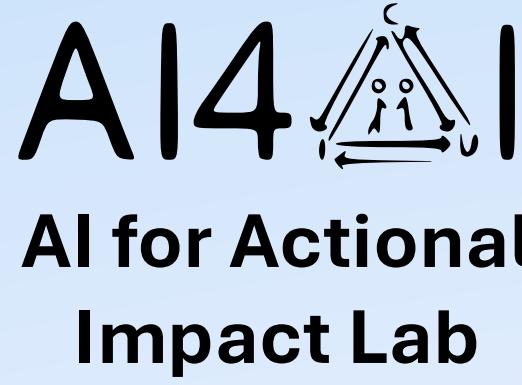


Improving ARDS Diagnosis Through Context-Aware Concept Bottleneck Models

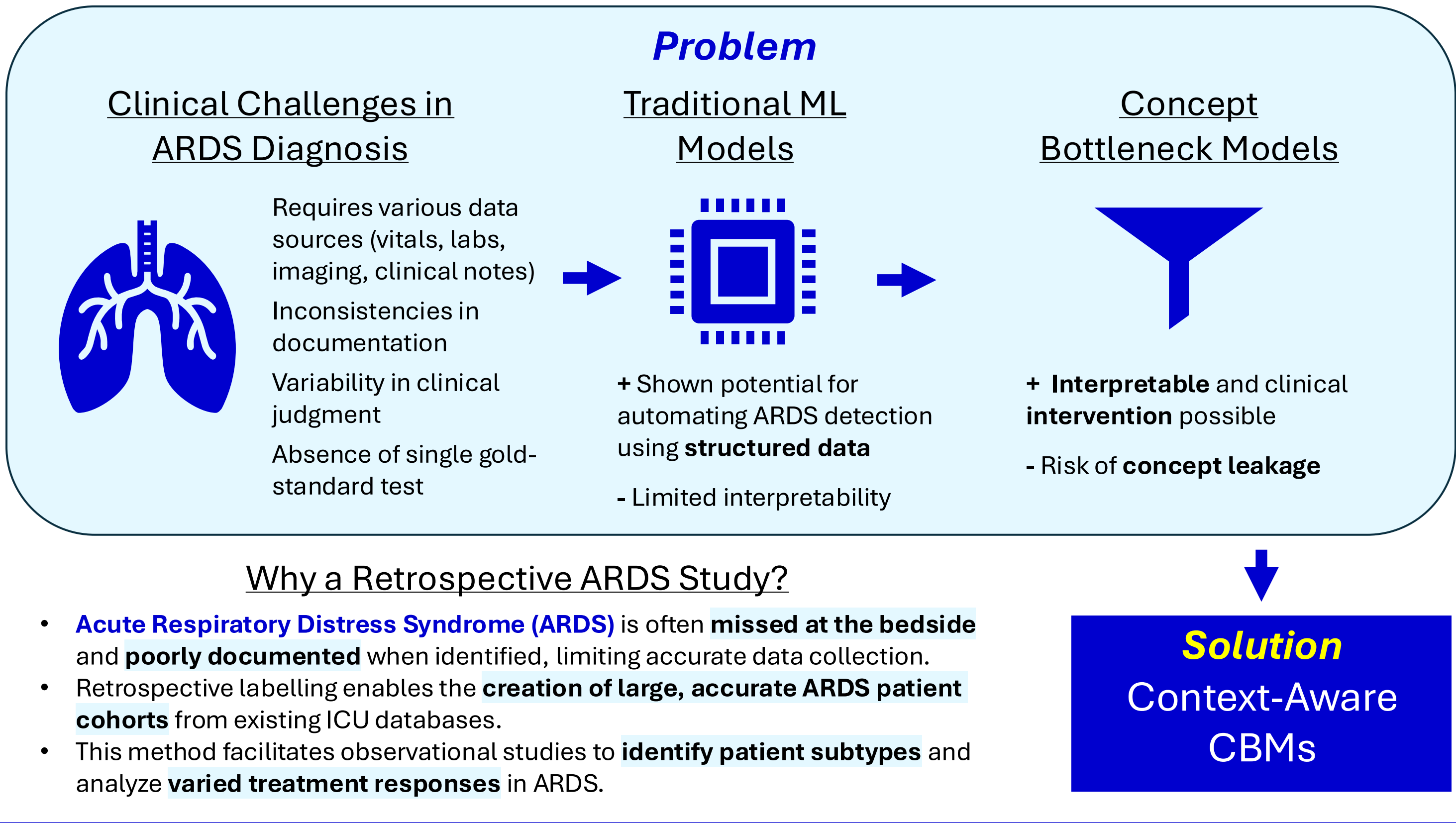
Authors: Anish Narain, Ritam Majumdar, Nikita Narayanan, Dominic C Marshall, Sonali Parbhoo
Contact: s.parbhoo@imperial.ac.uk

IMPERIAL

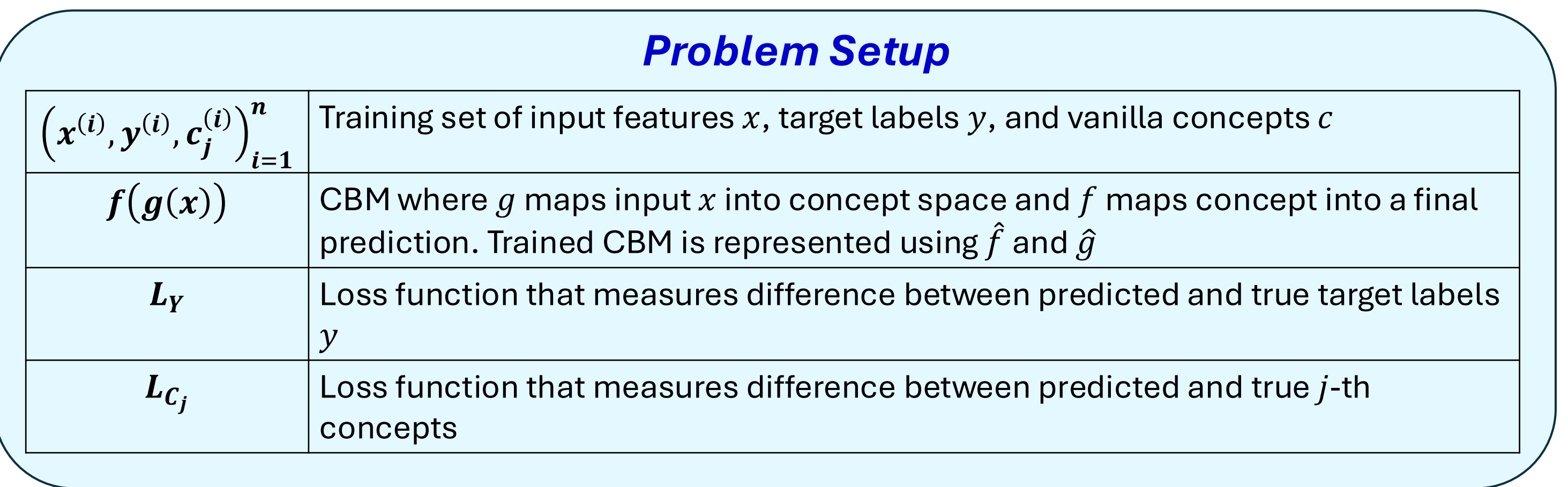


Project
GitHub

1. Overview



2. Context-Aware CBM



Context-Aware CBM Loss Function

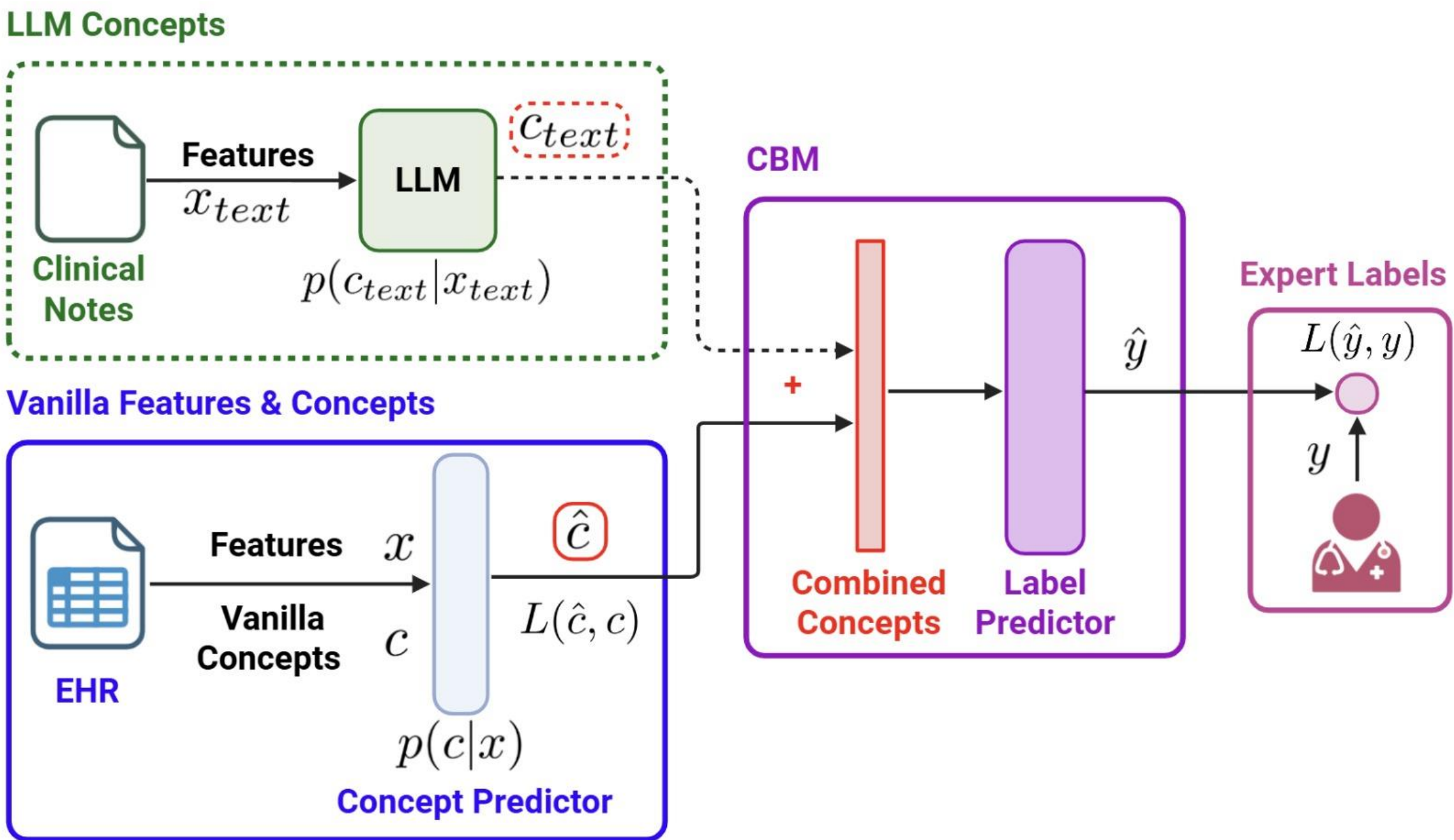
In Context-Aware CBMs, the initial predicted concepts come from distribution $p(c|x)$ and the LLM concepts come from a different distribution $p(c_{text}, x_{text})$. The loss function thus becomes:

$$\hat{g}, \hat{f} = \arg \min_{f, g} \sum_i \left[L_Y(f(g(x^{(i)}), c_{text}^{(i)}); y^{(i)}) + \sum_j \lambda L_{C_j}(g(x^{(i)}); c^{(i)}) \right]$$

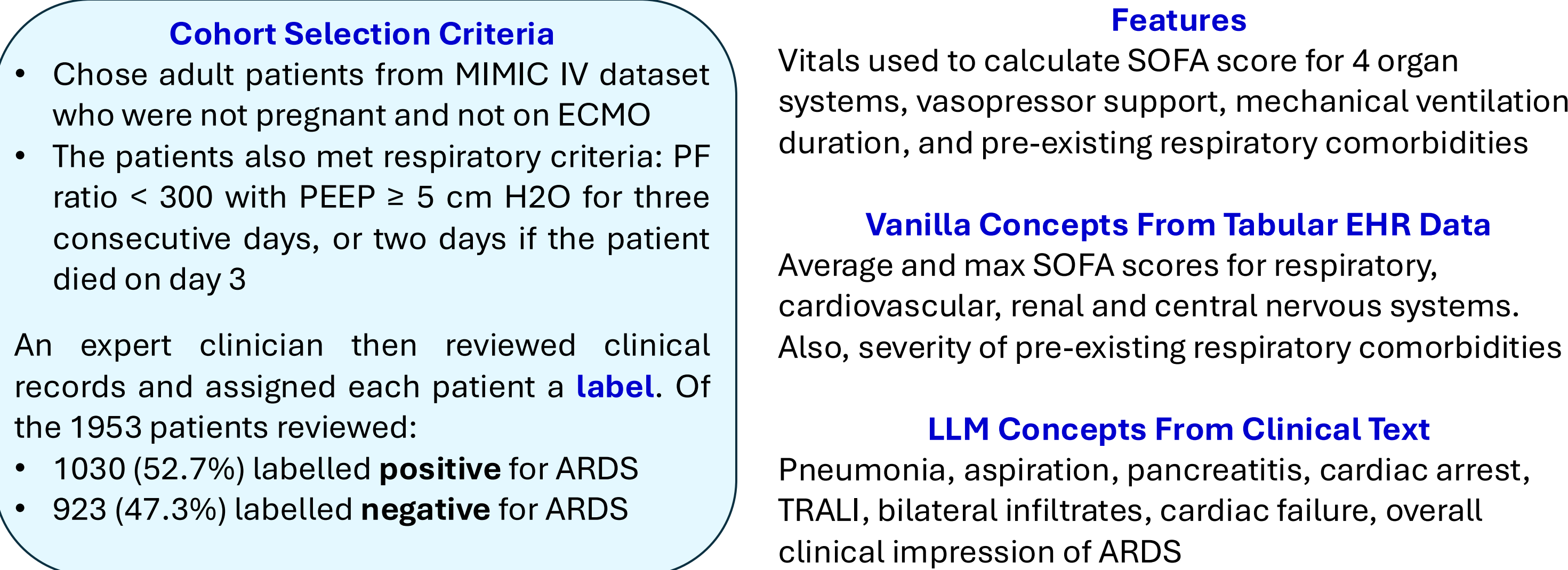
LLM Prompt Template

```
prompt_template=(
    "Context: You are a clinician receiving chunks of clinical text for patients in an ICU. Please do the reviewing as quickly as possible."
    "Task: Determine if the patient suffered from pancreatitis."
    "Instructions: Answer with 'Yes' or 'No'. If there is not enough information, answer 'No'."
    "Discharge Text: {discharge_text}"
    "Query: Does the chunk of text mention that the patient suffered from pancreatitis? Answer strictly in 'Yes' or 'No'."
),
input_variables=["discharge_text"]
```

Training Pipeline for Context-Aware CBMs

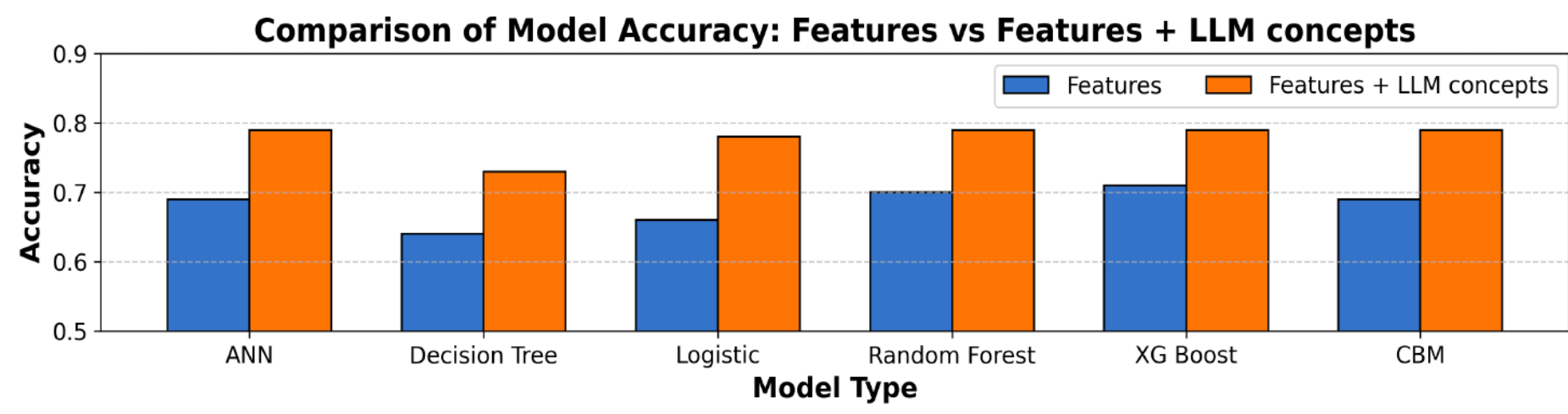


3. Cohort Selection and Dataset



4. Results

1. Context-Aware models augmented with LLM concepts show an **8-10% improvement over traditional models** across all metrics



Result #1: Comparative model performance between Features and Features augmented with LLM concepts.

2. Context-Aware CBMs **outperform vanilla CBMs by 10-20%** across all metrics, indicating that LLMs extract important concepts relevant for final ARDS prediction

3. Context-Aware CBMs product concepts which have **higher mutual information**, indicating **leakage reduction**

4. **Interventions** on misclassified concepts further **boost the performance** of Context-Aware CBMs by **12-20%**

Model Type	Acc.	Prec.	Recall	F1	M1
Vanilla CBM	0.69	0.71	0.72	0.71	0.08
Context-Aware CBM	0.79	0.76	0.88	0.81	0.21
Vanilla CBM + Vanilla concept interventions	0.85	0.88	0.86	0.90	0.78
Context-Aware CBM + Vanilla concept interventions	0.87	0.88	0.91	0.92	0.84
Context-Aware CBM + LLM concept interventions	0.94	0.92	1.00	0.93	0.92
Context-Aware CBM + (LLM + Vanilla) concept int.	0.96	1.00	0.93	0.97	0.96

Result #2,3,4: Performance of Vanilla and Context-Aware CBMs augmented with interventions over misclassified patients

5. Context-Aware CBMs **improve prediction** of concepts **critical to the outcome**, thereby **reducing leakage**

6. **Interventions** based on **correlated concepts** are **more effective** than interventions on individual concepts

7. Context-aware CBMs **outperform** vanilla CBMs on **out-of-distribution patients**

5. Case Studies

Context-aware CBMs are **interpretable** and **intervenable**. These properties allow us to identify key diagnostic factors in patient cases and explore how intervening on the concepts affect the diagnosis.

Study 1: Intervention to Correct False Negatives

- **Problem:** Vanilla CBM (using structured EHR data only) missed ARDS in two patients
- **Fix:** The **Context-Aware CBM**, with LLM access to clinical notes, identifies **key textual evidence** (e.g., “bilateral infiltrates,” “pneumonia,” absence of “cardiac failure”) that points to ARDS.
- **Result:** The diagnosis was corrected to ARDS-positive.

Study 2: Intervention to Correct LLM-Induced Errors

- **Problem:** Context-Aware CBM wrongly **over-relies on LLM-inferred concepts** that are not actually present or relevant.
 - **Patient 1:** LLM erroneously picks up “pneumonia/ARDS” → false positive.
 - **Patient 2:** LLM invents a “cardiac arrest” concept → false negative.
- **Fix:** Manual **intervention at the concept level** (adjusting or removing erroneous LLM concepts).
- **Result:** Restored correct predictions made by the vanilla CBM.

Study 3: Concept-Level Debugging

- **Problem:** Both vanilla and Context-Aware CBMs misclassify two patients—one ARDS-positive and one ARDS-negative—due to **inaccurate intermediate concept predictions** (e.g., wrong morbidity scores, missing comorbidities).
- **Fix:** The **interpretability of CBMs** was used to identify which **specific intermediate concepts** (e.g., morbidity, cardiac failure) are incorrect.
- **Result:** Intervening on those concepts **corrected the predictions**.

6. Conclusion

Paper Contributions

- ✓ 1. Proposed a general framework for **enhancing CBMs** using context from **unstructured data**, applicable to clinical use cases requiring multi-modal reasoning.
- ✓ 2. Achieved an **8-10% improvement** in **retrospective ARDS diagnosis**.
- ✓ 3. Enriched concepts with LLM-derived context. This **increased concept completeness**, thereby **mitigating concept leakage**.
- ✓ 4. Enabled transparent concept-level reasoning, allowing for **interventions to correct errors**, thus improving reliability for clinical use.

Paper Limitations

- **Delayed Diagnosis:** The current study relied on retrospective access to complete patient data. In practice, ARDS diagnosis requires rapid integration of diverse data soon after onset, so future models must use data available within a specific, real-time window.
- **Noisy LLM Concepts:** Concepts generated by large language models may introduce irrelevant or misleading signals into predictions. Incorporating human oversight could help reduce errors before clinical application.