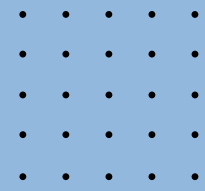


MedCoach

Interim Presentation

Mai Werthaim & Maya Kimhi

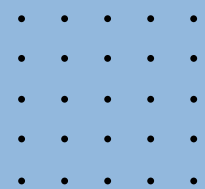


Project Description

Medical students are often overwhelmed by theoretical study and lack practical diagnostic experience.

MedCoach offers an AI-driven learning environment with realistic, interactive case simulations that advance medical decision-making with immediate, relevant feedback

Dataset: Diseases and their Symptoms 
Labels: Prognosis



Case Preparation

Input: Diseases and their symptoms

Output: Patient cases and their diagnosis

Task: Text generation

Student Examination

Input: Partial patient case

Output: Student diagnosis and questions

Task: Question-Answering

Feedback Generation

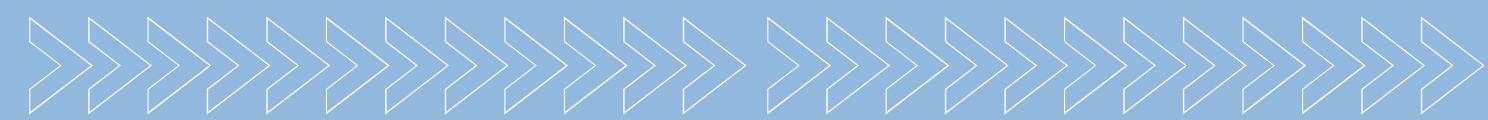
Input: Doctor questions and diagnosis, Student questions and diagnosis,

Output: Student evaluation (avg questions to diagnosis & similarity between doctor & student)

Task: Text similarity

Prior Art

Name	Med-PaLM 2	AMIE	ClinicalGPT-R1
Source	Singhal, K., et al. (2025). Toward expert-level medical question answering with large language models. Nature.. 	Tu, T., et al. (2025). Towards conversational diagnostic artificial intelligence. Nature. 	Lan, W., et al. (2025). ClinicalGPT-R1: Pushing reasoning capability of generalist disease diagnosis with large language model. arXiv. 
Goal	Enhance reasoning and grounding in long-form medical question answering through ensemble refinement and chain-of-retrieval strategies	Conduct AI-driven diagnostic dialogue by simulating clinician–patient interactions	Improve generalist disease diagnosis
Approach	Transformer+ fine-tuning on medical data; uses prompt tuning & ensemble refinement for reliable answers	Vignette generator Dialogue simulator Self-play loops	Synthetic Data Generation Two-Stage Fine-Tuning
Data	USMLE-style questions (MedQA), medical research (PubMedQA), MedMCQA, and clinical topics in MMLU	Real-world transcripts (~99 K conversations from MIMIC-III) and a self-play multi-agent to synthesize new case	Real EHR records with long-chain CoT prompts
Metrics	Accuracy	Clinicians scored AMIE’s history-taking and diagnostic reasoning using PACES-style criteria	Accuracy
Results	86.5 % accuracy on MedQA (+19 % over Med-PaLM)	Generated ~12K dialogues AMIE matched or exceeded benchmarks on key axes	Outperforms GPT-4o in Chinese diagnosis tasks and matches GPT-4o in English on MedBench-Hard



NLP Pipeline

Case Preparation	Benchmark Generation	Student Stimulation	Evaluation
<p>Input: Raw Data: 100 random sample of Symptoms and diagnoses.</p> <p>Output: Table of 4 columns - Diagnosis, full patient case, 80% case, 50% case.</p> <p>Task: Patient Case Creation</p> <p>Model: MedLlama2</p> <p>Metric: Model-based evaluation (PubMedBERT)</p>	<p>Input: A table of 4 columns - Diagnosis, full patient case, 80% case, 50% case.</p> <p>Output: K pairs of columns (doctor's question, doctor's diagnosis)</p> <p>Task: Doctor & patient role playing.</p> <p>Model: Me-LLaMA 13B as doctor & MedLlama2 as patient</p> <p>Metric: Accuracy, AUC</p>	<p>Input: Full patient case, 80% patient case, 50% patient case.</p> <p>Output: K pairs of columns (student question and diagnosis)</p> <p>Task: Student & patient role playing.</p> <p>Model: DeepSeek-R1 as student & MedLlama2 as patient</p> <p>Metric: Accuracy, AUC</p>	<p>Input: Doctors question & diagnosis, student question & diagnosis</p> <p>Output: Similarity between doctor questions & student questions</p> <p>Task: Comparing student to doctor</p> <p>Model: None (NLP Metrics)</p> <p>Metric: average questions to diagnosis, questions cosine similarity.</p>

Data exploration

Raw data - Diseases and their Symptoms

- 2564 rows
- 400 symptoms
- 133 unique diseases
- 13 duplicate rows

[illegible]

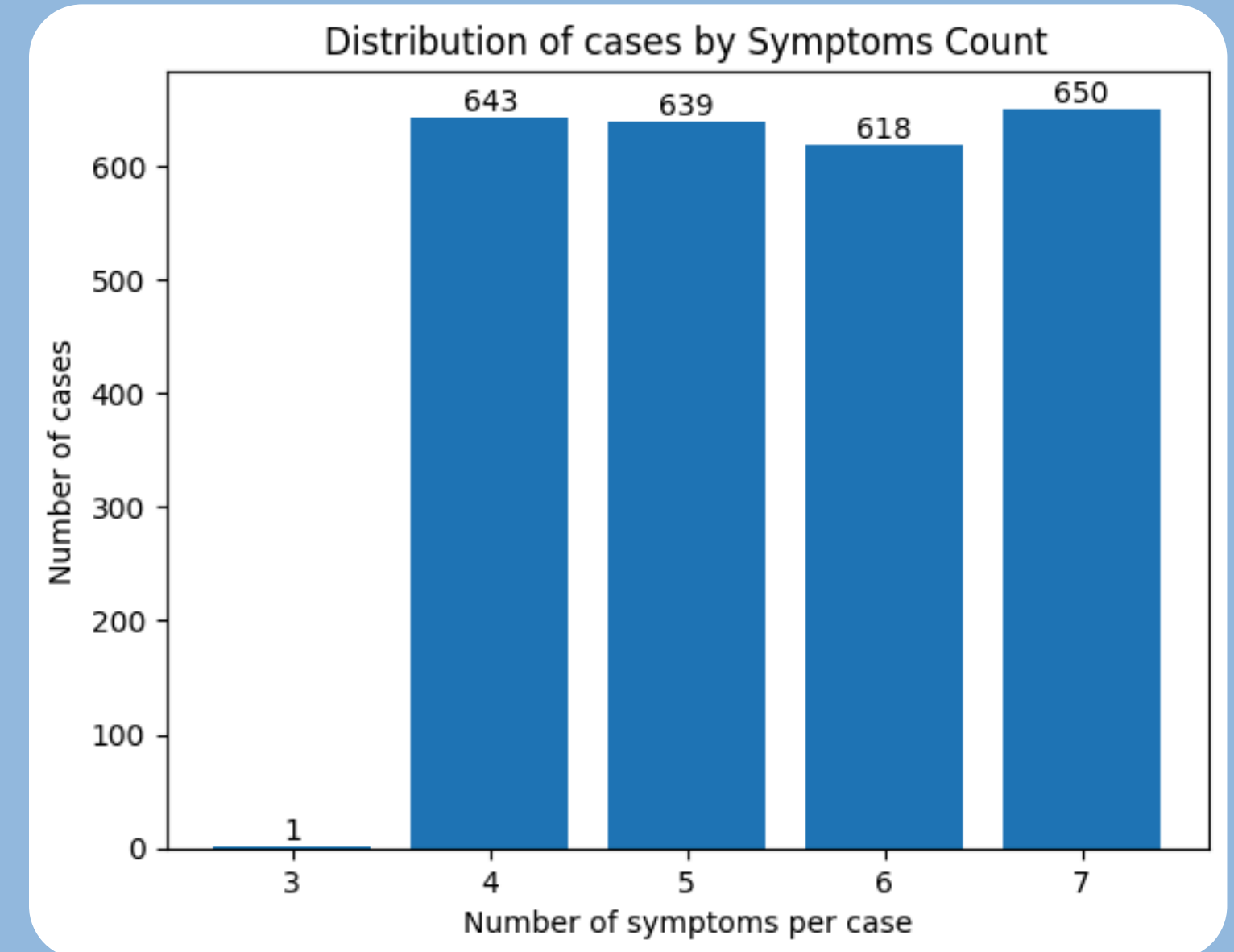
Data exploration

Raw data statistics:

- Average rows per disease: 19.18
- Disease with most rows: bipolar disorder (43 rows)
- Disease with fewest rows: decubitus ulcer (3 rows)
- Each disease have 3-7 symptoms
- The most common symptom is pain (323 cases)
- The least common symptom is dizzy spells (1 case)

Data Treating:

- Duplicate deletion
- Removal of symptoms not associated with any disease
- Selection of cases with ≥ 4 symptoms



Baseline

Random Sampling:

A random sample of 100 rows is selected from the original dataset.
Each row represents a real disease profile with associated symptoms.

Patient Case Generation:

For each selected disease instance, a synthetic patient case is generated using a language model (MedLlama2).
Each case includes:

- Full Case: All symptoms associated with the disease.
- 80% Case: Approximately 80% of the symptoms.
- 50% Case: Approximately 50% of the symptoms.

Text-Based Diagnosis Modeling as doctor:

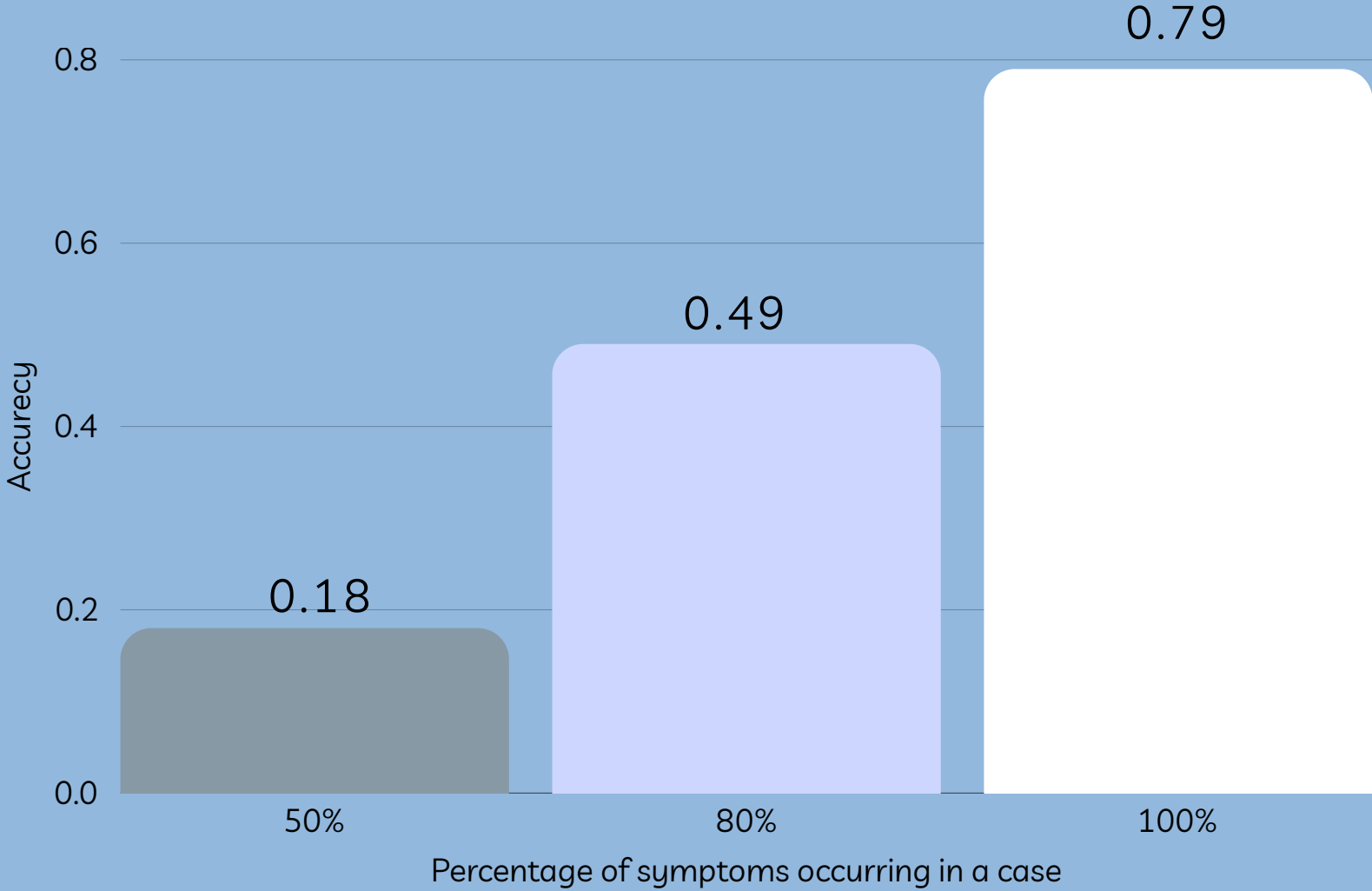
PubMedBERT is fine-tuned on the generated case descriptions to simulate a doctor's diagnosis.

Accuracy Comparison Across Case Levels:

Accuracy is measured for each level to assess how case completeness affects diagnosis quality.

Baseline

	Example case (Asthma)
100%	A 20-year-old female presents with shortness of breath, productive cough, distress respiratory, symptom aggravating factors.
80%	A 20-year-old female presents with symptom aggravating factors, distress respiratory, productive cough.
50%	A 20-year-old female presents with distress respiratory, shortness of breath.



As expected, the accuracy of the diagnosis increases as the percentage of available data in the case rises.

Insights

The data source is rich enough to provide good patient cases for diagnosis

There is a relationship between the amount of exposure and accuracy.

Recommendations

Assessment whether dataset size can be reduced.

Zero-shot diagnosis for further evaluate the robustness of the generated cases.

