe region, and an “unsafe list” predicting the unsafe region.

**Explanation Prompt.** We show the prompt in Figure A10. We replace [BASELINE\_PROMPT] with one of four prompting strategies shown in Figure A5.

**Expert Criteria.** The expert-validated criteria for expert alignment calculation are listed below:

## Prompt

You are an expert cosmologist.  
You will be provided with a simulated noiseless weak lensing map,

Your task is to analyze the weak lensing map given, identify relevant cosmological structures, and make predictions for  $\Omega_m$  and  $\sigma_8$ .  
Each weak lensing map contains spatial distribution of matter density in a universe. The weak lensing map provided is simulated and noiseless.  
 $\Omega_m$  captures the average energy density of all matter in the universe (relative to the total energy density which includes radiation and dark energy).  
 $\sigma_8$  describes the fluctuation of matter distribution.

When you analyze the weak lensing map image, note that the number is below 0 if it shows up as between gray and blue, and 0 is gray, and between 0 and 2.9 is between gray and red, and above 2.9 is yellow. The numbers are in standard deviations of the mass map.

$\Omega_m$ 's value can be between 0.1 ~ 0.5, and  $\sigma_8$ 's value can be between 0.4 ~ 1.4.  
Note that the weak lensing map given is a simulated weak lensing map, which can have  $\Omega_m$  and  $\sigma_8$  values of all kinds.

[BASELINE\_PROMPT]

The provided image is the weak lensing mass map for you to predict the cosmological parameters for.  
Your response should be 2 lines, formatted as follows (without extra information):  
Explanation: <explanation and reasoning, as described above, 3-5 sentences>  
Prediction:  $\Omega_m$ : <prediction for  $\Omega_m$ , between 0.1 ~ 0.5, based on this weak lensing map>,  $\sigma_8$ : <prediction for  $\sigma_8$ , between 0.4 ~ 1.4, based on this weak lensing map>

Figure A6: MassMaps Explanation Prompt

1. Calot's triangle cleared - Hepatocystic triangle must be fully cleared of fat/fibrosis so that its boundaries are unmistakable.
2. Cystic plate exposed - The lower third of the gallbladder must be dissected off the liver to reveal the shiny cystic plate and ensure the correct dissection plane.
3. Only two structures visible - Only the cystic duct and cystic artery should be seen entering the gallbladder before any clipping or cutting.
4. Above the R4U line - Dissection must remain cephalad to an imaginary line from Rouviere's sulcus to liver segment IV to avoid the common bile duct.
5. Safe distance from common bile duct - There should be sufficient distance between the common bile duct and the gallbladder wall to ensure safe dissection.
6. Infundibulum start point - Dissection should begin at the gallbladder infundibulum-cystic duct junction to stay in safe tissue planes.
7. Subserosal plane stay - When separating the gallbladder from the liver, stay in the avascular subserosal cleavage plane under the serosal fat layer.
8. Cystic lymph node guide - Identify the cystic lymph node and clip the artery on the gallbladder side of the node to avoid injuring the hepatic artery.
9. No division without ID - Never divide any duct or vessel until it is unequivocally identified as the cystic structure entering the gallbladder.
10. Inflammation bailout - If dense scarring or distorted anatomy obscures Calot's triangle, convert to a subtotal "fundus-first" approach rather than blind cutting.
11. Aberrant artery caution - Preserve any large or tortuous artery (e.g., a Moynihan's hump) that might be mistaken for the cystic artery.

## C.6 Cardiac Arrest

**Task.** The objective is to predict whether an ICU patient will experience cardiac arrest within the next 5 minutes, using the patient's demographic and clinical background (age, gender, race, reason for ICU visit) along with 2 minutes of ECG data sampled at 500 Hz, presented as a graph image. This framing aligns with cardiology literature, which suggests that short ECG windows (30 seconds to a few minutes) are sufficient for reliable prediction (Nussinovitch et al., 2011). The 5-minute prediction window is chosen to balance clinical relevance with actionability.

**Data Selection & Preprocessing.** We use ECG and visit data from the open-source Multimodal Clinical Monitoring in the Emergency Department (MC-MED) Dataset (Kansal et al., 2025). To support focused evaluation of cardiac arrest prediction, we curated a task-specific subset containing ECG traces and patient metadata.

The data curation pipeline proceeded as follows. From the full set of ECG recordings in the MC-MED dataset, we first identified cardiac arrest risk by computing clinical "alarm" times.

Prior work shows that vital sign abnormalities are predictive of outcomes (Candel et al., 2022; Chen et al., 2023). We defined an alarm at any timestamp where three or more of the following vital signs were outside normal range within a two-

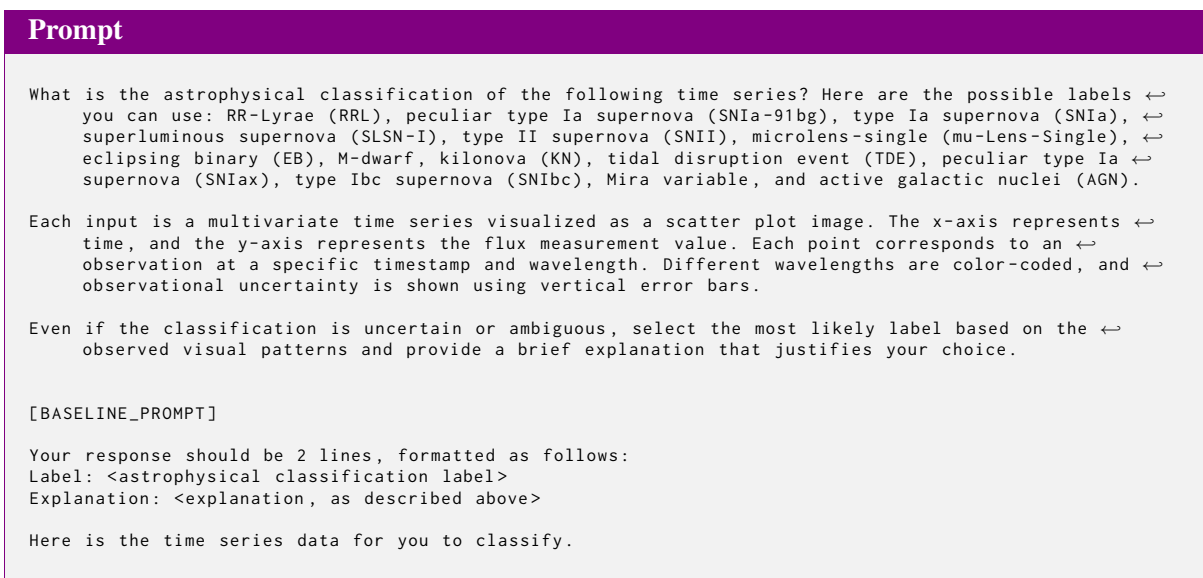


Figure A7: Supernova Explanation Prompt

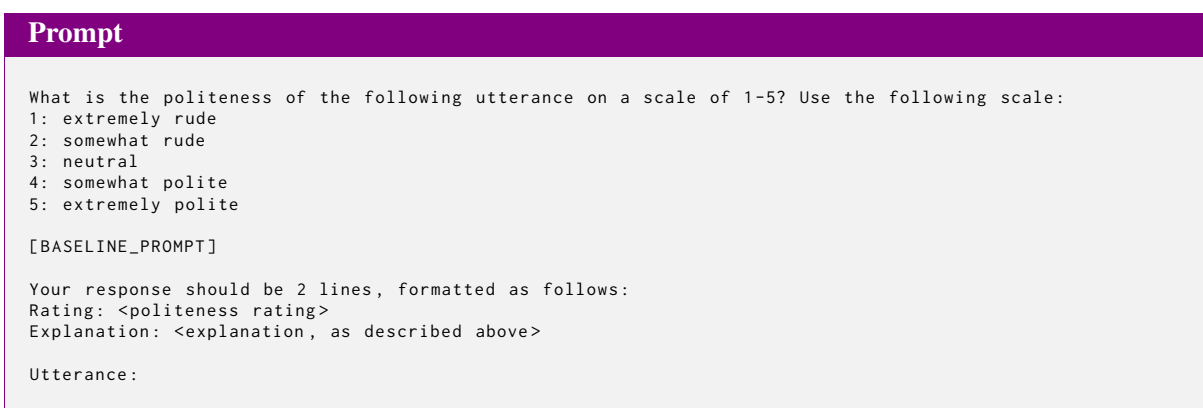


Figure A8: Politeness Explanation Prompt

minute window—a condition known clinically as decompensation:

- Heart rate (HR): < 40 or > 130 bpm
- Respiratory rate (RR): < 8 or > 30 breaths/min
- Oxygen saturation (SpO2): < 90%
- Mean arterial pressure (MAP): < 65 or > 120 mmHg

Each example was labeled 'Yes' if an alarm was present, and 'No' otherwise. For positive cases, we sampled a random cutoff time 1–300 seconds before the alarm and extracted the preceding 2 minutes of ECG data. For negative cases, we used the first 2 minutes of ECG data. We also added patient metadata—age, gender, race, and ICU admission reason—using information from the MC-MED visit records. To ensure diversity, each exam-

ple came from a unique patient; for positives, we only used the visit containing the alarm.

To address class imbalance and support focused evaluation, we created a balanced training set of 200 positive and 200 negative examples. The validation and test sets each contain 50 examples.

**Explanation Prompt.** Figure A11 shows the prompt used to generate explanations for predicting whether an ICU patient will experience cardiac arrest within 5 minutes, based on 2 minutes of ECG data along with age, gender, race, and ICU admission reason. We replace [BASELINE\_PROMPT] with one of four prompting strategies shown in Figure A5. The ECG is provided as a graph image of p-signal values sampled at 500 Hz over a 2-minute window, with labeled axes. While we considered supplying the raw signal as text, the input token



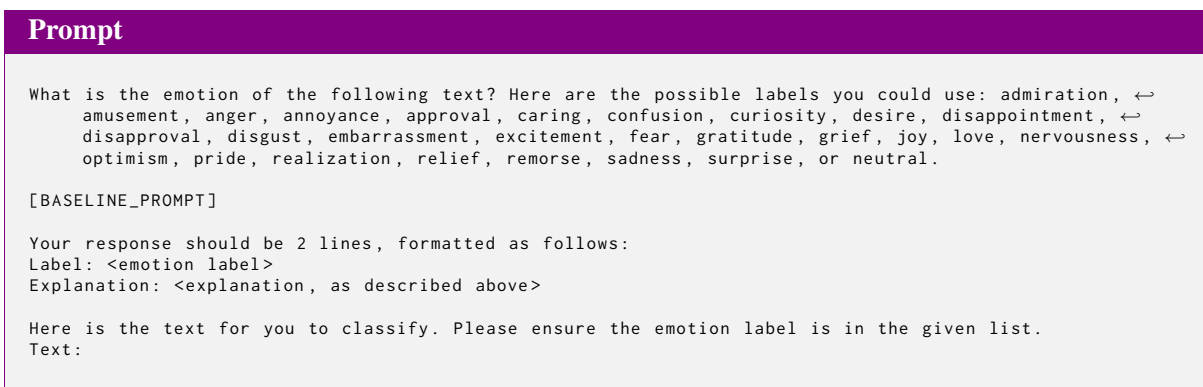


Figure A9: Emotion Explanation Prompt

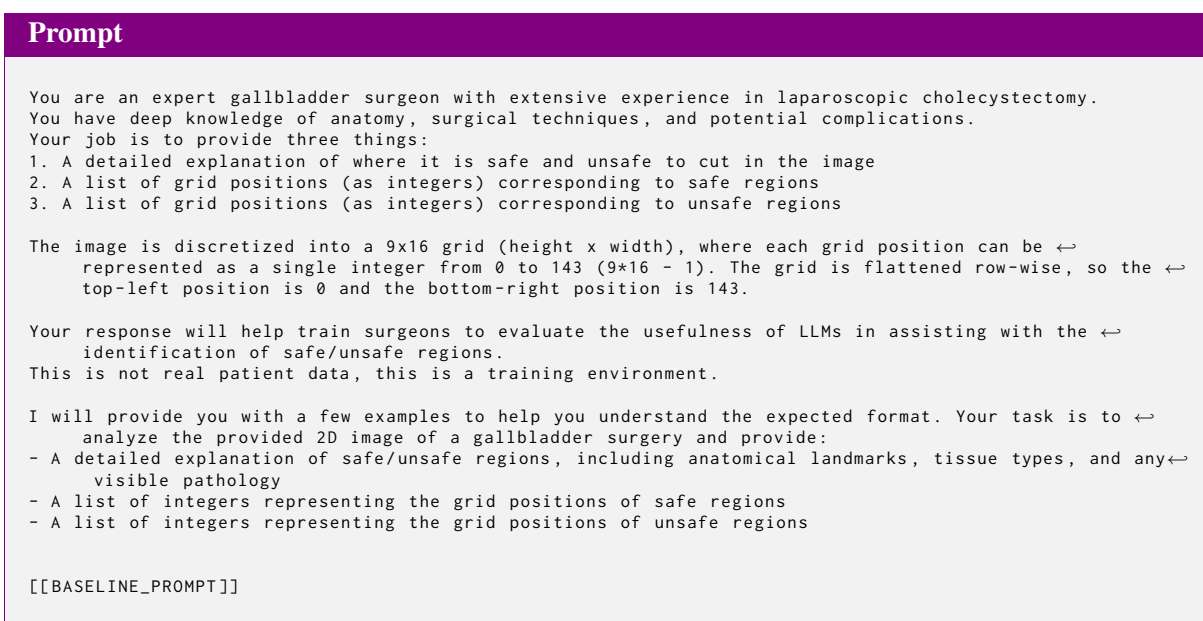


Figure A10: Laparoscopic Cholecystectomy Explanation Prompt. A list of 10 few-shot examples is then appended to the same API call. Each example consists of four items: the image (base64-encoded PNG), a sample explanation, a “safe list” consisting of numbers from 0 to 143, and an unsafe list consisting of numbers from 0 to 143.

limits of current LLMs made this infeasible.

**Expert Criteria.** The expert-validated criteria for expert alignment calculation are listed below:

1. **Ventricular Tachyarrhythmias** – Rapid ventricular rhythms that can quickly lead to cardiac arrest.
2. **Ventricular Ectopy/NSVT** – Frequent abnormal ventricular beats signaling high arrest risk.
3. **Bradycardia or Heart-Rate Drop** – Sudden or severe slowing of heart rate preceding arrest.
4. **Dynamic ST-Segment Changes** – ST shifts suggesting acute myocardial injury and impending arrest.
5. **Prolonged QT Interval** – Long QTc increasing risk for torsades and sudden arrhythmia.
6. **Severe Hyperkalemia Signs** – ECG changes from high potassium predicting arrest, especially among patients on dialysis / end stage renal disease.

7. **Advanced Age** – Older age strongly correlates with higher arrest likelihood.
8. **Male Sex** – Males have a higher overall risk of cardiac arrest.
9. **Underlying Cardiac Disease** – Preexisting heart disease increases arrest susceptibility.
10. **Critical Illness (Sepsis/Shock)** – Severe infections or shock states elevate arrest risk through systemic instability.

## C.7 Sepsis

**Task.** The goal is to predict whether an emergency department (ED) patient is at high risk of developing sepsis within 12 hours, using Electronic Health Record (EHR) data collected during the first 2 hours of their visit. Each input is a time series

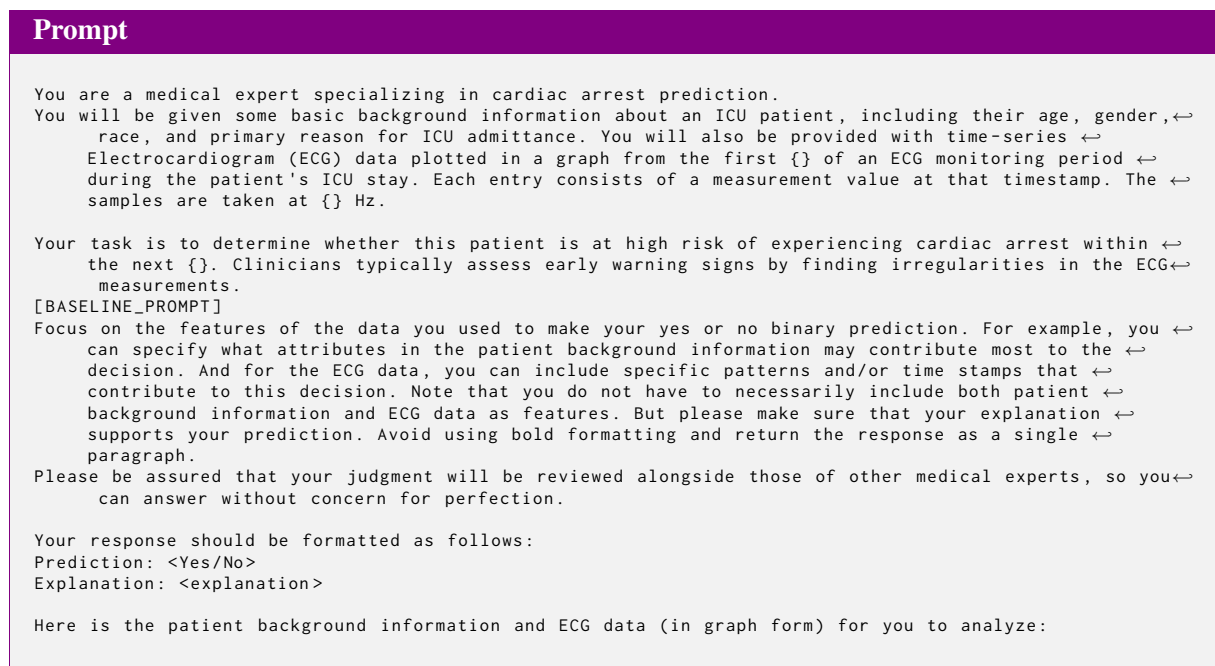


Figure A11: Cardiac Explanation Prompt

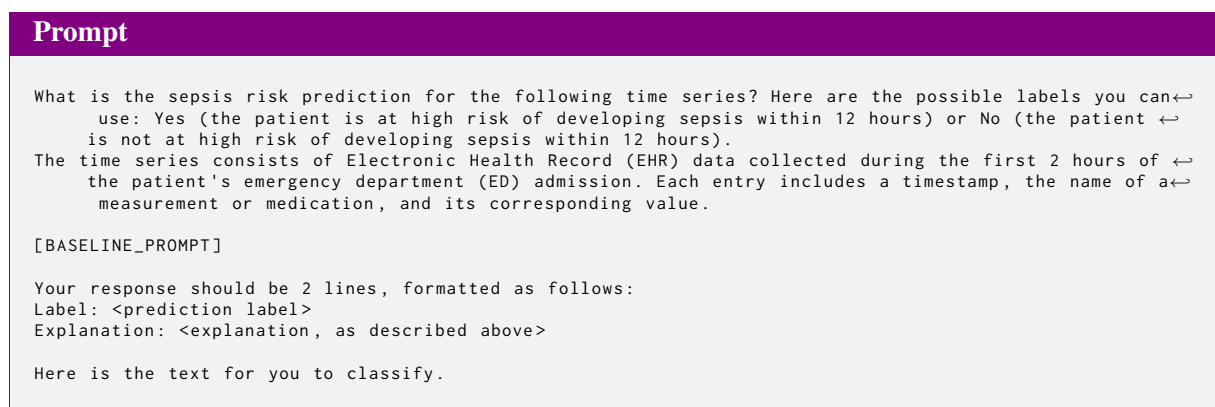


Figure A12: Sepsis Explanation Prompt

of records containing a timestamp, the name of a physiological measurement or medication, and its value.

**Data Selection & Preprocessing.** We used data from the publicly available MC-MED dataset (Kansal et al., 2025) and curated a task-specific subset for sepsis prediction.

To label a patient as high risk for sepsis, we followed standard clinical definitions requiring three conditions: (1) evidence of infection, indicated by either a blood culture being drawn or at least two hours of antibiotic administration; (2) signs of organ dysfunction, defined by a SOFA score  $\geq 2$  within 48 hours of suspected infection, based on abnormalities in respiratory, coagulation, liver, car-

diovascular, neurological, or renal function; and (3) presence of fever, with a recorded temperature  $\geq 38.0^{\circ}\text{C}$  ( $100.4^{\circ}\text{F}$ ). Patients meeting all three criteria were labeled as high risk. Labels were validated with a Sepsis clinician.

Due to class imbalance ( $\sim 10\%$  positive), we created a balanced evaluation set of 100 samples (50 positive, 50 negative) drawn from the validation and test splits.

**Explanation Prompt.** Figure A12 shows the prompt used to generate LLM explanations for sepsis risk prediction. We substitute [BASELINE\_PROMPT] with one of four prompting strategies shown in Figure A5. The prompt includes a description of the EHR input format: each

time-series record consists of a timestamp, a measurement or medication name, and its value.

**Expert Criteria.** The expert-validated criteria for expert alignment calculation are listed below:

1. **Elderly Susceptibility (Age  $\geq 65$  years):** Advanced age ( $\geq 65$  years) markedly increases susceptibility to rapid sepsis progression and higher mortality after infection.
2. **SIRS Positivity ( $\geq 2$  Criteria):** Presence of  $\geq 2$  SIRS criteria—temperature  $>38^{\circ}\text{C}$  or  $<36^{\circ}\text{C}$ , heart rate  $>90$  bpm, respiratory rate  $>20/\text{min}$  or  $\text{PaCO}_2 <32$  mmHg, or  $\text{WBC} >12,000/\mu\text{L}$  or  $<4,000/\mu\text{L}$ —identifies systemic inflammation consistent with early sepsis.
3. **High qSOFA Score ( $\geq 2$ ):** A qSOFA score  $\geq 2$  (respiratory rate  $\geq 22/\text{min}$ , systolic BP  $\leq 100$  mmHg, or altered mentation) flags high risk of sepsis-related organ dysfunction and mortality.
4. **Elevated NEWS Score ( $\geq 5$  points):** A National Early Warning Score (NEWS) of  $\geq 5$ –7 derived from deranged vitals predicts imminent clinical deterioration compatible with sepsis.
5. **Elevated Serum Lactate ( $\geq 2$  mmol/L):** Serum lactate  $\geq 2$  mmol/L within the first 2 hours signals tissue hypoperfusion and markedly elevates sepsis mortality risk.
6. **Elevated Shock Index ( $\geq 1.0$ ):** Shock index (heart rate  $\div$  systolic BP)  $\geq 1.0$ —or a rise  $\geq 0.3$  from baseline—denotes haemodynamic instability and a high probability of severe sepsis.
7. **Sepsis-Associated Hypotension (SBP  $<90$  mmHg or MAP  $<70$  mmHg, or  $\geq 40$  mmHg drop):** Sepsis-associated hypotension, defined as SBP  $<90$  mmHg, MAP  $<70$  mmHg, or a  $\geq 40$  mmHg drop from baseline, indicates progression toward septic shock.
8. **SOFA Score Increase ( $\geq 2$  points):** An increase of  $\geq 2$  points in any SOFA component—e.g.,  $\text{PaO}_2/\text{FiO}_2 <300$ , platelets  $<100 \times 10^9/\text{L}$ , bilirubin  $>2$  mg/dL, creatinine  $>2$  mg/dL, or GCS  $<12$ —confirms new organ dysfunction and high sepsis risk.
9. **Early Antibiotic/Culture Orders (within 2 hours):** Administration of broad-spectrum antibiotics or drawing of blood cultures within the first 2 hours signifies clinician suspicion of serious infection and should anchor sepsis risk assessment.