
COMPLLLM: Fine-tuning LLMs to Discover Complementary Signals for Decision-making

Anonymous Authors¹

Abstract

Multi-agent decision pipelines can outperform single agent workflows when *complementarity* holds, i.e., different agents bring unique information to the table to inform a final decision. We propose COMPLLLM, a post-training framework based on decision theory that fine-tunes a decision-assistant LLM using complementary information as reward to output signals that complement existing agent decisions. We validate COMPLLLM on synthetic and real-world tasks involving domain experts, demonstrating how the approach recovers known complementary information and produces plausible explanations of complementary signals to support downstream decision-makers.

1. Introduction

Multi-agent collaboration (including with humans) is increasingly adopted in complex decision workflows. For example, a clinician may consult a computer vision model and a written report from a radiologist to decide whether to order a biopsy for the patient; a moderator decides whether approve a post based on reviews from different human raters; a program chair decides paper acceptance based on reviewer comments and automated assessments from an LLM. A central challenge to the downstream decision-maker who must integrate inputs from upstream agents is determining complementary information: when and how does each source of information contribute information that could improve the final decision over and above existing agents’ decisions?

Most machine-learning interpretability methods are not designed to address complementarities. Explanations typically characterize why a single model produced its output, often with respect to that model’s own inputs (e.g., feature attributions, saliency, rationales). But in collaborative workflows, the bottleneck is different. For a clinician using a vision model in diagnosis, for example, what is often needed is

not simply a restatement of the model’s reasoning, but a list of features of patient information that are inconsistent with existing agent recommendations.

Our work formalizes available-but-overlooked evidence as **complementary signals**. We consider settings with two “agents”: an upstream agent that produces a recommendation Z (e.g., a vision model’s risk score on an X-ray image), and a supervisor agent that has access to potentially complementary unstructured information T (e.g., a text such as a radiology report). The goal is to identify a set of discrete, interpretable *signals* (i.e., findings) \mathbf{S} in T that capture complementary decision-relevant information not already conveyed by Z , i.e., \mathbf{S} that can improve the best-attainable decision conditioned on Z .

We propose COMPLLLM, a framework that (1) defines complementary value using the best-attainable performance on the decision problem, and (2) trains a language model to act as a complementary signal extractor from unstructured text. Formally, COMPLLLM learns a mapping from unstructured textual information (T) and recommendations (Z) to a set of structured signals, $\text{LLM} : \mathcal{T} \times \mathcal{Z} \rightarrow \mathcal{S}$, which prioritizes signals that meaningfully improve best-attainable decision quality relative to using the recommendation alone, rather than signals that are merely frequent or important. This shifts the role of “explanation” from justifying the agent recommendation to surfacing actionable complementary information that a supervisor should consider precisely because it is not already reflected in Z . COMPLLLM can be used for cases where the two agents are distinct, or where they represent the same human or model, i.e., an agent makes an initial decision using its available information, and then uses an LLM to do a second pass to surface overlooked cues in that same information.

We validate COMPLLLM across multiple domains. We first use a synthetic setting where complementary signals are controlled by construction to test the recoverability of the complementary signals. We then demonstrate use of the framework on three real-world decision-making tasks: identifying complementary signals in radiology reports that improve upon a vision model’s recommendation regarding cardiac dysfunction; identifying the signals that are more or less focal to a specific group compared to the average

¹Anonymous Institution, Anonymous City, Anonymous Region, Anonymous Country. Correspondence to: Anonymous Author <anon.email@domain.com>.

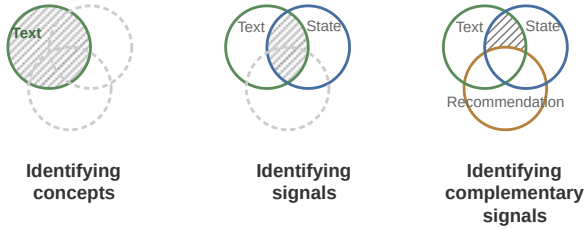


Figure 1. A graphical representation of how our work goes beyond simply identifying concepts that appear in text or signals that are predictive of a target state, by also ensuring that signals complement the existing recommendation. The diagonal stripes (twill pattern) represent the target of the corresponding methods.

human annotator in content moderation; and identifying the information that an LLM-authored review misses in human-written reviews. We further evaluate the relevance of the signals in the radiology diagnosis task by eliciting qualitative feedback from two medical domain experts (one cardiologist and one internist). For the paper reviewing task, we validate that offering the complementary signals improves the accuracy of the LLM (Gemini 2.0 Flash) judgment on the paper acceptance.

2. Related Work

2.1. Topic & Hypothesis Generation

A large literature studies topic (or *concept*) generation as a way to summarize or organize corpora: classical probabilistic models such as LDA infer latent topics as word distributions (Blei et al., 2003), while neural and embedding-based variants improve coherence and interpretability (e.g., ProdLDA/AVITM (Srivastava & Sutton, 2017), Embedded Topic Model (ETM) (Dieng et al., 2020), and other neural variational topic models such as Miao et al. (2017)). Recent representation-first pipelines treat topic modeling as clustering in embedding space (e.g., Top2Vec (Angelov, 2020), BERTopic (Grootendorst, 2022), and newer embedding-centric formulations (Angelov & Inkpen, 2024)), and LLMs increasingly support prompt-based topic generation/labeling with greater steerability (e.g., TopicGPT (Pham et al., 2023)). These methods primarily target concepts or topics (i.e., patterns grounded in text) without necessarily requiring a link to a target state.

A parallel line of work targets hypothesis (or *signal*) generation, which produces interpretable candidate factors that are grounded in text and also correlated with a state of interest (Figure 1). Examples include using SAEs to generate hypotheses from latent features with LLM interpretations (e.g., HypotheSAEs (Movva et al., 2025)), learning natural-language descriptions of distributional differences and using them in goal-driven discovery (Zhong et al., 2024; Zhou et al., 2024), and using LLM-proposed differences followed by statistical validation for causal inference on text-derived outcomes (Modarressi et al., 2025). As shown in Figure 1,

our work shifts focus to identify complementary signals, i.e., signals grounded in the supervisor information that add incremental decision value relative to an existing recommendation (not merely frequent, salient, or globally predictive signals).

2.2. Human-AI Decision-making

A growing body of work studies AI-assisted human decision-making based on its importance for legal and ethical accountability (Bo et al., 2021; Boskemper et al., 2022; Bondi et al., 2022; Schemmer et al., 2022). A recent meta-analysis (Vaccaro et al., 2024) finds that, on average, human-AI teams perform worse than the better of the two agents alone. A growing body of work seeks to evaluate and enhance complementarity in human-AI systems (Bansal et al., 2019; 2021b;a; Wilder et al., 2021; Hemmer et al., 2022; Holstein et al., 2023; Rastogi et al., 2023; Mozannar et al., 2024b; Guo et al., 2025). Some approaches explicitly incorporate human expertise in developing machine learning models or human-AI collaboration pipelines, such as by learning to defer (Mozannar et al., 2024a; Raghu et al., 2019; Keswani et al., 2022; 2021; Okati et al., 2021; Chen et al., 2022). Others develop alternative algorithms, e.g., with provable guarantees, that exploit cases where humans have additional contextual knowledge (Alur et al., 2024; Corvelo Benz & Rodriguez, 2023; Straitouri et al., 2023; De Toni et al., 2024; Arnaiz-Rodriguez et al., 2025). Closest to our work, Guo et al. (2025) provide a framework for assessing the complementary information value of arbitrary signals in a decision context. We demonstrate using complementary information as a fine-tuning objective, enabling LLM-based explanations of unexploited signals in unstructured text available at decision time.

3. Theoretical Framework

We consider a decision-making task with an upstream, recommending agent and a downstream, supervisor agent where, for each realization of the process, the recommending agent makes a recommendation $Z \in \mathcal{Z}$ based on their own features $X \in \mathcal{X}$, and the supervisor aggregates their information $T \in \mathcal{T}$ with Z to make a final decision $D \in \mathcal{D}$. The goal of the supervisor is to determine whether there is more information in T —in the form of a set of signals $\mathcal{S} \in \mathcal{S}$ —over Z that can improve their utility.

We do not assume any relationship between the recommending agent’s feature representation X and the supervisor’s information T . For example, in the radiology diagnosis task, the features of the recommending agent (vision model) X are the X-ray image, and the features of the supervisor (clinician) T are the radiology report and the patient’s medical history. In the content moderation task, the features of the recommending agent (a specific group of human annotators)

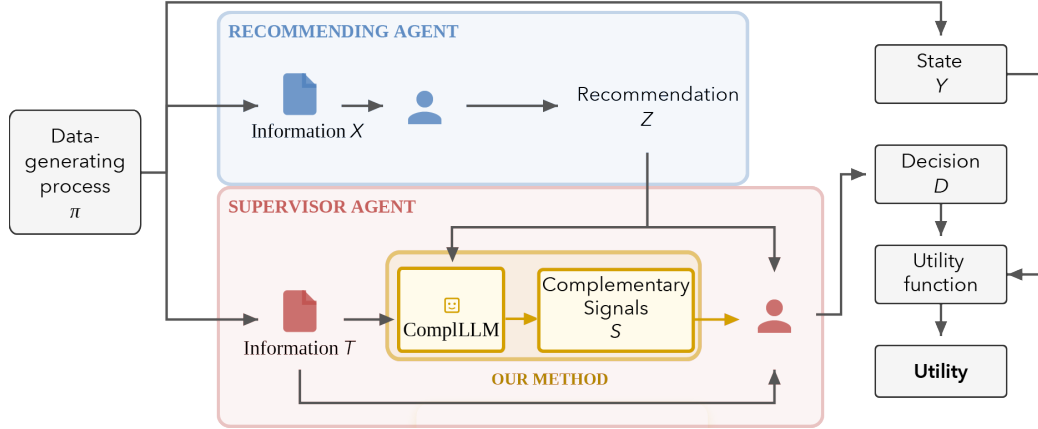


Figure 2. The “two agents” setting in our framework. The recommending agent makes a recommendation Z based on their own features X , and the supervisor agent aggregates their own information T with Z to make a final decision D .

X and the features of the supervisor (moderator) T are both the text content.

We use decision theory to characterize the decision-making process as a *decision problem* and an *information structure*. The decision problem consists of three components: the state space \mathcal{Y} , the decision space \mathcal{D} , and the utility function $U : \mathcal{D} \times \mathcal{Y} \rightarrow \mathbb{R}$. The state $y \in \mathcal{Y}$ represents the underlying state of the world that is relevant to the supervisor’s utility. Since the utility function U implicitly identifies the other two components \mathcal{Y} and \mathcal{D} , henceforce, we will use U to denote a decision problem.

The information structure is a joint distribution over the space of features of the recommending agent \mathcal{X} , the features of the supervisor \mathcal{T} , and the state \mathcal{Y} , denoted as $\pi \in \Delta(\mathcal{X} \times \mathcal{T} \times \mathcal{Y})$. Given a set of observations $\{(x_i, t_i, y_i)\}_{i \in [N]}$, we can estimate π on recommendation Z and derived discrete signals \mathbf{S} , assuming that the state is discrete, e.g., $Y \in \{0, 1\}$. We denote the probability of observing $Z = z$, $\mathbf{S} = \mathbf{s}$, and $Y = y$ as $\pi(z, \mathbf{s}, y)$. We denote the set of all possible signals as a combination of M basic signals, i.e., $\mathcal{S} = \{0, 1\}^M$. Each of the basic signals is a binary random variable, indicating the presence or absence of a certain signal, e.g., Table 1. We use upper-case letters to represent the random variables, e.g., \mathbf{S} , Z , Y , and lower-case letters to represent a value, i.e., \mathbf{s} , z , y . We denote the j -th basic signal as $\mathbf{S}_j \in \{0, 1\}$, and its absence or presence as $\mathbf{s}_j \in \{0, 1\}$.

The value of a signal is defined as the improvement that can be attained in the best attainable performance when that signal is provided. Concretely, a signal realization \mathbf{s} induces posterior distribution $\pi(y|\mathbf{s})$. The best attainable performance (\hat{U}) with the signal is the expected payoff of the decision that best-responds to the posterior distribution:

$$\hat{U}(\pi(y|\mathbf{S}), y) = \max_{d \in \mathcal{D}} \mathbb{E}_{y \sim \pi(y|\mathbf{S})} [U(d, y)].$$

We denote the language model that extracts the signals as $\text{LLM}(\cdot) : \mathcal{T} \rightarrow \mathcal{S}$. The **complementary value** of $\text{LLM}(T)$

Notation	Radiology Diagnosis instantiation
Y	Payoff state $\in \{0, 1\}$ for the presence of cardiac dysfunction.
Z	Prediction from the vision model on the probability of the presence of cardiac dysfunction $\in [0, 1]$.
X	Features of the vision model (X-ray image).
T	Features of the clinician (radiology report and patient’s medical history).
D	Final decision by the clinician $\in \{0, 1\}$ (e.g., whether to order a biopsy).
Signals in T	
\mathbf{S}_1	presence of enlarged cardiac silhouette $\in \{0, 1\}$
\mathbf{S}_2	presence of pleural abnormalities $\in \{0, 1\}$.
...	

Table 1. Notation summary (instantiated with the example of radiology-diagnosis).

over the agent’s decision Z in a decision problem U is the difference between the expected best attainable payoff on observing $\text{LLM}(T)$ and Z , and the expected best attainable payoff on observing Z :

$$\mathcal{V}^{U,Z}(\text{LLM}) := \mathbb{E}_{t,z,y \sim \pi} [\hat{U}(\pi(y|\text{LLM}(t), z), y)] - \mathbb{E}_{z,y \sim \pi} [\hat{U}(\pi(y|z), y)] \quad (1)$$

We demonstrate our framework in Table 1, using a radiology diagnosis task.

4. Methods

Given a training dataset $\{(t_i, z_i, y_i)\}_{i \in [N]}$, where t_i is the supervisor’s information, z_i is the agent’s recommendation, and y_i is the state, and a decision problem U , we fine-tune a large language model LLM parameterized by θ to extract a set of signals $\mathbf{s}_i = \text{LLM}_\theta(t_i)$ that maximizes the complementary value over the existing agent’s recommendation z_i , as defined in Equation (1).

This approach involves three steps. First, we estimate the

data-generating process as the joint distribution between the signals, agent decision, and the state. We use this estimated distribution as the reference model to define the complementary value. Second, we identify a set of complementary signals on each instance with the reference model and do supervised fine-tuning with the identified complementary signals. Last, we use reinforcement learning to further enhance the complementary value extracted by the LLM without using the labeled complementary signals.

4.1. Estimating the Data-Generating Process

We estimate the data-generating process by prompting an LLM model in two rounds. In the first round, we identify the space of possible signals, denoted as $\mathcal{S} = \{0, 1\}^M$, and in the second round, we identify whether a signal occurs in each instance (t_i, z_i, y_i) , denoted as $\{s_i^{(0)} : s_i^{(0)} \in \mathcal{S}\}_{i \in [N]}$.

We identify the space of signals by initializing a set of signals on each instance. We use an LLM reference model LLM_{ref} to output the signals that occur in the supervisor’s information t_i . The union of the signals across all instances forms the space of all possible signals in our estimate, $\mathcal{S} = \{0, 1\}^M$, where M is the number of signals in the union set. To improve stability, we sample $\zeta = 7$ outputs at temperature 0.7 and only keeps the signals with frequency larger than $N_\tau * \zeta$ ¹.

In the second round, for each signal $j \in [M]$ and each instance $i \in [N]$, we prompt LLM_{ref} again to only output whether signal j occurs in t_i . We sample $\zeta = 7$ outputs at temperature 0.7 and use a majority vote, i.e., we mark the signal as occurred, $s_{i,j}^{(0)} = 1$, if more than half of the samples indicate the signal occurs. We estimate the posterior distribution $\pi(y|\cdot)$ by a regression model (Algorithm 1) using $\{s_i^{(0)}, y_i\}_{i \in [N]}$, which we subsequently use to generate training data for SFT and the reward function for RL.

4.2. Supervised Fine-tuning (SFT) with Generated Complementary Signals

Generating Complementary Signals. We generate the training data for SFT as the complementary signals $s_i^{(1)} \in \{0, 1\}^M$ for each instance (t_i, z_i, y_i) . Given $s_i^{(0)}$ denoting the occurrence of basic signals in t_i , we generate $s_i^{(1)}$ by selecting those basic signals that both occur and have larger best-attainable payoff than the agent decision z_i on instance i (more than a minimum threshold ϵ), as these signals convey

¹ $N_\tau = \frac{z_{1-\delta/2}^2 p(1-p)}{\epsilon^2}$. $z_{1-\delta/2}$ is the $1 - \delta/2$ quantile of the standard normal distribution and p is the prior $\pi(y)$. This is to ensure ample sample size for each signal, such that the estimation error of the posterior distribution of $y|s_j$ is bounded by ϵ with confidence $1 - \delta$ for the binary state space.

decision-relevant information that is unexploited by z_i .

$$\begin{aligned} s_i^{(1)} &= (s_{i,1}^{(1)}, \dots, s_{i,M}^{(1)}) \\ \text{s.t. } s_{i,j}^{(0)} &= 1 \text{ for any } s_{i,j}^{(1)} = 1, \text{ and} \\ \hat{U}(\pi(y|s_i^{(1)}, z_i), y_i) &> \hat{U}(\pi(y|z_i), y_i) + \epsilon \end{aligned} \quad (2)$$

Supervised Fine-tuning. We fine-tune the LLM LLM_θ using the generated complementary signals $s_i^{(1)}$ as ground-truth labels. To improve the efficiency of this process ((Wang et al., 2025; Gandhi et al., 2025)), we use LLM_{ref} to generate a Chain-of-Thought (CoT) given the supervisor’s information t_i , the agent’s decision z_i , and the ground truth of $s_i^{(1)}$. We require the format of the CoT to identify the following: evidence for signals, the relevance to the state, and the complementary value relevant to the agent decision.

4.3. Reinforcement Learning

We sequentially fine-tune LLM_θ after SFT using **Group Relative Policy Optimization (GRPO)** (Shao et al., 2024) to maximize the objective given in Equation (1).

Reward function. We design the reward function using the best-attainable payoff. The reward function $R : \mathcal{S} \times \mathcal{Z} \times \mathcal{Y} \rightarrow \mathbb{R}$ rewards signals only when they increase best-attainable payoff beyond the agent’s recommendation, i.e., when they provide complementary information value. We only score signals that are *supported* by the supervisor’s information to prevent rewarding hallucinations, where s is supported if every basic signal in s occurs in the supervisor’s information, i.e., $s_{i,j}^{(0)} = 1$ for all $j \in [M]$ such that $s_j = 1$.

$$R(s, z_i, y_i) = \begin{cases} 1 & \text{if both } s \text{ and } s_i^{(1)} \text{ are empty} \\ (\hat{U}(\pi(y|s, z_i), y_i) - \hat{U}(\pi(y|z_i), y_i)) / \alpha_i & \text{else if } s \text{ are supported} \\ 0 & \text{otherwise} \end{cases} \quad (3)$$

where $\alpha_i = \hat{U}(\pi(y|s_i^{(1)}, z_i), y_i) - \hat{U}(\pi(y|z_i), y_i)$ is the best attainable payoff with the identified complementary signals in SFT training data, which we use as a normalizer.

Group Relative Policy Optimization (GRPO). We fine-tune LLM_θ using the GRPO algorithm (Shao et al., 2024). At each training step, the model generates G candidate signals s_i^1, \dots, s_i^G for each instance i . The advantage function of GRPO is defined as the reward difference between the candidate signal and the average reward of the candidates.

$$A_j = R(s_i^j, z_i, y_i) - \frac{1}{G} \sum_{k=1}^G R(s_i^k, z_i, y_i) \quad (4)$$

With the above advantage function, we train LLM_θ using the objective function defined in (Shao et al., 2024).

5. Experiments

We evaluate COMPLLLM by first testing its ability to **recover complementary signals** given by construction, then apply the approach to **evaluate practical utility** in three real-world scenarios: medical diagnosis, context moderation, and scientific paper reviewing. Across all experiments, we instantiate the “two agents” setting with the supervisor’s information t , recommending agent’s features x , and payoff state y .

5.1. Model and Training Details

We use *Qwen3-8B* to generate the signals for estimating the data-generating process and as the backbone language model for fine-tuning. For SFT, we train with 2 epochs, with a learning rate of 5×10^{-6} and a cosine learning rate scheduler. For GRPO, we initialize the training with 10 epochs but early stop when the improvement of the reward on validation set is less than 0.01, with a learning rate of 1×10^{-6} and a cosine learning rate scheduler. We sample 12 candidate completions for each instance in GRPO. For the hyperparameters that determine the threshold of the sample size N_τ for the data-generating process estimation, we use $\epsilon = 0.1$ and $\delta = 0.05$. We use the output token length of 1,500 tokens.

5.2. Baselines & Benchmark

We compare COMPLLLM with the following methods:

ZERO-SHOT and FEW-SHOT learning. We provide LLMs with task-specific instructions (zero-shot) and optionally with three demonstration examples (few-shot). We choose the demonstration examples from the generated SFT training dataset described in Section 4.2.

Topic-generation: BERTOPIC (Grootendorst, 2022) We compare with topic-generation methods that are agnostic to the existing agent decisions and the decision problems. BERTOPIC is a neural topic modeling method that produces topics by clustering text embeddings. We fit a multivariate logit model to predict the state y from the topic ID and the recommendation z . We assign each topic a natural language description by prompting Qwen3-8B with a sample of documents.

Hypothesis-generation: HYPOTHESAE (Movva et al., 2025) We compare with recent hypothesis-generation methods that take into account the decision problem and payoff-relevant state, but are agnostic to the existing agent decisions. HYPOTHESAE clusters the documents by the predictive power of selected neuron activations on the state y , and then produces hypotheses by prompting an LLM with a sample of documents in each cluster.

Benchmark To benchmark the value that text t and recommendation z have for predicting the ground-truth label y , we fine-tune Qwen3-8B with GRPO to predict the state,

serving as a non-interpretable (i.e., without discrete signals) benchmark for achievable performance.

To ensure comparability, we use the same training, validation, and test splits for all the methods. We use the same model as used in COMPLLLM (Qwen3-8B), to summarize and annotate the topics in the BERTopic method and the hypotheses of the HypotheSAE method. See the hyperparameters of baseline & methods in Appendix E.

5.3. Tasks, Datasets, and Metrics

5.3.1. COMPLEMENTARITY BY CONSTRUCTION

We test how well COMPLLLM recovers artificially induced complementary signals, using the medical decision problem of diagnosing cardiac dysfunction defined in Table 1.

Dataset and task. We use radiology reports as t from the MIMIC-CXR dataset (Johnson et al., 2023) with a ground-truth signal set s^* derived from the CheXpert labels (Irvin et al., 2019) on these reports. To generate the state y and upstream recommendation z for each instance, we fit a logistic regression model that regresses real-world cardiac dysfunction labels on the CheXpert labels. We derived the cardiac dysfunction labels from blood tests related to cardiac dysfunction (*troponin* and *NT-proBNP*) in MIMIC-IV, using domain-suggested age-cutoffs (Mueller et al., 2019; Heidenreich et al., 2022) to threshold into binary values. We generate y by inputting the full set of report labels from CheXpert to this model, and generate z by holding out a subset of labels (*Edema* and *Pleural Effusion*) so they are, by construction, predictive of y but not represented in z . To be closer to the realistic setting where doctors are assisted by a probabilistic prediction, we threshold the prediction by 0.5 to get a binary $y \in \{0, 1\}$ and keep $z \in [0, 1]$ as the probabilistic predictions. We use 6K/2K/4K instances for training/validation/test.

Metrics. We evaluate how well the output signals recover the ground-truth complementary signals, and the complementary information value of the output signals over the existing agent’s decisions.

- **Surface Similarity:** For each test instance, we prompt Qwen3-14B to score surface similarity for every pair of presented ground-truth signal (i.e., a ground-truth signal with label 1 for that instance) and output signal. Specifically, we prompt Qwen3-14B to assign scores of 1, 0.5, and 0 for signals that are the same, related, or distinct, respectively. For stability, we report the average similarity score over all matched pairs given 7 repetitions and temperature 0.7. For each presented ground-truth signal, we take the maximum similarity over output signals (i.e., the best match), and report these per-ground-truth maxima averaged across all instances.
- **F1 Similarity:** For each instance, we compute the F1 score between the presence of the ground-truth signals

and whether a output signal is matched with it (i.e., has surface similarity ≥ 0.5). We report the average F1 score over all instances and all ground-truth signals.

- **Complementary Information Value (Improvement on accuracy by signals):** We report the complementary information value $\mathcal{V}^{U,Z}(\text{LLM})$ from Equation (1) with accuracy $U(d, y) = \mathbb{1}[d = y]^2$ as the decision problem. This represents the improvement in best-attainable performance from observing extracted signals $\text{LLM}(T)$ in addition to Z . We use a multivariate logit model to fit the payoff state y on the extracted signals $\text{LLM}(T)$ and the existing agent’s decision Z to compute the best-attainable performance.

The first and second metrics are derived from Zhong et al. (2024), who similarly evaluate recovery of ground-truth signals.

5.3.2. COMPLEMENTARITY IN DECISION-MAKING

We test COMPLLLM on three realistic decision-making tasks: medical diagnosis, content moderation, and scientific paper reviewing. More details in Appendix A.

Medical Diagnosis: MIMIC-CXR (Johnson et al., 2023). We investigate what information in human-generated radiology reports can be used to improve a vision model’s predictions of cardiac dysfunction on X-ray images. We use the radiology reports from MIMIC-CXR as t , the chest X-ray images from MIMIC-CXR as x , and the cardiac dysfunction labels derived from blood tests in MIMIC-IV as $y \in \{0, 1\}$, as described in the above experiment. We fine-tune the CXR foundation model (Sellergren et al., 2023) on the cardiac dysfunction labels based on the X-ray images and the blood test results and use its probabilistic predictions as the agent recommendation $z \in [0, 1]$. We use the same training/validation/test as above.

Content Moderation: DICES (Aroyo et al., 2023). We investigate what toxicity cues in human-LLM conversations a specific group of human annotators use differently from the majority-vote average annotations. We use the DICES dataset, which contains 115K human annotations of toxicity on 1.3K human-LLM conversations with demographic information about annotators. We use annotations from one demographic group (i.e., 7K annotations from the largest demographic group of annotators: Asian millennial women with college degree or higher) as the recommending agent $z \in \{0, 1\}$ for non-toxicity and toxicity, and the average human annotation as the state $y \in \{0, 1\}$. We use the conversation text as t . We use 4K/1K/2K instances for training/validation/test respectively.

Scientific Paper Reviewing: Review5K (Weng et al., 2025). We investigate what information in human-written reviews is missed by an LLM review’s decision. We use the

²We threshold the probabilistic decision by 0.5 into binary.

Method	SURF.	F1
COMPLLLM	0.98	0.67
ZERO-SHOT	0.91	0.42
FEW-SHOT	0.95	0.38
BERTOPIIC	0.84	0.17
HYPOTHESAE	0.90	0.38

Table 2. Results comparing average surface similarity and F1 score of the extracted signals by each method on the synthetic dataset.

Review5K dataset, which contains 5K scientific papers with human-written reviews and the final decision for each paper. We generate the LLM review by prompting Gemini 2.0 Flash to review the paper and provide a judgement on the acceptance of the paper between “Accept”, “Unsure”, and “Reject”. We use the human-written review text as t and the final judgement of the LLM review as the recommending agent $z \in \{0, 0.5, 1\}$ for “Reject”, “Unsure”, and “Accept” respectively. We use the final decision (made by human AC in the real review process) on the paper as the state $y \in \{0, 1\}$ for “Reject” and “Accept” respectively. We use 3K/1K/1K instances for training/validation/test respectively.

Metrics.

- **Complementary Information Value:** Same as the metric for synthetic data.
- **Breadth:** We fit a multivariate logit model of y on z and $\text{LLM}(t)$, and report the number of extracted signals whose coefficients are significantly non-zero³ when controlling for z (with multiple-comparisons correction). This captures how many extracted signals add unique information beyond the recommendation.
- **Qualitative assessment:** We conducted interviews with two practicing physicians—one cardiologist (P1) and one internist (P2)—both of whom are also professors of medicine. During the sessions, we walked through four patient cases. For each case, we first showed the radiology report, image, and vision model prediction, and asked for their agreement and reasons from the report or image on why they agreed or disagreed. We then revealed the complementary signals, and asked them assess their relevance and whether they aligned or not with their domain knowledge. Finally, we present an updated prediction based on a multivariate logit model predicting cardiac dysfunction from the surfaced signals and agent decisions, and asked if they trusted this prediction more or less than the original.

6. Results

6.1. Complementarity by Construction

COMPLLLM recovers the complementary signals on the synthetic dataset. Table 2 shows that **COMPLLLM beats**

³We used a Bonferroni-correct p-value threshold of 5×10^{-3} .

Discovering Complementary Signals for Decision-making

Table 3. Signals on MIMIC-CXR dataset that are significant with $p < 0.05$ after Bonferroni correction. ‘ Δ Acc.’ represents the marginal accuracy improvement each signal s_j marginally provides over the other signals and the agent recommendation. ‘Total Acc.’ shows the combined accuracy when all significant signals are included, i.e., sum of the signals’ marginal gain. ‘# Sig’ indicates the number of statistically significant signals discovered by each method.

Source	Total Acc.	# Sig	Significant Signals	Δ Acc.
Agent Decision	0.819	-	-	-
COMPLLLM	0.839	12	<i>negative pneumothorax</i>	+0.0048
			<i>positive pleural effusion</i>	+0.0040
			<i>negative pleural effusion</i>	+0.0032
			<i>positive pulmonary congestion</i>	+0.0024
			<i>positive cardiac silhouette enlargement</i>	+0.0016
			<i>positive cardiac enlargement</i>	+0.0008
			<i>uncertain pleural effusion</i>	+0.0008
			<i>uncertain cardiomegaly</i>	+0.0006
			<i>positive pulmonary vascular congestion</i>	+0.0006
			<i>negative pulmonary edema</i>	+0.0002
			<i>negative cardiac enlargement</i>	+0.0002
			<i>positive cardiomegaly</i>	+0.0002
HYPOTHESAE	0.831	6	<i>mentions no pleural effusion</i>	+0.0032
			<i>mentions absence of pleural effusion or pneumothorax</i>	+0.0024
			<i>mentions pulmonary edema</i>	+0.0020
			<i>mentions presence of pleural effusion</i>	+0.0020
			<i>mentions presence of pulmonary edema</i>	+0.0018
			<i>mentions pulmonary edema or interstitial opacities</i>	+0.0008
ZERO-SHOT	0.825	5	<i>cardiac enlargement</i>	+0.0014
			<i>pleural effusion</i>	+0.0012
			<i>bilateral pleural effusion</i>	+0.0012
			<i>interstitial edema</i>	+0.0010
			<i>pulmonary vascular engorgement</i>	+0.0010
FEW-SHOT	0.823	3	<i>no pleural effusion</i>	+0.0014
			<i>positive interstitial edema</i>	+0.0012
			<i>positive cardiac silhouette enlargement</i>	+0.0010
BERTOPIC	0.818	1	<i>endotracheal tube position</i>	+0.0005

all baseline methods in surface similarity and F1 score.

As shown in Table 4, the signals extracted by COMPLLLM cover the ground-truth complementary signals. We also find that fine-tuning helps recover ground-truth complementary signals, i.e., though few-shot learning outputs signals that are more similar to the ground-truth signals than zero-shot learning, it does not predict those signals at the correct instance (i.e., low F1 score). The signals extracted by COMPLLLM also provide the highest complementary information value over the existing agent’s decision recommendation in the synthetic dataset across all the methods (except for the non-interpretable benchmark) in Figure 3.

6.2. Complementarity in Decision-making

Figure 3 shows the complementary information value of the signals for the real-world datasets. Tables 3, 5 and 6 shows all significant signals on the MIMIC-CXR, Review5k, and DICES datasets.

6.2.1. IDENTIFYING COMPLEMENTARY SIGNALS

COMPLLLM identifies the signals with the highest complementary information value among the methods in our experiments. Comparing to the baseline methods, COMPLLLM is the only method that extracts signals with significant complementary information value (i.e., with a non-overlapping confidence interval with the accuracy of the recommending agent) in all three datasets (as shown in Figure 3). This implies that designing for complementarity is necessary for improving over the existing agent’s decision recommendation. We also find that COMPLLLM achieves a comparable performance with the non-interpretable benchmark in MIMIC-CXR dataset and DICES dataset (Figure 3), suggesting that extracting complementary signals does not harm the predictive power of the same LLM model in the decision problem in these two tasks. We also find that COMPLLLM extracts more signals with complementary information value than the baseline methods in all three datasets (Table 3), suggesting that the signals extracted by COMPLLLM are more diverse and comprehensive.

385
386
387
388
389
390
391
392
393
394
395
396
397
398
399
400
401
402
403
404
405
406
407
408
409
410
411
412
413
414
415
416
417
418
419
420
421
422
423
424
425
426
427
428
429
430
431
432
433
434
435
436
437
438
439

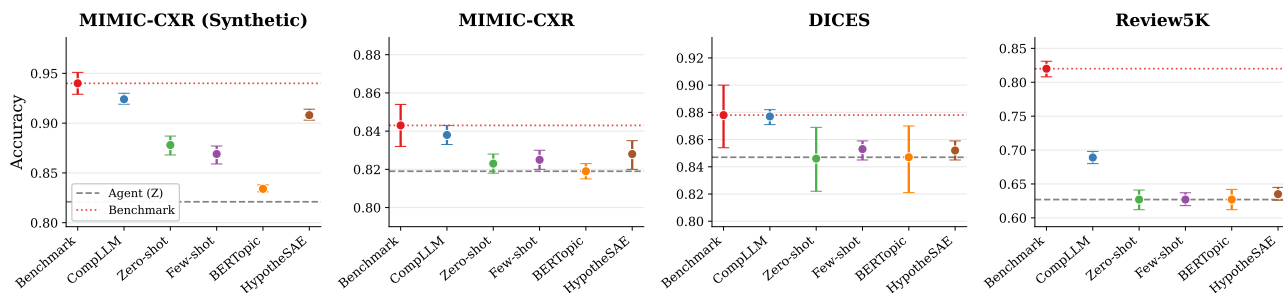


Figure 3. Expected accuracy given the extracted signals and the agent’s recommendation by each method. Dashed lines represent agent decision accuracy and the accuracy of the benchmark method. Error bars depict bootstrapped 95% confidence intervals (N=5000).

6.2.2. QUALITATIVE FEEDBACK ON COMPLLLM FOR MEDICAL DIAGNOSIS

The physicians viewed multiple overlapping cases (3 and 4 respectively, for a total of 6). Both **appraised the complementary signals as aligning with their domain knowledge**, with two exceptions. The first was a case where P1 noticed a signal that they believed provided relevant complementary information (the patient’s history of cardiac problems), but was missed by COMPLLLM. Additionally, both physicians indicated that one of the complementary signals for a case (Figure 10) both saw, which, though named as *negative_pleural_effusion*, referred to absence of pleural effusion and pneumothorax (Figure 11), was only half relevant, as the second condition neither associated strongly with cardiac dysfunction or the most likely alternative condition. In one other case, P2 did not dispute the complementary signals, but felt they might be consistent with—rather than adding new information over—the original prediction.

In an exit interview, both physicians described seeing value for COMPLLLM in practice, including to assist in creating lists of supporting versus contradictory evidence that clinicians already make (P1, P2), and to support ER doctors and general internists who have less specific training in cardiology (P1).

6.2.3. IMPROVED LLM DECISIONS ON PAPER REVIEWING

We validate the usefulness of the complementary signals by giving the signals to a proxy, LLM decision-maker (Gemini 2.0 Flash) and seeing if they improve its performance in the scientific paper reviewing task. Figure 4 shows the accuracy of the LLM’s paper review decision relative to the human AC’s ground truth decision on the Review5K dataset. We test two settings over the baseline of giving the LLM only the paper text: also giving the LLM the human review text, and also giving the LLM the human review text and the complementary signals. We find that the LLM with the complementary signals as input achieves higher accuracy than the LLM only taking the paper text and the human review text as input, (79.7% [95% CI: 78.9%, 80.5%] vs.

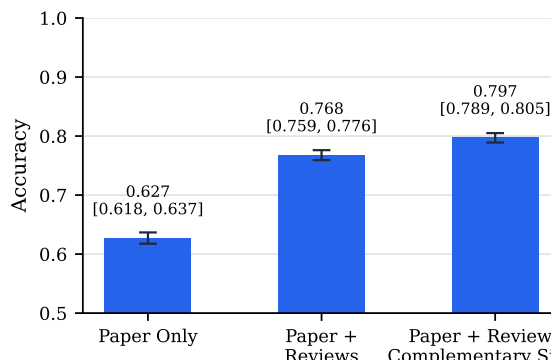


Figure 4. Accuracy of LLM paper review decisions on the Review5K dataset. Error bars depict bootstrapped 95% confidence intervals (N=5000).

78.8% [95% CI: 75.9%, 77.6%]).

7. Limitations

COMPLLLM uses the estimated data-generating process as the reference model to fine-tune and reinforce the LLM. We expect this approach to identify all important signals. However, there is a risk that COMPLLLM drops rare but important signals, i.e., signals that have frequency lower than N_τ in the dataset but would be considered important to domain experts.

We identify the signals and their decision-relevant value by the posterior distribution of the state and best-attainable performance given the occurrence of the signal as determined by the reference model. However, when the recommending agent and supervisor are the same, as in AI-assisted human decision workflows where a human first arrives at an independent judgment, then consults an AI, providing complementary signals may induce learning such that the predicted complementary value of signals no longer holds in hindsight. Future work could extend this by formalizing a continual learning problem that accounts for changes to human beliefs about the state after viewing the complementary signals and updates the LLM model correspondingly.⁴

⁴The data and code to reproduce our experimental results are available at https://osf.io/z5gyj/overview?view_

Impact Statement

Our work advances the field of human-AI or multi-agent collaboration and decision-making, which stands to contribute to a number of public-facing and scientific domains. To the best of our knowledge, there are no particular negative social consequences imposed by our work compared to machine learning research in general.

We use the DICES dataset to study differences in labeling patterns. Any deployment should include domain-appropriate oversight and evaluation to avoid over-reliance on automatically surfaced signals.

References

Alur, R., Raghavan, M., and Shah, D. Distinguishing the indistinguishable: Human expertise in algorithmic prediction. *arXiv preprint arXiv:2402.00793*, 2024.

Angelov, D. Top2vec: Distributed representations of topics. *arXiv preprint arXiv:2008.09470*, 2020.

Angelov, D. and Inkpen, D. Topic modeling: Contextual token embeddings are all you need. In *Findings of the Association for Computational Linguistics: EMNLP 2024*, pp. 13528–13539, 2024.

Arnaiz-Rodriguez, A., Benz, N. C., Thejaswi, S., Oliver, N., and Gomez-Rodriguez, M. Towards human-ai complementarity in matching tasks. *arXiv preprint arXiv:2508.13285*, 2025.

Aroyo, L., Taylor, A., Diaz, M., Homan, C., Parrish, A., Serapio-García, G., Prabhakaran, V., and Wang, D. Dices dataset: Diversity in conversational ai evaluation for safety. *Advances in Neural Information Processing Systems*, 36:53330–53342, 2023.

Bansal, G., Nushi, B., Kamar, E., Weld, D. S., Lasecki, W. S., and Horvitz, E. Updates in human-ai teams: Understanding and addressing the performance/compatibility tradeoff. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, pp. 2429–2437, 2019.

Bansal, G., Nushi, B., Kamar, E., Horvitz, E., and Weld, D. S. Is the most accurate ai the best teammate? optimizing ai for teamwork. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, pp. 11405–11414, 2021a.

Bansal, G., Wu, T., Zhou, J., Fok, R., Nushi, B., Kamar, E., Ribeiro, M. T., and Weld, D. Does the whole exceed its parts? the effect of ai explanations on complementary team performance. In *Proceedings of the 2021*

only=c32fcf0a14b54b9fa6ed22ddb8d1f774.

CHI Conference on Human Factors in Computing Systems, CHI '21, New York, NY, USA, 2021b. Association for Computing Machinery. ISBN 9781450380966. doi: 10.1145/3411764.3445717. URL <https://doi.org/10.1145/3411764.3445717>.

Blei, D. M., Ng, A. Y., and Jordan, M. I. Latent dirichlet allocation. *Journal of machine Learning research*, 3(Jan): 993–1022, 2003.

Bo, Z.-H., Qiao, H., Tian, C., Guo, Y., Li, W., Liang, T., Li, D., Liao, D., Zeng, X., Mei, L., et al. Toward human intervention-free clinical diagnosis of intracranial aneurysm via deep neural network. *Patterns*, 2(2), 2021.

Bondi, E., Koster, R., Sheahan, H., Chadwick, M., Bachrach, Y., Cemgil, T., Paquet, U., and Dvijotham, K. Role of human-ai interaction in selective prediction. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 36, pp. 5286–5294, 2022.

Boskemper, M. M., Bartlett, M. L., and McCarley, J. S. Measuring the efficiency of automation-aided performance in a simulated baggage screening task. *Human factors*, 64(6):945–961, 2022.

Chen, C., Feng, S., Sharma, A., and Tan, C. Machine explanations and human understanding (2022). URL: <http://arxiv.org/abs/2202.04092>, 2022.

Corvelo Benz, N. and Rodriguez, M. Human-aligned calibration for ai-assisted decision making. *Advances in Neural Information Processing Systems*, 36:14609–14636, 2023.

De Toni, G., Okati, N., Thejaswi, S., Straitouri, E., and Rodriguez, M. Towards human-ai complementarity with prediction sets. *Advances in Neural Information Processing Systems*, 37:31380–31409, 2024.

Dieng, A. B., Ruiz, F. J., and Blei, D. M. Topic modeling in embedding spaces. *Transactions of the Association for Computational Linguistics*, 8:439–453, 2020.

Gandhi, K., Chakravarthy, A., Singh, A., Lile, N., and Goodman, N. D. Cognitive behaviors that enable self-improving reasoners, or, four habits of highly effective stars. *arXiv preprint arXiv:2503.01307*, 2025.

Grootendorst, M. Bertopic: Neural topic modeling with a class-based tf-idf procedure. *arXiv preprint arXiv:2203.05794*, 2022.

Guo, Z., Wu, Y., Hartline, J., and Hullman, J. The value of information in human-ai decision-making. *arXiv preprint arXiv:2502.06152*, 2025.

- 495 Heidenreich, P. A., Bozkurt, B., Aguilar, D., Allen, L. A.,
 496 Byun, J. J., Colvin, M. M., Deswal, A., Drazner, M. H.,
 497 Dunlay, S. M., Evers, L. R., et al. 2022 aha/acc/hfsa
 498 guideline for the management of heart failure: a report
 499 of the american college of cardiology/american heart as-
 500 sociation joint committee on clinical practice guidelines.
 501 *Journal of the American College of Cardiology*, 79(17):
 502 e263–e421, 2022.
- 503 Hemmer, P., Schemmer, M., Kühn, N., Vössing, M., and
 504 Satzger, G. On the effect of information asymmetry in
 505 human-ai teams. *arXiv preprint arXiv:2205.01467*, 2022.
- 506 Holstein, K., De-Arteaga, M., Tumati, L., and Cheng, Y. To-
 507 ward supporting perceptual complementarity in human-
 508 ai collaboration via reflection on unobservables. *Pro-
 509 ceedings of the ACM on Human-Computer Interaction*, 7
 510 (CSCW1):1–20, 2023.
- 511 Irvin, J., Rajpurkar, P., Ko, M., Yu, Y., Ciurea-Ilicus, S.,
 512 Chute, C., Marklund, H., Haghighi, B., Ball, R., Shpan-
 513 skaya, K., et al. Chexpert: A large chest radiograph
 514 dataset with uncertainty labels and expert comparison. In
 515 *Proceedings of the AAAI conference on artificial intelli-
 516 gence*, volume 33, pp. 590–597, 2019.
- 517 Johnson, A. E., Pollard, T. J., Berkowitz, S. J., Greenbaum,
 518 N. R., Lungren, M. P., Deng, C.-y., Mark, R. G., and
 519 Horng, S. Mimic-cxr, a de-identified publicly available
 520 database of chest radiographs with free-text reports. *Sci-
 521 entific data*, 6(1):317, 2019.
- 522 Johnson, A. E., Bulgarelli, L., Shen, L., Gayles, A., Sham-
 523 mout, A., Horng, S., Pollard, T. J., Hao, S., Moody, B.,
 524 Gow, B., et al. Mimic-iv, a freely accessible electronic
 525 health record dataset. *Scientific data*, 10(1):1, 2023.
- 526 Keswani, V., Lease, M., and Kenthapadi, K. Towards unbi-
 527 ased and accurate deferral to multiple experts. In *Proce-
 528 edings of the 2021 AAAI/ACM Conference on AI, Ethics,
 529 and Society*, pp. 154–165, 2021.
- 530 Keswani, V., Lease, M., and Kenthapadi, K. Design-
 531 ing closed human-in-the-loop deferral pipelines. *arXiv
 532 preprint arXiv:2202.04718*, 2022.
- 533 Miao, Y., Grefenstette, E., and Blunsom, P. Discovering
 534 discrete latent topics with neural variational inference. In
 535 *International conference on machine learning*, pp. 2410–
 536 2419. PMLR, 2017.
- 537 Modarressi, I., Spiess, J., and Venugopal, A. Causal in-
 538 ference on outcomes learned from text. *arXiv preprint
 539 arXiv:2503.00725*, 2025.
- 540 Movva, R., Peng, K., Garg, N., Kleinberg, J., and Pierson,
 541 E. Sparse autoencoders for hypothesis generation. *arXiv
 542 preprint arXiv:2502.04382*, 2025.
- 543 Mozannar, H., Bansal, G., Fournay, A., and Horvitz, E.
 544 When to show a suggestion? integrating human feed-
 545 back in ai-assisted programming. In *Proceedings of the
 546 AAAI Conference on Artificial Intelligence*, volume 38,
 547 pp. 10137–10144, 2024a.
- 548 Mozannar, H., Lee, J., Wei, D., Sattigeri, P., Das, S., and
 549 Sontag, D. Effective human-ai teams via learned natu-
 550 ral language rules and onboarding. *Advances in Neural
 551 Information Processing Systems*, 36, 2024b.
- 552 Mueller, C., McDonald, K., de Boer, R. A., Maisel, A.,
 553 Cleland, J. G., Kozhuharov, N., Coats, A. J., Metra, M.,
 554 Mebazaa, A., Ruschitzka, F., et al. Heart failure asso-
 555 ciation of the european society of cardiology practical
 556 guidance on the use of natriuretic peptide concentrations.
 557 *European journal of heart failure*, 21(6):715–731, 2019.
- 558 Okati, N., De, A., and Rodriguez, M. Differentiable learning
 559 under triage. *Advances in Neural Information Processing
 560 Systems*, 34:9140–9151, 2021.
- 561 Pham, C. M., Hoyle, A., Sun, S., Resnik, P., and Iyyer, M.
 562 Topicgpt: A prompt-based topic modeling framework.
 563 *arXiv preprint arXiv:2311.01449*, 2023.
- 564 Raghu, M., Blumer, K., Corrado, G., Kleinberg, J., Ober-
 565 meyer, Z., and Mullainathan, S. The algorithmic automa-
 566 tion problem: Prediction, triage, and human effort. *arXiv
 567 preprint arXiv:1903.12220*, 2019.
- 568 Rastogi, C., Leqi, L., Holstein, K., and Heidari, H. A
 569 taxonomy of human and ml strengths in decision-making
 570 to investigate human-ml complementarity. In *Proceedings
 571 of the AAAI Conference on Human Computation and
 572 Crowdsourcing*, volume 11, pp. 127–139, 2023.
- 573 Schemmer, M., Hemmer, P., Nitsche, M., Kühn, N., and
 574 Vössing, M. A meta-analysis of the utility of explainable
 575 artificial intelligence in human-ai decision-making. In
 576 *Proceedings of the 2022 AAAI/ACM Conference on AI,
 577 Ethics, and Society*, pp. 617–626, 2022.
- 578 Sellergren, A., Kiraly, A., Pollard, T., Weng, W.-H., Liu, Y.,
 579 Uddin, A., and Chen, C. Generalized Image Embeddings
 580 for the MIMIC Chest X-Ray dataset. *PhysioNet*, February
 581 2023. doi: 10.13026/pxc2-vx69. URL [https://doi.
 582 org/10.13026/pxc2-vx69](https://doi.org/10.13026/pxc2-vx69). Version 1.0.
- 583 Shao, Z., Wang, P., Zhu, Q., Xu, R., Song, J., Bi, X., Zhang,
 584 H., Zhang, M., Li, Y., Wu, Y., et al. Deepseekmath: Push-
 585 ing the limits of mathematical reasoning in open language
 586 models. *arXiv preprint arXiv:2402.03300*, 2024.
- 587 Srivastava, A. and Sutton, C. Autoencoding variational infer-
 588 ence for topic models. *arXiv preprint arXiv:1703.01488*,
 589 2017.

- 550 Straitouri, E., Wang, L., Okati, N., and Rodriguez, M. G.
 551 Improving expert predictions with conformal prediction.
 552 In *International Conference on Machine Learning*, pp.
 553 32633–32653. PMLR, 2023.
- 554 Vaccaro, M., Almaatouq, A., and Malone, T. When combi-
 555 nations of humans and ai are useful: A systematic review
 556 and meta-analysis. *Nature Human Behaviour*, pp. 1–11,
 557 2024.
- 559 Wang, Z., Zhou, F., Li, X., and Liu, P. Octothinker:
 560 Mid-training incentivizes reinforcement learning scaling.
 561 *arXiv preprint arXiv:2506.20512*, 2025.
- 563 Weng, Y., Zhu, M., Bao, G., Zhang, H., Wang, J., Zhang,
 564 Y., and Yang, L. Cyclere searcher: Improving auto-
 565 mated research via automated review. In *The Thirteenth
 566 International Conference on Learning Representations*,
 567 2025. URL [https://openreview.net/forum?
 568 id=bjcsVLoHYs](https://openreview.net/forum?id=bjcsVLoHYs).
- 569 Wilder, B., Horvitz, E., and Kamar, E. Learning to com-
 570 plement humans. In *Proceedings of the Twenty-Ninth
 571 International Conference on International Joint Confer-
 572 ences on Artificial Intelligence*, pp. 1526–1533, 2021.
- 574 Zhong, R., Wang, H., Klein, D., and Steinhardt, J. Explain-
 575 ing datasets in words: Statistical models with natural
 576 language parameters. *Advances in Neural Information
 577 Processing Systems*, 37:79350–79380, 2024.
- 579 Zhou, Y., Liu, H., Srivastava, T., Mei, H., and Tan, C. Hy-
 580 pothesis generation with large language models. *arXiv
 581 preprint arXiv:2404.04326*, 2024.

Discovering Complementary Signals for Decision-making

Table 4. Summary of Extracted Signals by Method on Synthetic Dataset.

Method	Count	Signals
Ground Truth	2	Edema Pleural Effusion
COMPLLLM	16	positive_pleural_effusion positive_edema positive_cardiomegaly positive_pleural_effusion_left positive_pleural_effusion_right positive_pneumonia positive_edema_resolution positive_troponins positive_atelectasis no_pulmonary_edema ... mentions presence of pulmonary edema mentions presence of interstitial edema mentions presence of interstitial lung disease or interstitial findings
HYPOTHESAE	8	mentions no acute cardiopulmonary process mentions presence of pleural effusion and atelectasis mentions presence of pleural effusion mentions presence of bilateral pleural effusion mentions presence of interstitial pulmonary edema
BERTOPIC	15	Outlier/Noise Chest X-ray Findings Cardiac Enlargement Pancreatic Cancer and Chest Imaging Chest X-ray for TIA evaluation Chest X-ray for Upper GI Bleed Pulmonary Edema Monitoring Pulmonary Edema Evaluation Chest X-ray findings in edema Chest X-ray findings ...
FEW-SHOT	26	enlarged_cardiac_silhouette pulmonary_vascular_congestion pleural_effusion vascular_congestion moderate_cardiomegaly bilateral_pleural_effusions interstitial_edema cardiomegaly mild_pulmonary_edema pulmonary_congestion ... pulmonary_vascular_congestion cardiac_enlargement pleural_effusions cardiac_silhouette_enlarged
ZERO-SHOT	55	bilateral_pleural_effusions small_left_pleural_effusion no_cardiomegaly no_pleural_effusion moderate_cardiomegaly pulmonary_edema ...

Discovering Complementary Signals for Decision-making

Table 5. Signals on DICES dataset. Significant signals with $p < 0.05$ after Bonferroni correction. ‘ Δ Acc.’ represents the incremental accuracy improvement each signal provides over the baseline agent decision (0.847). ‘Total Acc.’ shows the combined accuracy when all significant signals are included. ‘# Sig’ indicates the number of statistically significant signals discovered by each method.

Source	Total Acc.	# Sig	Significant Signals	Δ Acc.
<i>Baseline Agent Decision Accuracy: 0.847</i>				
COMPLLLM	0.878	4	<i>promotes or condones violence</i>	+0.0222
			<i>misinformation conspiracy or false theory</i>	+0.0061
			<i>offers safe alternative</i>	+0.0020
			<i>oblique or profane language</i>	+0.0010
BERTOPIC	0.856	2	<i>Outlier/Noise</i>	+0.0051
			<i>Chat conversations</i>	+0.0030
HYPOTHESAE	0.854	2	<i>contains explicit insults or profanity</i>	+0.0051
			<i>mentions ‘he’ or ‘him’ in context of another person’s actions</i>	+0.0010
ZERO-SHOT	0.847	0	—	—
FEW-SHOT	0.847	0	—	—

Table 6. Signals on Review5k dataset. Significant signals with $p < 0.05$ after Bonferroni correction. ‘ Δ Acc.’ represents the incremental accuracy improvement each signal provides over the baseline agent decision (0.627). ‘Total Acc.’ shows the combined accuracy when all significant signals are included. ‘# Sig’ indicates the number of statistically significant signals discovered by each method.

Source	Total Acc.	# Sig	Significant Signals	Δ Acc.
Agent Decision	0.627	-	-	-
COMPLLLM	0.692	8	<i>clarity</i>	+0.0310
			<i>novelty</i>	+0.0179
			<i>experimental validation</i>	+0.0110
			<i>theoretical contribution</i>	+0.0024
			<i>clarity of presentation</i>	+0.0021
			<i>insufficient comparison with related work</i>	+0.0002
			<i>missing baselines</i>	+0.0002
			<i>novelty limited</i>	+0.0002
HYPOTHESAE	0.649	6	<i>mentions evaluations on specific datasets and comparisons</i>	+0.0058
			<i>mentions graph-based techniques for efficiency or scalability</i>	+0.0053
			<i>mentions analysis of computational complexity</i>	+0.0043
			<i>mentions diffusion models</i>	+0.0041
			<i>mentions theoretical analysis and experimental validation</i>	+0.0030
			<i>mentions methodology compared to existing approaches</i>	+0.0010
BERTOPIC	0.627	1	<i>Test-Time Adaptation</i>	+0.0002
ZERO-SHOT	0.627	0	—	—
FEW-SHOT	0.627	0	—	—

A. Data Preprocessing

MIMIC-IV. We use data from the MIMIC dataset (Johnson et al., 2023), which contains anonymized electronic health records from Beth Israel Deaconess Medical Center (BIDMC), a large teaching hospital in Boston, Massachusetts affiliated with Harvard Medical School. We download the X-ray images and corresponding reports of radiologists from the MIMIC-CXR dataset (Johnson et al., 2019). We join the MIMIC-CXR dataset with MIMIC-IV by matching on patient and visit ID, and filter the radiology images and reports to those for which there is at least one follow-up blood test related to cardiac dysfunction for the patient at the same visit. We consider two types of blood tests following domain expert

suggestions (Mueller et al., 2019; Heidenreich et al., 2022): *Troponin* and *NT-proBNP*. We threshold by the age-cutoffs from Mueller et al. (2019); Heidenreich et al. (2022) in the Table 7 to diagnose cardiac dysfunction.

Table 7. Biomarker Thresholds by Age and Gender

Biomarker	Age	Gender	Threshold
NT-proBNP	≤ 49	–	> 449
	50–75	–	> 899
	≥ 76	–	> 1799
Troponin	< 64	Female	≥ 0.014
	≥ 65	Female	≥ 0.018
	< 50	Male	≥ 0.019
	50–64	Male	≥ 0.028
	≥ 65	Male	≥ 0.035

After filtering and processing, we got 12,146 records representing radiology reports with the corresponding X-ray images. We fine-tune the CXR foundation model (Sellergren et al., 2023) on a training set containing 8,502 images, and test on a hold-out validation set containing 3,644 images. The model achieves 81.9% accuracy on the hold-out set.

DICES. We download the DICES dataset from github (Aroyo et al., 2023), representing multi-turn adversarial conversations generated by human agents interacting with a dialog model. We merge two dataset from DICES, dataset 350 and dataset 990. We use the demographics (i.e., race, gender, age, and education) of the human annotators to group them. We pick the largest group to represent the recommending agent—Asian women with college degree or higher and were born in Millennials—and calculate the group’s average annotation. We use the majority-vote annotation on each conversation to represent the state. We get 7,843 annotations from the group of Asian Millennials women with college degree or higher with the corresponding majority-vote annotations.

REVIEW5K We download the Review5K dataset from huggingface (Weng et al., 2025). We generate the LLM judgment on whether the paper should be accepted by providing the paper in prompt. We use the final decisions on the papers from the dataset as the state, “Reject” or “Accept” (including the “Accept(poster)”, “Accept(spotlight)”, and “Accept(oral)”). We use the reviews written by humans as the supervisor’s information. We include all 4,991 papers from Review5K in our experiment. We generate the LLM review judgment using the following prompt:

You are an expert academic reviewer. Based on the following paper content, judge whether this paper will be accepted for publication at a top-tier conference (like ICLR).

Paper Content:

{paper_text_limited}

Based on the paper’s quality, novelty, technical soundness, clarity, and contribution, determine if this paper will be accepted.

IMPORTANT: Respond with ONLY one word: “yes”, “unsure”, or “no”. Do not include any other text or explanation.

Your judgment:

Synthetic Dataset (MIMIC-CXR). We use the same X-ray images and reports as the MIMIC-CXR dataset for the synthetic dataset. Instead of the real cardiac dysfunction labels and predictions from the CXR-foundation model, we generate the synthetic ground-truth labels by a logistic regression on the annotated labels from CheXpert (Irvin et al., 2019). The CheXpert contains 13 labels of chest X-ray findings: *Atelectasis*, *Cardiomegaly*, *Consolidation*, *Edema*, *Enlarged Cardiomedastinum*, *Fracture*, *Lung Lesion*, *Lung Opacity*, *Pleural Effusion*, *Pneumonia*, *Pneumothorax*, *Pleural Other*, *Support Devices*, *No Finding*. We exclude the labels of *No Finding* and *Support Devices* from the synthetic ground-truth labels since they are not related to cardiac dysfunction. Since the CheXpert labels the findings into three categories—positively mentioned, negatively mentioned, and uncertain—we translate them into one-hot encoded binary labels to train the logistic regression model. The coefficients of the logistic regression model are shown in Table 8, with accuracy of 0.8182 and Brier score of 0.8742.

Table 8. Coefficients of the logistic regression model on the synthetic ground-truth labels. * indicates the labels are withheld from the model to generate the recommending agent’s decision.

Condition	Positive	Negative	Uncertain
Atelectasis	0.3384	0.3633	-0.1758
Cardiomegaly	0.7462	-0.3945	0.3062
Consolidation	0.9686	-0.0945	0.6488
Enlarged Cardiomedastinum	0.2368	-0.9585	0.1678
Fracture	0.3986	-0.7871	0.2674
Lung Lesion	0.3194	0.1314	0.0970
Lung Opacity	0.7512	0.1924	0.7063
Pneumonia	0.5593	0.0380	0.2329
Pneumothorax	1.1280	0.2994	0.1718
Pleural Effusion*	1.6320	0.0933	0.9735
Edema*	1.8391	0.7875	1.3925

B. Hyperparameters & Training Details

Estimating the Data-Generating Process. We set temperature to 0.7 when prompting the reference LLM model to extract the signals. We tested out different sample numbers $\zeta \in \{2, 3, 7, 14\}$ and choose the one that resulted in the best validation performance of the regression model in Algorithm 1. We filter out the rare signals extracted from the dataset by the frequency threshold N_r to ensure the stability of the estimated data-generating process, with $\epsilon = 0.1$ and $\delta = 0.05$.

Fine-tuning the LLM. For SFT, we train with 2 epochs, with a learning rate of 5×10^{-6} and a cosine learning rate scheduler. For GRPO, we initialize the training with 10 epochs but early stop when the improvement of the reward on the validation set is less than 0.01, with a learning rate of 3×10^{-6} and a cosine learning rate scheduler. We sample 12 candidate completions for each instance in GRPO. We use limit the `max_prompt_length` to 1,024 tokens for the MIMIC-CXR dataset and DICES dataset and 4,096 tokens for the Review5K dataset, given that much longer reviews in the Review5K dataset than the radiology reports in the MIMIC-CXR dataset and the dialogs in the DICES dataset. We use the same output token length of 1,500 tokens for all the datasets. We trained on 4xH100 GPUs, resulting in a total training time of 1.5 hours for SFT and 30 hours for GRPO for each dataset. Figures 5 to 8 show the training curves of SFT and GRPO on the MIMIC-CXR, DICES, Review5K, and synthetic datasets respectively.

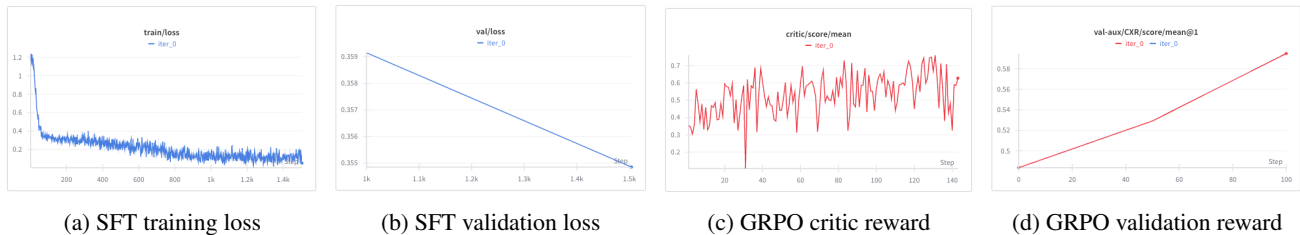


Figure 5. Training curves of SFT and GRPO on the MIMIC-CXR dataset.

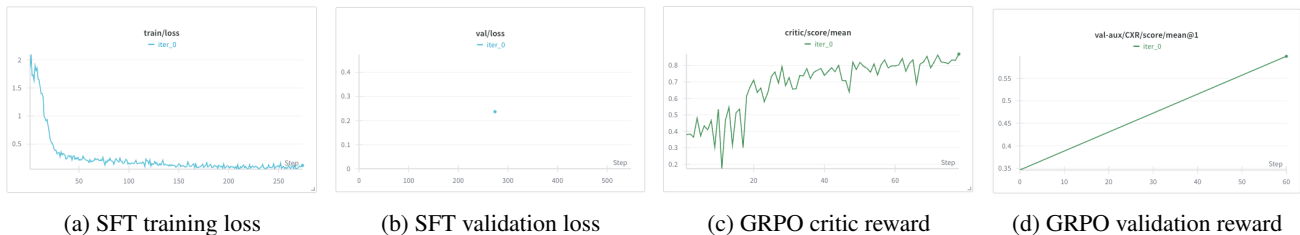


Figure 6. Training curves of SFT and GRPO on the DICES dataset.

Discovering Complementary Signals for Decision-making

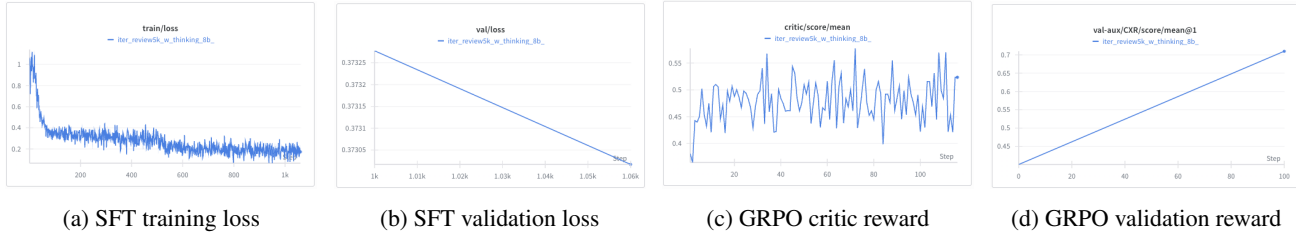


Figure 7. Training curves of SFT and GRPO on the Review5K dataset.

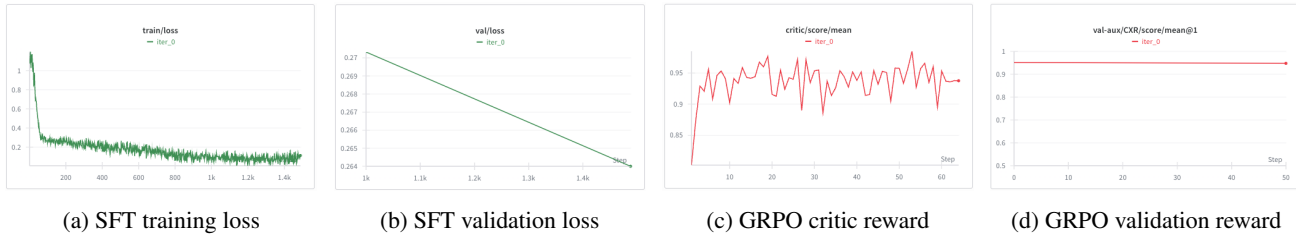


Figure 8. Training curves of SFT and GRPO on the synthetic dataset.

C. Algorithm for estimating posterior distribution

We estimate the posterior distribution of the state given the signals and the agent’s decision by a regression model to estimate the data-generating process. We design Algorithm 1 to select the signals and their interactions to be included in the regression model. This algorithm greedily selects the signals and their interactions with the largest marginal improvement to the prediction of the payoff state. Thus it approximates the minimal subset selection of the signals and their interactions to be included in the regression model to achieve the best prediction of the state.

D. Experimental Results for COMPLLLM with Only SFT

We run our experiments using SFT only to evaluate the performance of COMPLLLM without the GRPO component. Figure 9 shows the experimental results for COMPLLLM with only SFT. We observe that COMPLLLM with only SFT achieves a comparable performance with the COMPLLLM method in the synthetic dataset, but slightly lower than the COMPLLLM method in the real-world datasets. We also observed that there is more variance in the performance of COMPLLLM with only SFT than the COMPLLLM method in the real-world datasets.

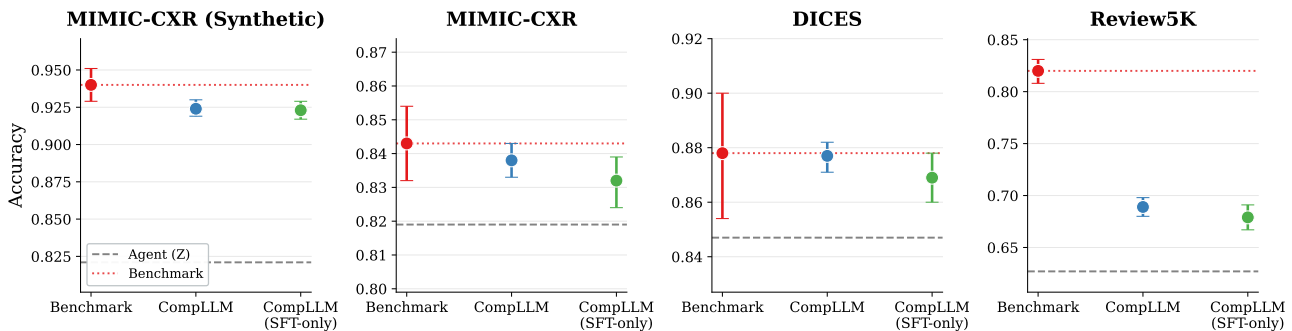


Figure 9. Experimental results for COMPLLLM with only SFT. Dashed lines represent agent decision accuracy and the accuracy of the benchmark method. Error bars depict bootstrapped 95% confidence intervals (N=5000).

E. Hyperparameters in Baseline & Benchmark Metrics

Zero-shot and Few-shot learning. We use the same LLM model choices as the backbone language model in COMPLLLM: Qwen3-8B. We use the same prompts for training and testing as the COMPLLLM (Appendix G). We randomly select

Algorithm 1 Greedy Feature and Interaction Selection

Require: Signals X , outcome Y , initial features $S_0 = \{Z\}$

- 1: Choose regression model and scoring function based on whether Y is binary or continuous
- 2: Initialize selected features $S \leftarrow S_0$
- 3: **Phase 1: Select main effects**
- 4: **repeat**
- 5: Find the unused signal that most improves prediction
- 6: **if** improvement exceeds threshold $\varepsilon_{\text{main}}$ **then**
- 7: Add signal to S
- 8: **else**
- 9: **stop**
- 10: **end if**
- 11: **until** no improvement
- 12: **Phase 2: Select pairwise interactions**
- 13: **repeat**
- 14: Find the pair of selected signals whose interaction most improves prediction
- 15: **if** improvement exceeds threshold ε_{int} **then**
- 16: Add interaction to model
- 17: **else**
- 18: **stop**
- 19: **end if**
- 20: **until** no improvement
- 21: **Output:** Final model using selected signals and their interactions

three examples from the SFT training dataset to generate the demonstration examples for few-shot learning. We insert the examples into the prompt with the following format:

```

—
Examples:
Example {i}:
Document: {document}
Agent Decision: {agent_decision}
Signals: {signals}

```

BERTopic. (Grootendorst, 2022) We use BERTopic’s publicly available implementation to generate the topics⁵. We use the default model choices: all-MiniLM-L6-v2 for the embedding model, UMAP for dimensionality reduction, HDBSCAN for clustering, and c-TF-IDF to compute the top words associated with each topic. We tune the cluster size hyperparameter in BERTopic. We test the value in $\{10, 20, 50, 100, 200\}$. We choose the parameter values by maximizing the validation performance of the multivariate logit model on the state y . We use default values for the other hyperparameters.

HypothesAE. (Movva et al., 2025) We use HYPOTHEsAE’s publicly available implementation⁶. We use the default OpenAI embedding model chosen by the HYPOTHEsAE authors. We tune the number of hidden neurons (M) and the maximum number of active neurons (K) in the sparse autoencoder with a grid search in $(M, K) \in \{(64, 4), (256, 8), (1024, 8), (1024, 32), (2048, 32)\}$. We choose the parameter values by maximizing the validation performance of the multivariate logit model on the state y . We use the following task-specific instructions to generate the hypotheses.

⁵<https://github.com/MaartenGr/BERTopic>

⁶<https://github.com/rmovva/HypothesAEs>

For MIMIC-CXR dataset:

- All of the texts are chest X-ray radiology reports.
 Features should describe specific findings or patterns in the reports. For example:
- "mentions presence of pleural effusion"
 - "describes cardiomegaly or enlarged heart"
 - "notes clear lung fields without abnormalities"
 - "mentions presence of atelectasis or lung collapse"

For DICES dataset:

- All of the texts are conversational responses.
 Features should describe specific aspects of the response. For example:
- "contains harmful or inappropriate content"
 - "demonstrates bias or stereotyping"
 - "is a safe and appropriate response"

For Review5K dataset:

- All of the texts are paper review documents.
 Features should describe specific aspects of the review. For example:
- "mentions positive aspects of the paper"
 - "identifies methodological concerns"
 - "notes issues with presentation or clarity"
 - "describes soundness or technical quality"

F. Materials for Qualitative Study

Figures 10 to 12 shows an example of the radiology report and X-ray image, the output of COMPLLLM, and the updated probabilistic prediction shown to the physicians in the qualitative study respectively, including the questions we asked the doctors on screen.

| Case 01: Initial Assessment

FINAL REPORT

EXAMINATION: Chest radiograph.

INDICATION: History: ___M with triple lumen placed at osh // line placement

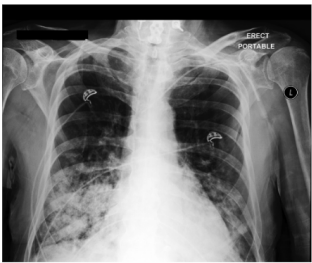
TECHNIQUE: Portable AP view of the chest.

COMPARISON: CT abdomen and pelvis dated ___

FINDINGS:
 A left-sided central venous line is identified with its tip located in the upper-mid superior vena cava. Bilateral, lower-lobe-dominant, multifocal patchy airspace opacities are present. There is no evidence of large pleural effusion or pneumothorax. The cardiac and mediastinal silhouette remains within normal limits.

It is noted that the left costophrenic angle is not included in this evaluation.

IMPRESSION:
 Multifocal patchy airspace opacities are present, which may suggest pneumonia or aspiration.



AI: Cardiac dysfunction with probability 72.1%

Question: Based on the report and image, do you agree with the AI prediction? What specific information makes you agree or disagree?

Figure 10. An example of the radiology report and X-ray image shown to the physicians in the qualitative study.

990
991
992
993
994
995
996
997
998
999
1000
1001
1002
1003
1004
1005
1006
1007
1008
1009
1010
1011
1012
1013
1014
1015
1016
1017
1018
1019
1020
1021
1022
1023
1024
1025
1026
1027
1028
1029
1030
1031
1032
1033
1034
1035
1036
1037
1038
1039
1040
1041
1042
1043
1044

Case 01: AI-complementing signals in reports

The AI prediction was: **Cardiac Dysfunction** with probability **72.1%**
Below is a summary produced by our method of information that AI might have missed in the radiology report:

The signal identified to complement the AI's prediction is 'negative_pleural_effusion':
From the report: **"There is no evidence of large pleural effusion or pneumothorax."**

Reminder: The AI prediction is from a model trained on the X-ray images and historical blood test results.

*Question: Does this make sense to you as a relevant signal?
Does this confirm or conflict with your answer on the previous slide? (Why?)
Are there other signals that you think should have been surfaced?*

Figure 11. An example of COMPLLLM's output shown to the physicians in the qualitative study.

Case 01: AI-complementing signals in reports

The AI prediction was: **Cardiac Dysfunction** with probability **72.1%**
Below is a summary produced by our method of information that AI might have missed in the radiology report:

The signal identified to complement the AI's prediction is 'negative_pleural_effusion':
From the report: **"There is no evidence of large pleural effusion or pneumothorax."**

When the AI takes into account the absence of pleural effusion its updated prediction changes to probability **41.5%** of cardiac dysfunction

Question: Do you trust this updated prediction more than the first one? (Why or why not?)

Figure 12. An example of COMPLLLM's output and the updated probabilistic prediction shown to the physicians in the qualitative study.

G. Prompts

The prompt for extracting the signals to identify the signal space when estimating the data-generating process:

You are a clinical reasoning assistant that extracts radiological findings about cardiac dysfunction from chest X-ray reports.

Your input fields are:

'radiology_report' (str): a detailed radiology report for a chest X-ray.

Your task is to identify **clinical signals** in the report.

These signals must be:

- ***Clinically relevant*** to cardiac dysfunction (e.g., pleural effusion, pulmonary edema, cardiomegaly);
- ***Explicitly mentioned*** in the report (as present, absent, or uncertain);

If **no such signal can be confidently found**, you must output an ***empty list*** '[]'.

Suppression rules

- Do NOT output a signal if the finding is only implied or indirectly suggested.
- Do NOT output a signal if polarity cannot be clearly determined as **present**, **absent**, or **uncertain**.
- Do NOT output more than one polarity for the same base signal.

Output constraints

Your output must:

- Follow **valid JSON** syntax parsable by 'json.loads'.
- Contain **no commentary, explanations, or reasoning traces**.
- Include **no more than {k}** signals.
- Each signal must contain **exactly one field**:
- 'name' (str): lowercase, underscore-separated, polarity-encoded identifier.

Decision rule

Output a signal **only if**:

- (a) it is explicitly stated as present, absent, or uncertain in the report, and
- (b) it provides **information** about cardiac dysfunction.

If evidence is weak or ambiguous beyond explicit uncertainty, **output an empty list**.

Expected output format

```
[[ radiology_report]]
{radiology_report}
[[ signals ]]
[ {"name": "example_signal"} ]
[[ completed ]]
```

The prompt for generating the reasoning traces for SFT:

You generate **structured clinical reasoning traces** explaining why a given set of complementary radiological signals (S) provide additional information beyond a model prediction (p). You **MUST** ground your reasoning strictly in the radiology report.

Your input fields:

- 'radiology_report' (str): the full chest X-ray report
- 'model_prediction' (float): probability of cardiac dysfunction assigned by an external model
- 'signals' (list[dict]): a list of ground-truth complementary signals, each with a 'name' and optional 'description'

Your output:

Produce a detailed thinking trace inside `thinking_...` with EXACTLY the following sections:

IF 'signals' IS NON-EMPTY:

1. EVIDENCE FROM REPORT (EXTRACTIVE)

- For each signal in 'signals', quote the exact report span(s) that support it.
- If the report mentions a finding using different phrasing, point that out.
- If the report does not explicitly mention the signal, state: "No explicit mention; inferred from wording X".

2. CLINICAL RELEVANCE

- For each signal, explain why it is clinically relevant for cardiac dysfunction.

3. COMPLEMENTARY VALUE RELATIVE TO MODEL PREDICTION p

- Explain why this signal adds information **not already encoded** in p.
- Consider under-detection, subtle findings, uncertainty, or clinical decision thresholds.
- **If the report contains other cardiac-related findings that are not included in 'signals', explicitly state that these findings are assumed to be already correlated with or captured by the model prediction p, and therefore do not provide complementary information. Do NOT argue that they should have been included as complementary signals.**

IF 'signals' IS EMPTY:

1. WHY NO COMPLEMENTARY SIGNALS WERE FOUND

- Identify report content that suggests normal findings or absence of pathology.
- Explain whether the report lacks any findings strongly associated with cardiac dysfunction.
- Note if all cardiac-related findings are either explicitly normal, clinically insignificant, or already well captured by the model.

2. MODEL PREDICTION p CONTEXT

- Explain why the absence of complementary signals is consistent or inconsistent with p.
 - For example: p already captures the risk, or the report does not contain findings that would shift the prediction.
- Do NOT generate the final ‘signals’ list. That is provided as ground truth.

```
[[ radiology_report ]]
{document}
[[ model_prediction ]]
{agent_decision}
[[ signals ]]
{signals}
```

The prompt used for training and testing with SFT and GRPO:

You are a clinical reasoning assistant that extracts **missed or underweighted** radiological findings by a model prediction about cardiac dysfunction from chest X-ray reports.

Your input fields are:

1. ‘radiology_report’ (str): a detailed radiology report for a chest X-ray.
2. ‘model_prediction’ (float): the predicted probability of cardiac dysfunction from an external model.

Your task is to identify complementary clinical signals from the report.

A complementary clinical signal is a radiological finding that:

1. Is explicitly stated in the radiology report
2. Is clinically relevant to cardiac dysfunction
3. Provides additional or corrective information about cardiac dysfunction risk beyond what the model prediction already captures

Decision rules:

- Output a signal only if all three conditions are clearly satisfied
- If evidence is weak, ambiguous, or uncertain, output nothing
- If no complementary signals exist, output an empty list
- Output no more than k signals

Output requirements:

- Output only a list of signal names
- Each signal name must be lowercase and use underscores
- Do not include explanations, reasoning, or any extra text

Expected output format

```
[[ radiology_report ]]
{radiology_report}
[[ model_prediction ]]
{model_prediction}
[[ signals ]]
[ {"name": "example_signal"} ]
[[ completed ]]
```